

Solving Nim by the Use of Machine Learning

*Exploring how well Nim can be played
by a computer*

Mikael Nielsen Røykenes



Thesis submitted for the degree of
Master in Informatics: Programming and
Networks
60 credits

Department of Informatics
Faculty of mathematics and natural sciences

UNIVERSITY OF OSLO

Autumn 2019

Solving Nim by the Use of Machine Learning

*Exploring how well Nim can be played
by a computer*

Mikael Nielsen Røykenes



© 2019 Mikael Nielsen Røykenes

Solving Nim by the Use of Machine Learning

<http://www.duo.uio.no/>

Printed: Reprocentralen, University of Oslo

Contents

1	Intro	3
1.1	My Goal with This Thesis	3
1.2	Hypothesis	3
2	Ethics	3
3	The Game Nim	4
3.1	Rules and How to Play It	4
3.2	Nimrod	5
3.3	Impartial Games	5
3.4	The Sprague-Grundy Theorem	6
4	Machine Learning	6
4.1	Reinforcement learning	7
4.1.1	The Principle	7
4.1.2	The Proper	7
4.1.3	Action Selection/Policy	8
4.1.4	Sarsa and Q-learning	8
4.2	Supervised Learning	9
4.2.1	The Perceptron	10
4.2.2	The Multilayer Perceptron	11
5	The Algorithms for the Problem	14
5.1	Preexisting Algorithm	14
5.1.1	Nim-sum	14
5.1.2	Nim-sum of the State	15
5.1.3	Finding the Correct Move	16
5.2	Reinforcement Learning	16
5.2.1	Representation of Nim as a graph	16
5.2.2	How the Graph Is Represented in the Program	18
5.2.3	How the Graph Is Made Smaller	18
5.2.4	Training	20
5.2.5	Terminating the Program	21
5.2.6	One-dimensional Program	21
5.3	Supervised Learning	21
5.3.1	Acquiring Data	21
5.3.2	Training	21
6	Describing the Code	23
6.1	The Deterministic Algorithm	23
6.2	Programs for Reinforcement Learning.	24
6.2.1	Initialization	24
6.2.2	Setup	25
6.2.3	Training	29

6.2.4	Making Moves	32
6.2.5	<i>Run.py</i>	35
6.2.6	<i>Play.py</i>	38
6.3	Programs for Supervised Learning	40
6.3.1	<i>MakeData.py</i>	40
6.3.2	The Multilayer Perceptron	42
6.3.3	<i>RunMlp.py</i>	44
6.3.4	Program With Stochastic Termination	46
7	Comparing the Algorithms with Time Complexity	49
7.1	The Deterministic Algorithm	49
7.1.1	Time-complexity of the Nim-sum Operation	49
7.1.2	The Rest of the Algorithm	50
7.1.3	Practically, Then?	50
7.2	The Reinforcement Learning Algorithm	51
7.2.1	Setup of the Matrices	51
7.2.2	Training the Matrix	52
7.2.3	Giving an Answer	52
8	Comparing the Algorithms with Actual Time-use	52
8.1	Reinforcement Learning	53
8.1.1	The naming of the programs in the graph	53
8.1.2	The Settings Used by the Programs Running, and Other Relevant Information	53
8.1.3	General Statistics	54
8.1.4	The Untrimmed Sarsa	55
8.1.5	The Untrimmed Q-learning	56
8.1.6	The Trimmed Sarsa	57
8.1.7	The Trimmed Q-learning	58
8.1.8	The One Dimensional Program	59
8.1.9	Phenomena: Total Time Use Is More than the Sum of the Training and Setup Time.	60
8.1.10	Phenomena: Irregular Decrease in Success-rate	62
8.1.11	Training From All States	63
8.1.12	Difference in Time-use, Comparing the Algorithms	64
8.1.13	Stochastic Termination	66
8.1.14	Playing the Game	66
8.2	Supervised Learning	68
9	Conclusion	69
9.1	Where Now?	69
10	The Code in Full	70

1 Intro

The game of Nim is rather old, how old is uncertain, but certainly over a hundred years old. It is also a pretty simple game, and it has a complete mathematical theorem, see “Nim, a game with a complete mathematical theorem”.¹ There is a fairly simple algorithm for making winning moves, so simple that a human can learn it by heart and use it. The author of this paper being an example of this. Machine learning is a type of stochastic algorithms that try to find a solution based upon statistical data. A stochastic algorithm is an algorithm with some random elements, a deterministic one does not have any random elements. It can often find good solutions quite quickly, certainly faster than conventional deterministic algorithms. Machine learning as well as other stochastic types of algorithms can often be good at finding good enough solutions, that is solutions which aren't perfect, yet still good. In this paper the problem which will be explored is how proficient it will be at playing perfectly, that is, having a chance at winning against a human player that knows how to play perfectly.

1.1 My Goal with This Thesis

What I am trying to do is to have a computer play Nim perfectly, with as little fore-knowledge as possible. This is reminiscent of computers playing other games, like chess. For example the chess-match between the Grandmaster Garry Kasparov playing against the computer Deep Blue. While having problems at first, it would eventually win a game, and then a full match against Kasparov. I am going to use stochastic methods to achieve this. This is different from Deep Blue, which used deterministic algorithms to choose moves.

1.2 Hypothesis

Making a program that uses reinforcement learning to play Nim should be able to do quite well, given enough time. I expect this because Nim, like other games, can be expressed in the form of a tree or a graph, and reinforcement learning is usually good at navigating graphs. On the other hand I believe it will be significantly slower than the pre-existing algorithm, at least in the worst case². Furthermore, I do not believe that the algorithm will be able to train enough with a static amount of episodes.

2 Ethics

Clearly, solving nim with for example, reinforcement learning, does not in itself have much of an impact on anything, ethically. This is because there is already

¹Charles L. Bouton. “Nim, A Game with a Complete Mathematical Theory”. In: *Annals of Mathematics* 3.1/4 (1901), pp. 35–39. ISSN: 0003486X. URL: <http://www.jstor.org/stable/1967631>.

²Worst case time use does not make much sense in the case of stochastic algorithms, like machine learning. I believe the difference here will be so large it is significant anyway.

an algorithm even a human could easily learn how to use. It would not change anything about the world, as there isn't a competitive Nim scene or anything like that, unlike for example chess. If one was able to solve chess, which is possible given the finite amount of positions, that could have a serious impact on professional chess-players. Not so here. Because of that, a look into possible indirect consequences of finding a solution here is warranted.

The long and short of it is that a solution to this problem might be a solution, or might inspire a solution for some unforeseen other problem. The traits of the problem of solving Nim, requiring perfection from the training, as well as finding the correct move in all cases. Not just playing the game from the starting state perfectly. There might be some other problem which share those traits.

3 The Game Nim

3.1 Rules and How to Play It

"The game is played by 2 players, A and B. Upon a table are placed three piles of objects of any kind, let us say counters. The number in each pile is quite arbitrary, except that is well to agree that no two piles shall be equal at the beginning. A play is made as follows:-The player selects one of the piles, and from it takes as many counters as he chooses; one, two, ..., or the whole pile. The only essential things about a play are that counters shall be taken from a single pile, and that at least one shall be taken. The players play alternately , and the player who takes up the last counter of counters from the table wins."³

A game of Nim shown as an algorithm Assuming A is the first player.

1. Player A removes some amount of counters from a pile. If there are still counters left, go to 2., if not, go to 3.
2. Player B removes some amount of counters from a pile. If there are still counters left, go to 1., if not, go to 4.
3. Player A removed the last counter/s and is the winner.
4. Player B removed the last counter/s and is the winner.

When playing Nim there are several ways to represent it, like in the movie "Last Year at Marienbad" where the main character shows how to play the game with cards, putting them in a pyramid fashion, and having the piles be the rows of cards. Alternatively literal piles of counters can also be used. What all of these have in common is that they are unary in representation, that is, they have n objects for a pile of size n .

³Bouton, "Nim, A Game with a Complete Mathematical Theory", op. cit., p. 35.

3.2 Nimrod

During the “Festival of Britain” in 1951 in the UK, in the science exhibition, there was a computer called Nimrod with which visitors could play a game of Nim. This might very well be the first time that a Digital computer(as they called it) could play Nim. It is also probably⁴ one of the first computer games to exist. It was a computer with 4 rows of lights indicating the size of 4 heaps. The computer also had a control panel which allowed human players to play with the computer. It had controls for different modes of play, as well as the same set of lights as on the computer, with companion buttons used to indicate the move the player wishes to make. An informational booklet was sold, which contained information about the computer and computers in general(at the time), the game Nim and how to play it, as well as an explanation on how the computer itself plays Nim. It contains the very algorithm on how to play Nim that will be presented later in this paper. The booklet was sold for 1 shilling and sixpence.⁵

3.3 Impartial Games

An impartial game is defined as a game where:

- There are 2 players.
- The players take alternate turns to make moves.
- There are a finite amount of unique legal moves for a player in any given state of the game, and there are a finite amount of possible moves in a game.
- All moves must be possible for both players to make. That is, in a given state, no matter which players turn it is, the same moves are possible for them.
- When a player cannot make a move, one of the players win.
- There is complete information, that is, both players know everything about the current state of the game.
- There are no random elements, like the use of dice or shuffled cards.

Nim is an example of such a game, because it fulfills all those requirements. It has, for example, always a finite amount of legal moves, because for each pile there are a finite amount of counters to take. This means there are only so many moves that can be made. Since counters have to be removed for each move, the game cannot go on forever, as there are a finite amount of counters to take.

⁴ It is hard to say exactly when and what the first computer game was, depending on several factors including the very definition of computer game. Considering it among the first seems likely however, given when it was available to the public. The question of what was the first computer game is a question worth a full article by itself.

⁵ *The Ferranti NIMROD Digital Computer.* URL: <https://www.goodeveca.net/nimrod>.

The longest possible game one could have is the game where each player takes one counter every turn, there will then be as many turns as there are counters. When a player takes the last counter they win, and the next player cannot make any moves, so that criteria is fulfilled as well. A player's moves are decided by the state of the game, not by the player whose turn it is, and so the requirement for moves being possible for both players is fulfilled as well. It seems obvious that Nim fulfills the rest of the requirements.

Games That Are Not Impartial Chess, for example, has the players control different pieces on the board which break the fourth rule, since each player can only move their own pieces. In Ludo there are more than 2 players, each player has their own pieces, and a random element, a die, is used. All of which individually puts it out of the category of impartial games. In Poker there is not complete information, because the cards in each hand are hidden. Furthermore the shuffling and dealing of the deck is a type of random element, if done correctly⁶.

3.4 The Sprague-Grundy Theorem

The Sprague-Grundy Theorem states that all impartial games are equivalent with a nimber. A nimber is defined as the value of a pile in a game of Nim.⁷ The Nimber of a game of Nim is simply the Nim-sum of the game, which will be explained later. Because of this theorem, it seems easy to conclude that whatever solutions are found, they should be equally as good at playing other impartial games.

4 Machine Learning

The book “Machine learning: an algorithmic perspective” states that “Machine learning, then, is about making computers modify or adapt their actions (whether these actions are making predictions, or controlling a robot) so that these actions get more accurate, where accuracy is measured by how well the chosen actions reflect the correct ones”.⁸ In this paper only reinforcement learning and supervised learning will be used, and as such those are the only parts of machine learning that will be discussed here. Note that this will not be an exhaustive explanation, only on the parts of the subjects which are relevant, mostly those in use.⁹

⁶Shuffling cards in such a way that the cards end up in a predetermined way is possible, but this would be considered cheating. You'd get shot for less in Texas.

⁷John H. Conway, Elwyn R. Berlekamp, and Richard K. Guy. *Winning Ways for Your Mathematical Plays*. Second. Vol. 1. CRC Press, 2001, p. 56.

⁸Stephen Marsland. *Machine learning: an algorithmic perspective*. Second. CRC Press/- Taylor & Francis Group, 2015, p. 4.

⁹The information in this section is taken freely from ibid.

4.1 Reinforcement learning

Reinforcement learning is a type of machine learning where a program tries to find a solution for a problem, and is taught how to do this, by getting rewards and punishments. It is based upon how one might train an animal with treats and such, to do a trick.

4.1.1 The Principle

The logic behind reinforcement training is simple enough. When training a dog, for example, teaching it to do a trick, you might follow the following procedure:

1. Tell the dog how to do the trick, likely showing it what it is supposed to do, somehow.
2. Give the order for the trick. If the dog does it well enough, it gets a reward. This reward could be a treat, for example.
3. Repeat the previous step until the dog manages to do it consistently.

This might teach a dog to sit, or roll around. For such a simple trick, that is fine, but what if you wanted to teach a dog to navigate a maze? Dogs have memories, so even if you only give them treats when they reach the end, they might remember which route they took, and take it again to get treats. All this is well and good, but a program is not a dog, so its methods will differ from a dogs.

4.1.2 The Proper

The way a program learns in reinforcement learning, is by trying somewhat randomly, and then compile statistics about what reward it can expect from a certain action. It might have, for example, a matrix containing the expected reward from moving between states. The algorithm works with the concept of episodes, which are the program finishing a run, when the program has found its way through a maze. A more proper procedure might look like this:

1. The amount of episodes is, or has been chosen.
2. Repeat for each episode.
 - (a) The program starts in some start position, like in the middle or at the entrance of a maze.
 - (b) A move is chosen, usually either a random move, or the move the program thinks is best. The one with the highest expected reward.
 - (c) The move is made, and the expected reward of the move is updated based upon the viability of the possible next move. If the program now sees the exit of the maze it was a good move.

- (d) Now the program is in a new state. If the program has reached the end, escaped the maze, then the episode is done, if not, the program is back to step (b) again.

4.1.3 Action Selection/Policy

The action selection, or policy, is how the program chooses which move to make when training. There are several ways to do this, but two of them are:

- Greedy, where the program always picks the move, or action, which gives the highest expected reward.
- ϵ -greedy, which is the same as greedy, except there is a chance to instead make a random move.

If only the Greedy selection is used then the program might get stuck at a less than optimal solution, because it never tries something other than what it knows. This is why ϵ -greedy is the one used further on in this paper.

4.1.4 Sarsa and Q-learning

There are two relevant specific algorithms for training in reinforcement learning, the two used in this paper. All of the algorithms are based upon the so called value function, of which there are also two types, the state-value function $V(s)$ which finds the average expected reward for a state, and the action-value function $Q(s, a)$ which finds the expected value of an action, or a move from a state s to state a . Both algorithms use the action-value function.

Q-learning First we have to Q-learning algorithm, which looks like this:¹⁰

- Initialize
 - set $Q(s, a)$ to small random values for all s and a
- Repeat:
 - initialize s
 - repeat:
 - * select action a using *epsilon*-greedy or another policy
 - * take action a and receive reward r
 - * sample new state s'
 - * update $Q(s, a) \leftarrow Q(s, a) + \mu(r + \gamma \max_{a'} Q(s', a'))$
 - * set $s \leftarrow s'$
 - For each step of the current episode
- Until there are no more episodes

¹⁰Ibid., p. 242.

Sarsa The algorithm assumes that the optimal move will be made in the next state, unrelated to what the actual policy in use is. It is thus considered off-policy. If one wants an algorithm which is instead on-policy, it must be modified slightly, into one which is called Sarsa:¹¹

- Initialize
 - set $Q(s, a)$ to small random values for all s and a
- Repeat:
 - initialize s
 - choose action a using the current policy
 - repeat:
 - * take action a and receive reward r
 - * sample new state s'
 - * choose action a' using the current policy
 - * update $Q(s, a) \leftarrow Q(s, a) + \mu(r + \gamma Q(s', a') - Q(s, a))$
 - * $s \leftarrow s'$, $a \leftarrow a'$
 - For each step of the current episode
- Until there are no more episodes

It is called Sarsa because the algorithm use variables: s , a , r , s' , and a' , thus spelling Sarsa.

4.2 Supervised Learning

Supervised learning is a type of machine learning that relies on the use of training data, that is data that states input and output. This means the program knows what it needs to answer for many inputs, and for it to work it needs to be able to generalize on that data, use it to make a model that will hopefully apply in all cases. The focus will be on neural nets, because that is what is used later on in this paper.¹²

Neurons A neuron is, put simply, a kind of valve for signals. It takes the sum of all its incoming signals, inputs, and if the sum is above some threshold it fires off, sending out a signal of its own, if the sum of the input is below the threshold, it does not fire. McCulloch and Pitts made a simple mathematical model of a neuron¹³, which is used as a basis to make more advanced neural nets. It looks like Figure 1, and calculates output in the following procedure:

¹¹Ibid., p. 243.

¹²The information in this section is taken freely from ibid.

¹³Note that it is a simple mathematical model, it isn't particularly accurate, and doesn't emulate the more strange behaviours of nerve-cells.

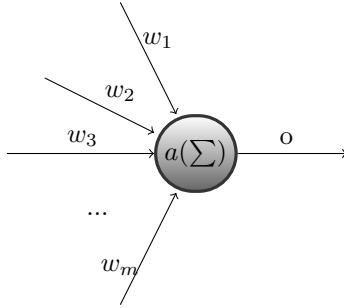


Figure 1: A model of a McCulloch Pitts neuron.

1. get the sum of the inputs times their weight. $\sum_{i=1}^m w_i x_i = s$
2. apply the activation function with the sum as input. $a(s) = o$
3. The result, o , is output through the synapse.

The activation function is typically a threshold function, that is, a function which returns almost exclusively a number close to either 1 or 0. This models biological neurons, which can either send a signal, or not send a signal. This activation function does not have to be a threshold function, however, in which case the neuron might behave somewhat differently than a biological neuron, for example sending output like 0.5, which does not make sense from a biological perspective.

4.2.1 The Perceptron

The perceptron is a simple neural net, that is it is a network of neurons. It essentially is a row of neurons, that all get the same inputs, but the weights are all different¹⁴. Figure 3 shows a model of a perceptron. Note that there is just 1 layer of neurons, the darker nodes, but those neurons are the McCulloch Pitt neurons mentioned previously.

Bias If the input for a perceptron was all 0, the output would be the same no matter what the weights were, since $x * 0 = 0$. The bias is included so the perceptron can change the output for input of all 0's. The bias is simply another input which always has the same value. The value can be anything non-zero, but 1 or -1 are usually chosen for simplicity's sake, though it does not matter as the weight could be adjusted accordingly.

Training The perceptron can be trained, for each weight it should be updated like in equation 1, where x_i is the input from node i , y_j is the output from neuron

¹⁴Usually. They can in theory all be the same, or at least some be the same, it is just rather unlikely unless they are initialized as such

j , and t_j is the expected output from neuron j . η is the training variable, which is just adjusted as needed to train as one wish.

$$w_{ij} \leftarrow w_{ij} - \eta(y_j - t_j) * x_i \quad (1)$$

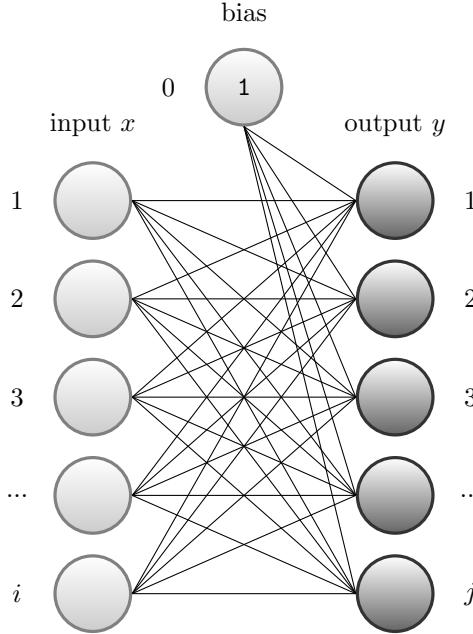


Figure 2: A model of a single layer Perceptron. Note that only the darker nodes are actual neurons, the lighter ones are there to show how the inputs are given to all neurons.

Limitations A perceptron can find patterns, and fire for something that belongs to a class, and not fire for something that doesn't belong to a class. It is constrained to separating them with a straight line, if in 2D, and thus is quite constrained in what it can identify.¹⁵ Suffice to say that it isn't very powerful, and not very useful. There is a solution to this however, simply have more layers. That solution does have some problems, specifically, how to train a so called Multilayer perceptron.

4.2.2 The Multilayer Perceptron

To make the perceptron more useful, we now have another layer. The problem with doing this is that the algorithm needed to train it is rather more complicated than before. The proper equations will be the same used later on when reviewing the code, so for now a more abstract explanation should suffice:

¹⁵Marsland, *Machine learning: an algorithmic perspective*, op. cit., p. 55.

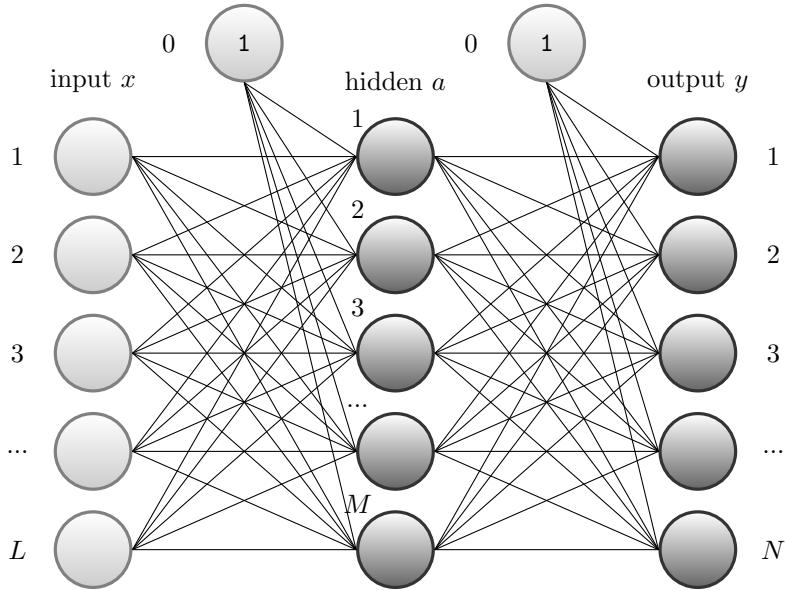


Figure 3: A model of a MultiLayer Perceptron.

1. An input is given to the input nodes.
2. The inputs are fed forwards through the network.
 - (a) The hidden layer get the inputs, and like the single layer perceptron, neurons take the inputs times the weights, and applies that sum to the function. The result is given as input to the output layer.
 - (b) The output layer does the same, only with the data from the hidden layer instead.
3. Calculate the error of the output.
 - (a) The weights between the output and hidden layer are updated, based upon the error, and their input.
 - (b) Then the weights between the hidden and the input layer are updated, using error inherited from the output layer, modified by the weights. The bigger value a neuron gave, the bigger amount of error it inherits.

The proper algorithm¹⁶ The actual algorithm is somewhat more complex, and use the following symbols: i, j, k is the indices of the layers respectively, their sizes being L, M, N respectively, with the corresponding greek letter ι, ζ, κ for fixed indices. Weights between the input and hidden layer is called v , and the

¹⁶Ibid., p. 77.

weights between the hidden and output layer is called w . The input nodes are x_i , the bias node is x_0 , the hidden nodes are a_j , the bias node a_0 and the output layer is y_k . The targets for the outputs, what the mlp should output, is t_k . The proper algorithm for the Multilayer perceptron will then be:

The Multilayer Perceptron Algorithm¹⁷

- **Initialization**

- initialize all weights to small (positive and negative) random values

- **Training**

- repeat:

- * For each vector:

Forwards phase:

- compute the activation of each neuron j in the hidden layer(s) using:

$$h_\zeta = \sum_{i=0}^L x_i v_{i\zeta} \quad (2)$$

$$a_\zeta = g(h_\zeta) = \frac{1}{1 + \exp(-\beta h_\zeta)} \quad (3)$$

- work through the network until you get to the output layer neurons, which have activations¹⁸

$$h_\kappa = \sum_{j=0}^M a_j w_{j\kappa} \quad (4)$$

$$y_\kappa = g(h_\kappa) = \frac{1}{1 + \exp(-\beta h_\kappa)} \quad (5)$$

¹⁷Ibid., p. 78.

¹⁸ The actual equation in the source is $h_\kappa = \sum_j a_j w_{j\kappa}$ but that is lacking a couple of symbols, as j needs constraints.

Backwards phase:

- compute the error at the output using:

$$\delta_o(\kappa) = (y_\kappa - t_\kappa)y_\kappa(1 - y_\kappa) \quad (6)$$

- compute the error in the hidden layer(s) using:

$$\delta_h(\zeta) = a_\zeta(i - a_\zeta) \sum_{k=1}^N w_\zeta \delta_o(k) \quad (7)$$

- update the output layer weights using:

$$w_{\zeta\kappa} \leftarrow w_{\zeta\kappa} - \eta \delta_o(\kappa) a_\zeta^{hidden} \quad (8)$$

- update the hidden layer weights using¹⁹:

$$v_{\zeta\kappa} \leftarrow v_{\zeta\kappa} - \eta \delta_h(\zeta) x_\kappa \quad (9)$$

* (if using sequential updating) randomise the order of the input vectors so that you don't train in exactly the same order each iteration

– until learning stops

- Recall

– use the Forwards phase in the training section above

5 The Algorithms for the Problem

5.1 Preexisting Algorithm

This is the algorithm which already exist for solving Nim, and this is how it works. This algorithm is deterministic, that is, it will give the same result and follow the same steps each time, if given the same input twice.

5.1.1 Nim-sum

The algorithm uses an operation called Nim-sum, which is really equivalent to the exclusive-or operation if the numbers are converted to binary. I shall use \oplus as the symbol for this operator for every digit in the two numbers, the corresponding digit in the Nim-sum will be 1 if they are different, and 0 if they are the same. Table 1 is an example of a nim-sum operation. When finding the Nim-sum of more than 2 numbers all that is done is to get the Nim-sum of the first 2 numbers, then find the Nim-sum of that and the third number, continuing that until finished.

The book “Winning ways for your mathematical plays”²⁰ has the following suggestion for how to calculate a nim-sum if you are a human, not a computer:

¹⁹ The actual equation in the source is $v_\kappa \leftarrow v_\kappa - \eta \delta_h(\kappa) x_\kappa$ but that doesn't make sense, as a weight is between 2 nodes, but here the weights aren't like that.

²⁰Conway, Berlekamp, and Guy, *Winning Ways for Your Mathematical Plays*, op. cit.

Table 1: A Nim-sum example

Binary	base10
00100110 \oplus	38 \oplus
00001101	13
=	
00101011	43

The Basics of Nim-addition You can use these two basic properties to find the nim-sum of any collection of numbers by writing each of them as a sum of distinct powers of 2 and then cancelling repetitions in pairs. For example,

$$5 \oplus 3 = (4 + 1) \oplus (2 + 1) = 4 \oplus 1 \oplus 2 \oplus 1 = 4 \oplus 2 = 4 + 2 = 6,$$

$$11 \oplus 22 \oplus 33 = (8 + 2 + 1) \oplus (16 + 4 + 2) \oplus (32 + 1) = 8 + 16 + 4 + 32 = 60.$$

Table 2: This table shows the decision process for each digit in the nim-sum operation.

only one digit	only one digit	different digits	same digits	different digits	same digits	
1	0	0	1	0	1	\oplus
1	0	1	0	1	0	=

5.1.2 Nim-sum of the State

the crux of the algorithm is to find the Nim-sum of the current state of the game. It will then be used to figure out which move is the correct to make, unless the Nim-sum is 0, because then there is no good move. The Nim-sum of a game-state is simply the Nim-sum of all the sizes of the heap's in the game. For the game-state 3 5 7 the Nim-sum is $3 \oplus 5 \oplus 7 = x$. Table 3 shows the calculation with the answer: 1.

Table 3: Nim-sum of the state 3 5 7

Binary	base10
011 \oplus	3 \oplus
101 \oplus	5 \oplus
111	7
=	
001	1

5.1.3 Finding the Correct Move

To find the correct move one should get the nim-sum of the state, and find the Nim-sum of that and the size of a heap. If the resulting sum is less than the heaps size, the correct move is the move reducing the heap to the Nim-sum. If the sum is not smaller than the heap size, the same procedure should be done with another heap. As long as the Nim-sum of the stat is not 0, a move can be found. In Table 4 a correct move is found in the case of game-state 3 5 7 which has a Nim-sum of 1. The Nim-sum is 2, which means a correct move would be to reduce the heap of size 3 to size 2. There are other correct moves here too, however. For example: $1 \oplus 5 = 4$, so reducing the heap of size 5 to size 4 is also a correct move.

Table 4: Nim-sum of Table 3s Nim-sum and the first heap in the game-state.

Binary	base10
$001 \oplus$	$1 \oplus$
011	3
=	
010	2

5.2 Reinforcement Learning

This second algorithm is a machine learning algorithm, and it is, as opposed to the previous one, stochastic.

5.2.1 Representation of Nim as a graph

Reinforcement learning usually works with a graph, to find the best way through the graph. Therefore the game of Nim should be represented as a graph so that the algorithm might be more easily applied to this problem. All games can easily be represented as trees, where all the children of a node represents all the legal moves possible to make from the parent node. See the first graph in Figure 4. This is not very space-efficient however, as there can be several nodes which are identical, which could be eliminated for a slimmer graph. As small a graph as possible is desired, specifically fewer edges, so less training is needed. The amount of nodes will mainly determine the space-use of the program. The graph is directed and has nodes and edges, where the a node is a state the game can be in, that is a set of heaps, each with some amount of counters in them, including 0 counters. There is an edge from one node to another, if there is a legal move in the first state which results in the second state. The node which does not have any edges going into it is the node containing the start state, the state which is given as the input for the algorithm. The node with only incoming edges is the winning condition, that is, no counters left in any heap.

This is the graph which the first iteration of the algorithm uses, see the second graph in Figure 4.

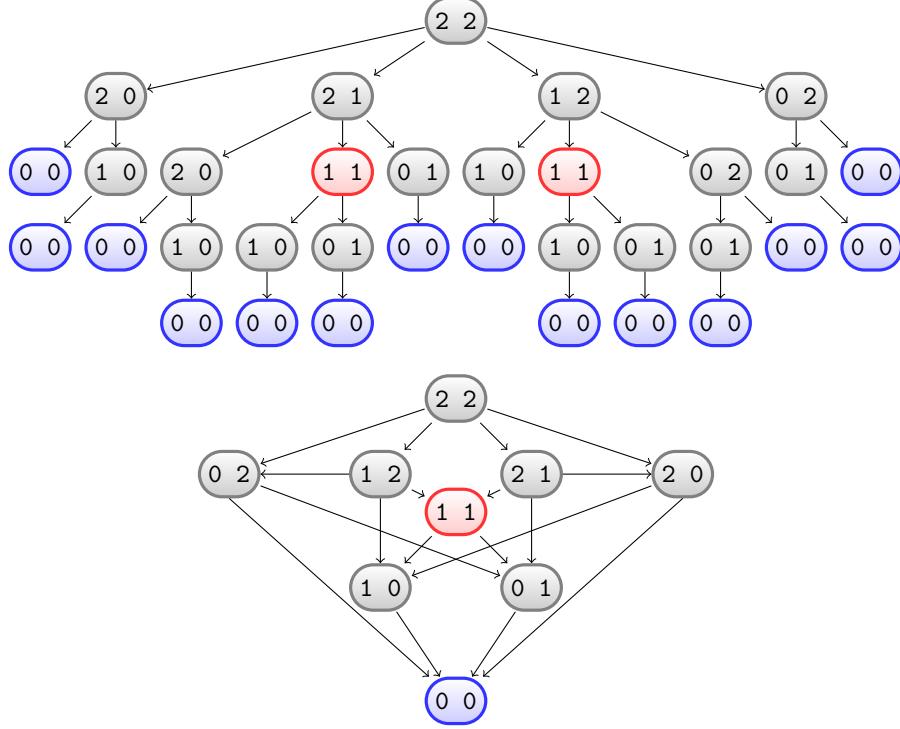


Figure 4: First graph is the graph for the game tree of a game of Nim, starting with the state 2 2. Notice how many times state 0 0 (blue) occurs. State 1 1 (red) occurs twice as well. The second graph is the resulting graph if all equal state-nodes in the game tree is combined. Notice how the states 0 0, and 1 1 occur just once in this graph.

In the first graph in Figure 4 the node 1 1, in red, occurs twice, and both have the same amount of edges going out from them. In the second graph in Figure 4, where there is only one occurrence of 1 1, there are also less edges in total going out of occurrences of 1 1, as well as there being less occurrences overall. This does not apply to all nodes however; node 0 0, in blue, never has any edges going away from it, so even though it is the most common node in the first graph Figure 4 removing duplicates of it will not as such remove edges, it will still make the total memory requirement smaller.

Amount of nodes the amount of nodes in a graph is calculated using equation 10 where h_k is the size of the k th heap

$$(h_0 + 1) * (h_1 + 1) * (h_2 + 1) * \dots * (h_k + 1) = \text{amount of nodes} \quad (10)$$

5.2.2 How the Graph Is Represented in the Program

This is a matrix showing all legal moves from a state in the vertical list of states, to a state in the horizontal list of states. The program gets the second graph in figure 4, modelled as a matrix, where the indices on both axis corresponds to one state of the game, or one node of the graph. Index 0 is the starting position which the program gets as an input, while index 1 is the state you get when one counter is removed from the last heap of the state in index 0. Each field represents whether there is an edge from the state on the index on the horizontal axis, to the state on the index in the vertical axis. The field is of the value 1 if there is an edge, and of value 0 if there isn't an edge. The program also has a second matrix of the same size, which will be used to record expected reward for a certain move.

Table 5: The graph as a matrix

states	2 2	2 1	2 0	1 2	1 1	1 0	0 2	0 1	0 0
2 2	0	1	1	1	0	0	1	0	0
2 1	0	0	1	0	1	0	0	1	0
2 0	0	0	0	0	0	1	0	0	1
1 2	0	0	0	0	1	1	1	0	0
1 1	0	0	0	0	0	1	0	1	0
1 0	0	0	0	0	0	0	0	0	1
0 2	0	0	0	0	0	0	0	1	1
0 1	0	0	0	0	0	0	0	0	1
0 0	0	0	0	0	0	0	0	0	0

5.2.3 How the Graph Is Made Smaller

The graphs are much better than the trees, but there is still overlap between the nodes. For example the nodes 2 1 and 1 2 are essentially the same, the only difference is that they are mirrored. So there is still something that can be done to make the graph even smaller, by only having one ordering of any state, that is, only 2 1 and not 1 2. I have decided that the correct ordering is from biggest to smallest. So 2 1 1 is allowed, but not 1 1 2, nor 1 2 1. The graph for 2 2 will then look like the right graph in Figure 5.

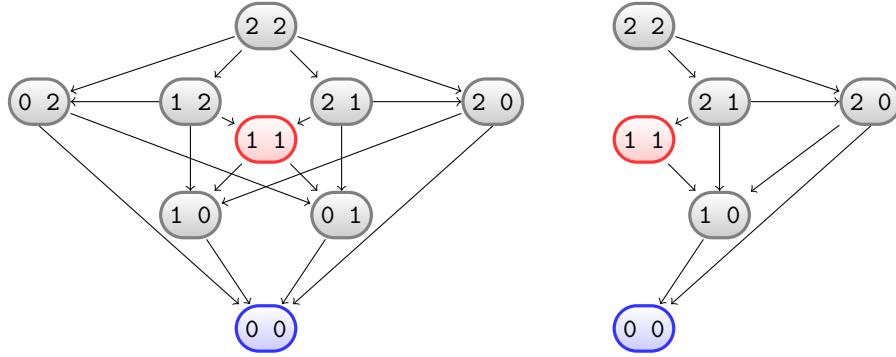


Figure 5: Graph to the left is the untrimmed graph, the one to the right the trimmed one. See how the trimmed one has almost been cut in two? It is much smaller, another improvement to the graph.

Equation 11 is used to calculate the amount of possible states reachable in a game of nim, provided all the heaps are of the same size. n is the size of the heaps, and k is the amount of heaps.

Amount of nodes in a trimmed game-tree

$$\frac{(n+k)!}{n! * k!} \quad (11)$$

Combinations with repetition

$$\frac{(n+k-1)}{k! * (n-1)!} \quad (12)$$

This equation is almost the equation for finding combinations with repetition, which it is based upon. This makes sense, because of the number needed can be considered to be the following: for each heap draw a number between and including the size of the heaps and 0. If the order is not important, how many different combinations of these are there. Since the amount of things to draw is actually 1 more than the size of the heaps, because 0 can also be drawn, the equation is changed a little, because n is *amount of heaps* + 1 and n is always preset with a -1 which nullifies each other.

Problems with the trimmed graph The big problem with the algorithms for the trimmed graph is that they can actually only take start-states with heaps of the same size, as it needs to know the size of the graph from the start. This is not so bad, however, the benefits of using this trimmed graph far outweighs any restrictions, see section 8.

5.2.4 Training

The algorithm learns what moves to make by training, which consists of trying moves, and updating the expected rewards for those moves. The algorithm essentially plays games of Nim against itself, and modifies the aforementioned second matrix to get a better estimation of how good certain moves are, but there is a difference from the common algorithm for reinforcement learning. The difference is that while usually the expected reward for a move is modified by adding it to the expected reward of some possible next move(depending on the variant), because moving to a good position is usually desired for it moves you closer to some goal. In this situation, and for a similar algorithm for other similar games, a move to a good position is not desirable, because that means the opponent gets to move from a good position. In this case instead of adding the expected reward for some next move, it is subtracted. This would only apply to games where both players can make the same moves, but not games like chess where each player have different pieces. The assignment looks like this: $Q(s, a) \leftarrow Q(s, a) - \mu(r + \gamma Q(s', a') - Q(s, a))$

To make it more clear, here is what happens, in the correct order.

1. The amount of episodes is, or has been chosen. For example 1000 episodes.
2. An episode is started, and the number of episodes is incremented.
 - (a) An episode starts on the predefined start position. In Table 5 it is 2 2.
 - (b) With the current state of the game, a move is chosen. It will usually pick the best move, that is the one with the highest expected reward, but there is a chance for it to be picked randomly.
 - (c) The move is made, and the expected reward of the move is updated, based upon the move the next player is expected to make. If that next move is good, the reward will be decreased if not appropriate.
 - (d) The new state, which is reached by the chosen move, is now considered the current state. If the current state is of a finished game, that is, it consists of only empty heaps, then the episode is done. If not, go back to step (b) with the new current state.
3. If the number of episodes has reached the predefined amount, in this case 1000, the program is done. If the amount is not reached, go back to step 2.

Terminology From this point on the programs that use the untrimmed graph, the one that combined the nodes in the game-tree that were the same,will be referred to as the untrimmed programs. The graph will also merely be called the game-tree. The programs that use the trimmed graph where the order is not enough for difference are called the trimmed programs. Any other further programs will be named for their unique features.

5.2.5 Terminating the Program

The program runs a finite amount of episodes before it terminates. The amount of episodes is configurable in the code. This wasn't the only possibility however, the program could also have ran on until it was sufficiently trained(whatever that is). These are the options of a logical or stochastic termination, the logical is predetermined, and the same each time the program is run, and the stochastic is(likely) different each time the program is run. The logical termination was chosen because of its predictability, as well as the sizes of the game-trees for larger games: the amount of episodes of training that is necesarry gets so large the program runs for an impractical amount of time. Since I wish to easily make statistics, logical termination seems appropriate, with the added benefit of the amont of episodes being controlled as well. A test-program that features stochastic termination will be featured later in the paper, and the results should support my decision.

5.2.6 One-dimensional Program

Since the states in Nim have the same legal moves for both players, and the goal for finding a good move is to find a good state to move to, perhaps an improvement would be to base a program on the use of the state-value function $V(s)$ instead of the action-value function $Q(s, a)$. Since it uses the one-dimensional $V(s)$, some parts of the code must be changed, but overall the program should be quite similar. It shall be called OneDSarsa, as it uses the Sarsa algorithm adapted to $V(s)$, which is one-dimensional.

5.3 Supervised Learning

This algorithm relies on using data on correct moves to learn how to play.

5.3.1 Acquiring Data

Acquiring data is simple enough, since there exists an algorithm for making correct moves, so a way would simply be to have a list of states, and record what the algorithm outputs for those states. Another way would be to use the reinforcement algorithm to get some positions.

5.3.2 Training

The program trains by modifying the weights in a graph, or network, called a neural network, as have already been explained. The algorithm works somewhat like this, see figure 6 for an example neural net:

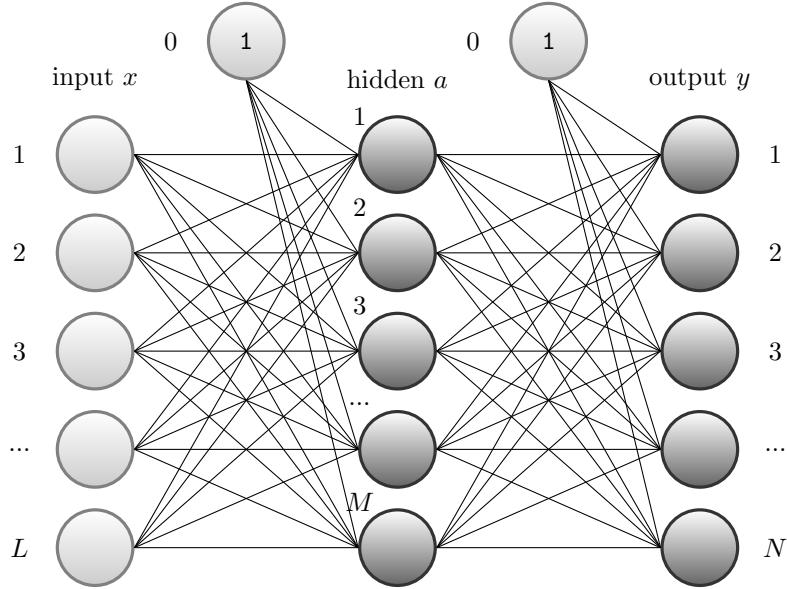


Figure 6: A simple model of a Multilayer Perceptron.

1. The weights are initialized with small random positive or negative numbers.
2. Repeat:
 - (a) Repeat once for each line of data in the training set:
 - i. Give the current line of training data as input.
 - ii. Calculate inwards to the hidden layer, through each synapse. From each synapse multiply the input value with the weight, and add it to the output node.
 - iii. Calculate inwards to the output layer, through each synapse. From each synapse multiply the input value with the weight, and add it to the output node.
 - iv. Compare the resulting output with the expected output and calculate error.
 - v. Go backwards from the output layer to the hidden layer and adjust weights based on their error.
 - vi. Go backwards from the hidden layer to the input layer and adjust weights based on their error.
 - (b) Check the total error of the neural net.
3. If the error is higher than last iteration, do not repeat. Training is done.

Total error To reiterate, the total error of the neural net is checked by comparing the result of some input and the expected correct result. Equation 13 is the equation for the total error of the output layer of a neural net, where N is the amount of nodes, or perceptrons in the output layer.

$$\sum_{k=1}^N (y_k - t_k)^2 \quad (13)$$

6 Describing the Code

Having given a more abstract explanation on how the algorithms work, it is time to take a more indepth look into the actual programs i have written²¹.

6.1 The Deterministic Algorithm

The code for the deterministic algorithm is simple enough. It takes the state of the game as input, and returns the correct move as an array of size 2 containing firstly how many counters should be removed, and secondly from which heap the counters should be removed from.

```

1  class PerfectPlayer:
2      #This method applies the deterministic algorithm for making
5         perfect moves in Nim, and returns the index of the move and the
6         amoount of counters removed.
7
8      def makeMove(self, state):
9          nimSum = self.findNimSum(state)
10         for i in range(len(state)):
11             remove = nimSum^state[i]
12             if(remove < state[i]):
13                 result = [state[i] - remove, i]
14             return result
15         for i in range(len(state)):
16             if(state[i] != 0):
17                 result = [1, i]
18             return result
19         return -1
20
21
22     #Since finding nimsums is relevant in the program Run.py this
23     #program finds the nimsum of the input state.
24
25     def findNimSum(self, state):
26         nimSum = 0
27         for i in range(len(state)):
28             nimSum = nimSum^state[i]
29         return nimSum

```

Listing 1: Algorithm for finding the correct move from a certain state in Nim

²¹There are some parts of the code where some outside inspiration has had its part, though those parts will be mentioned as such, when relevant.

The program could have returned the new size of the heap instead, however that seemed less natural, as it feels like a move should be considered to remove a certain amount of counters, as that is what one would do when playing the physical game.

6.2 Programs for Reinforcement Learning.

I have written five programs that use reinforcement learning. Two that use the untrimmed game-tree, and two that use the trimmed tree, and one that has a one-dimensional array for the expected rewards, that is it has an expected reward for a state, rather than a move. It also uses the trimmed game-tree. There are some functions they have in common, so they will be mentioned. I will mention them in the same order as here.

6.2.1 Initialization

There are two versions of the initialization. One in the programs using the untrimmed tree, and one in the programs using the trimmed tree.

Untrimmed The initialization for these programs looks like this:

```

8 def __init__(self, startBoard):
9     self.sum = startBoard[0]+1
10    self.board = startBoard
11    for i in range(1, len(self.board)):
12        self.sum *= self.board[i]+1
13    self.listOfStates = [0] * self.sum
14    self.listOfStates[0] = self.board[:]
15    self.T = []
16    self.Q = []

```

Listing 2: Initialization of the programs using the untrimmed game-tree

The important part is to find the sum, which is really the amount of nodes in the game tree, or how many possible states can be reached from the starting state. In the untrimmed game tree it is easy to calculate, as shown before. see Equation 14. T and Q are the matrices for legal moves and expected rewards from moves repsectivly, and listOfStates is appropriately enough the list of possible states, with the rule that any move that can be made from a state, is only to states below it in the list.

$$(h_0 + 1) * (h_1 + 1) * (h_2 + 1) * \dots * (h_n + 1) = \text{amount of nodes} \quad (14)$$

Trimmed This only accepts input boards of the type where all the heaps are of the same size. The initialization for these programs looks like this:

```

9 def __init__(self, startBoard):
10     self.board = startBoard
11     self.sum = self.getSum()
12     self.board.sort(reverse = True)
13     self.listOfStates = []
14     self.listOfStates.append(self.board[:])
15     self.T = []
16     self.Q = []

```

Listing 3: Initialization of the programs using the trimmed game-tree

This is almost the same as the last one, except that the sum is calculated differently, it has its own function. the function just uses Equation 15 to calculate the sum.

```

19 def getSum(self):
20     return int(math.factorial(len(self.board)+ self.board[0])/(math.
factorial(len(self.board))*math.factorial(self.board[0])))

```

Listing 4: A function for applying Equation 15 to find the amount of states reachable from the start state.

$$\frac{(n+k)!}{n! * k!} \quad (15)$$

6.2.2 Setup

The setup is also divided by the game-trees, but the program with the one-dimensional array for expected rewards are also separate here.

Untrimmed The setup is for properly making and filling the matrices T and Q, so that training can start. It runs the function to start training, it also records the time setup and training use for statistical purposes.

The List of States All the programs use a list of states for several purposes, for example, when it is filled, that also functions as a method for filling T with proper data on the legal moves. It is also necessary for when you make moves. The problem is that in the setup of all these programs, the list is filled by making all legal moves from each state, by making every legal move from each state. This also fills the list of states, which is the source of the problem. The problem is whether, when some index of the list is reached, has that state been filled in? The following theorem shall show that with a certain way of ordering the list, this is the case. The list shall be ordered in such a way that each heap, except the right most, the order is so that for state $h_0, \dots, h_i, \dots, h_n$, for each state h_i to decrement, state h_{i+1} must decrement until it is empty. since there is no state h_{n+1} , state h_n can decrement without any restrictions. The order created by that, when decrementing h_0 is the order list of states should have. In the case of the trimmed tree, the order is the same, except with any state out of order removed. Another way to state to question is whether some state does not have

states above it in the list, that has legal moves to it. The only exception would be the start state. I can do better, it can be proven that no state is, in the list, under any state it has a legal move to.

Theorem: All States Are Above Any State to Which They Have a Legal Move This can be proven by proof of contradiction. Let us assume that some state has a legal move that ends up in a state above it in the list. In that case we have a state $s: h_0, \dots, h_i, \dots, h_n$, and the result of the move, state $s' h_0, \dots, h_i - x, \dots, h_n$, where $x > 0 \ \& \ x = < h_i$, which is earlier in the list. This is a contradiction, because the order is decided by decrementing heaps, and since h_i is bigger than $h_i - x$, and all the other heaps are the same, state s must by definition be above state s' . This is a contradiction, thus the assumption is wrong. For when the list is the trimmed one, the same proof almost holds. The difference is that the result of the move does not have the heaps in the same place as they had earlier, for example the move from 2 2 1 to 2 1 0. So the logic is almost the same, you have state s and s' , in the formats $h_0, \dots, h_i, \dots, h_n$ and $h_0, \dots, h_i - x_0, h_{i+1} - x_1, \dots, h_n - x_{n-i}$ respectively, where $0 < \sum_{y=0}^{n-i} x_y \leq h_i$. This causes a contradiction, because there will be at least one case of $h_j > h_j - x_{j-i}$, which like the previous time, means the state s should be earlier than s' according to the predefined order.

```

23 def setup(self, type):
24     timeStart = time.time()
25     for i in range(self.sum):
26         self.T.append([0] * self.sum)
27         for x in range(len(self.listOfStates[i])):
28             for y in range(1, self.listOfStates[i][x] + 1):
29                 tmpValue = self.newIndex(x, y)
30                 self.T[i][i+tmpValue] = 1
31                 if(self.listOfStates[i+tmpValue] == 0):
32                     tmpList = self.listOfStates[i][:]
33                     tmpList[x] -= y
34                     self.listOfStates[i+tmpValue] = tmpList[:]

35         self.Q.append([])
36         for j in range(self.sum-1):
37             if(self.T[i][j] == 1):
38                 self.Q[i].append(random.random())
39             else:
40                 self.Q[i].append(-np.inf)
41             if (self.T[i][self.sum-1] == 1):
42                 self.Q[i].append(100.0)
43             else:
44                 self.Q[i].append(-np.inf)
45     timeEnd = time.time()
46     setupTime = timeEnd-timeStart
47     if(type == 0):
48         epsilon = 0
49         mu = 0.5
50         gamma = 0.5
51         amountEpisodes = 1000
52         timeStart = time.time()
53         self.train(amountEpisodes, epsilon, mu, gamma)
```

```

55     timeEnd = time.time()
56     trainTime = timeEnd-timeStart

58     times = []
59     times.append(setupTime)
60     times.append(trainTime)

62     return times
63     return -1
64

66 def newIndex(self, heapIndex, amount):
67     newMove = 1
68     for i in range(len(self.board)-1, heapIndex, -1):
69         newMove *= self.board[i]+1
70     newMove *= amount
71     return newMove

```

Listing 5: The setup function for the untrimmed programs.

It works by going through the list of all states, adding all states directly reachable from the current state in 1 move, and adding that as a legal move in T, and filling in the new states in the list listOfStates. Since all states in the list except the first one has to be reachable from some other state in the list, we can conclude that when an arbitrary states index is reached, it has been added by some previous state. knowing which index a state belongs to is quite simple, and works as shown in Equation 16 where k is the index of the last heap, b is the size of a heap before, and a is the size after the move has been made. s is a heap as it is in the startstate. The function *newIndex* simply gets a move as input, and calculates what the index of the new position is, relative to the state the move was made from. That is used to calculate where the new states should be, and if that position is empty, that is the state has not already been added, it will be added.

$$i = (b_k - a_k) + (b_{k-1} - a_{k-1}) * (s_k + 1) + \dots + (b_1 - a_1) * (s_k + 1) * (s_{k-1} + 1) * \dots * (s_1 + 1) \quad (16)$$

When the matrix T is made, the program iterates through it, and makes the Q matrix based upon it. For any legal move, a small random number between 1 and -1 is chosen, and for illegal moves the value negative infinity is chosen²². The expected reward to move to the winning state starts off as 100, so that the program knows to go there. The setup function then sets the values for the training variables mu and gamma used by the training function, as well as epsilon for the ϵ -greedy choice, and lastly the amount of episodes.

²²Negative infinity is a feature of numpy, it is only used so that when picking the move with the highest expected reward, an illegal move is never picked, as negative infinity is always smaller than any number.

Trimmed The setup function for these programs are similar, from the last line in the function in Listing 6 to the end of the function they are the same. It uses a function called *newState*, which takes a state and a move, and gives the resulting state, and gives it back, ordered correctly.

```

45 def setup(self):
46     timeStart = time.time()
47     emptyBoard = [0] * len(self.listOfStates[0])
48     for i in range(self.sum):
49         self.T.append([0] * self.sum)
50         for x in range(len(self.listOfStates[i])-1, -1, -1):
51             atEnd = 0
52             curIndex = i+1
53             for y in range(1, self.listOfStates[i][x]+1):
54                 newState = self.newState(x, y, self.listOfStates[i])
55                 if(atEnd):
56                     self.T[i][curIndex] = 1
57                     self.listOfStates.append(newState)
58                     curIndex += 1
59                 else:
60                     atEnd = 1
61                     for j in range(curIndex, len(self.listOfStates)):
62                         if (newState == self.listOfStates[j]):
63                             self.T[i][j] = 1
64                             curIndex = j+1
65                             atEnd = 0
66                             break
67                     if (atEnd):
68                         self.T[i][len(self.listOfStates)] = 1
69                         self.listOfStates.append(newState)
70                         curIndex = len(self.listOfStates)
71                     self.Q.append([])
72
73 def newState(self, heapIndex, amount, state):
74     curState = state[:]
75     curAmount = amount
76     for i in range(heapIndex, len(state)-1):
77         dif = curState[i] - curState[i+1]
78         if (dif >= curAmount):
79             curState[i] -= curAmount
80             curAmount = 0
81             break
82         else:
83             curState[i] -= dif
84             curAmount -= dif
85     if (curAmount > 0):
86         curState[len(curState)-1] -= curAmount
87     return curState

```

Listing 6: The setup function for the trimmed programs.

Since at the time of writing the code I knew of no way to calculate the index of a state when using the trimmed game-tree²³, the creation of the matrix T

²³See Equation 17 for how the calculation is done.

works a bit differently. It works like this:

- Repeat for all possible states:
 - Repeat for all heaps in current state, starting at the last heap, keeping track of the current index of where the possible moves could go, it is reset for each heap to be the index of the current state + 1:
 - * Repeat for all counters in current heap:
 1. Use the function *newState* to find what the resulting state of the new move will be.
 2. Is there a defined state at the current index in the current column in T? if no, add the resulting state at the end of *listOfStates*, and update the spot for the move in T to 1, if yes, go through the states on and after the index, to see if any of the states are the same as the resulting state. if they are, update T to 1 on the relevant spot. If no such state is found, the state is added at the end of *listOfStates*, and T is updated, once again in the same fashion.

The function *newState* could be made simpler by just having it make the move as if it was the untrimmed game-tree, and then sort the new state.

One-dimensional This program only has a small difference to the previous programs, since they both use the trimmed game-tree, as well as the same matrix T. Listing 7 shows the part of the code which is different from the trimmed programs, the first and last line it has in common.

```
70     curIndex = len(self.listOfStates)
71     if(i == (self.sum-1)):
72         self.Q.append(100.0)
73     else:
74         self.Q.append(random.random())
75     timeEnd = time.time()
```

Listing 7: The part of the setup function which is different in the One-dimensional program.

6.2.3 Training

Here the difference is not in the game tree, as the important part there is how the matrices T and Q are made, but rather the algorithm used to train, which is unrelated to how you make the matrices.

An Error in the Code There is an error in the code for the training, which is that for the variable *r*, which is the actual reward for making a move, the programs instead use the expected reward. This is not a big problem, because the actual reward for all moves except the winning one is 0, so the result is

just that some rewards are somewhat exaggerated, but it shouldn't impact the training. Having done some tests, the difference seems negligible in this case. This could easily hinder training in other problems, however, so I shall strive not to make that error again.

Q-learning These programs use this algorithm, as expected, except for one thing, that is the very adjustment of the expected reward. It isn't added to the current value, it is subtracted. This is based upon the logic that giving your opponent the option of a good move is bad, and forcing the opponent to make a bad move is good.

```

95 def train(self, amountEpisodes, epsilon, mu, gamma):
96     for neverUsed in range(amountEpisodes):
97         curState = 0
98         playing = 1
99         while(playing):
100             curAction = self.epsilonGreedyChoice(curState, epsilon)
101             curReward = self.Q[curState][curAction]
102             nextState = curAction
103             nextAction = np.argmax(self.Q[nextState])
104             self.Q[curState][curAction] -= mu * (curReward + gamma*self.Q
105 [nextState][nextAction] - self.Q[curState][curAction])
106             curState = nextState
107             if (nextAction == self.sum-1):
108                 playing = 0
109
110 def epsilonGreedyChoice(self, state, epsilon):
111     if(random.random()<epsilon):
112         tmpList = []
113         for i in range(state, self.sum):
114             if(self.T[state][i] == 1):
115                 tmpList.append(i)
116         return tmpList[random.randrange(0, len(tmpList))]
117     else:
118         return np.argmax(self.Q[state])

```

Listing 8: The code for the Q-learning algorithm. Note that linenumbers are from the untrimmed program.

Sarsa This is the code for the Sarsa algorithm, but like the Q-learning programs it subtracts, instead of adding the adjustment of expected reward.

```

95 def train(self, amountEpisodes, epsilon, mu, gamma):
96     for neverUsed in range(amountEpisodes):
97         curState = 0
98         curAction = self.epsilonGreedyChoice(curState, epsilon)
99         playing = 1
100        while(playing):
101            curReward = self.Q[curState][curAction]
102            nextState = curAction
103            nextAction = self.epsilonGreedyChoice(nextState, epsilon)
104            self.Q[curState][curAction] -= mu * (curReward + gamma*self.Q
105 [nextState][nextAction] - self.Q[curState][curAction])

```

```

105     curState = nextState
106     curAction = nextAction
107     if (curAction == self.sum-1):
108         playing = 0
109
110 def epsilonGreedyChoice(self, state, epsilon):
111     if(random.random()<epsilon):
112         tmpList = []
113         for i in range(state, self.sum):
114             if(self.T[state][i] == 1):
115                 tmpList.append(i)
116         return tmpList[random.randrange(0, len(tmpList))]
117     else:
118         return np.argmax(self.Q[state])

```

Listing 9: The code for the Sarsa algorithm. Note that linenumbers are from the untrimmed program.

One-dimensional This program uses the sarsa algorithm, but the code is slightly different, as the Q is One-dimensional. Thus for each occurance of Q it only uses one index, which is the difference. For the sake of completeness, Listing 10 has the code.

```

144 def train(self, amountEpisodes, epsilon, mu, gamma):
145     for neverUsed in range(amountEpisodes):
146         curState = 0
147         curAction = self.epsilonGreedyChoice(curState, epsilon)
148         playing = 1
149         while(playing):
150             curReward = self.Q[curAction]
151             nextState = curAction
152             nextAction = self.epsilonGreedyChoice(nextState, epsilon)
153             self.Q[curAction] -= mu * (curReward + gamma*self.Q[
154                 nextState] - self.Q[curAction])
155             curState = nextState
156             curAction = nextAction
157             if (curAction == self.sum-1):
158                 playing = 0
159
160 def epsilonGreedyChoice(self, state, epsilon):
161     if(random.random()<epsilon):
162         tmpList = []
163         for i in range(state, self.sum):
164             if(self.T[state][i] == 1):
165                 tmpList.append(i)
166         return tmpList[random.randrange(0, len(tmpList))]
167     else:
168         highVal = -np.inf
169         highValIndex = -1
170         for i in range(state+1, len(self.Q)):
171             if(self.T[state][i] == 1 and np.greater(self.Q[i], highVal)):
172                 highVal = self.Q[i]
173                 highValIndex = i

```

```
176     return highValIndex
```

Listing 10: The code for the Sarsa algorithm, in the One-dimensional program

6.2.4 Making Moves

The programs all have a function called *makeMove*. It takes a state, and returns how an array with 2 values, the first for how many counters should be removed, and the second is from which heap they should be removed from. There is a difference between the programs, once again based upon their use of game-trees, as well as the one-dimensional one being different.

Untrimmed The function in these programs is quite simple, the code in Listing 11 works like this:

1. If the input state has fewer heaps than the start-state the program was trained from, extra heaps of size 0 are added at the end.
2. Find the state we wish to make a move from in the list of states.
3. Find the move with the highest expected reward.
4. Find the index from which the counters were removed, and put the index, as well as the amount of counters removed, into an array which is then returned.

```
76 def makeMove(self, state):
77     curState = state[:]
78     nextState = []
79     move = []
80     if len(state) < len(self.listOfStates[0]):
81         dif = len(self.listOfStates[0]) - len(state)
82         for i in range(dif):
83             curState.append(0)
84     for i in range(len(self.listOfStates)):
85         if curState == self.listOfStates[i]:
86             nextState = self.listOfStates[np.argmax(self.Q[i])]
87     for i in range(len(curState)):
88         if curState[i] != nextState[i]:
89             move.append(curState[i] - nextState[i])
90             move.append(i)
91     return move
```

Listing 11: The *makeMove* function in the untrimmed programs.

Trimmed The function is somewhat different in these programs, as finding the correct move is simple enough, but the input state is sorted to fit the program, so finding the correct index for moving is somewhat more difficult, so it is longer and more complex. In addition, finding the state in the list of states is done by calculation, which is faster, though uses more code. The process is this:

Figure 7: The equation for finding the index of a state, in relation to the start state, for a specific move. n is the size of the heaps in the starting state, x is the amount of counters being removed, and k is the amount of piles right if the pile being removed, which would be $\text{length}(\text{state}) - k'$ where k' is the pile being removed from.

$$x * \left(\frac{1}{n!} * (k+1) * (k+2) * \dots * (k+n) \right) \quad (17)$$

1. Sort the input state, biggest first, then if the input state has fewer heaps than the start-state the program was trained from, extra heaps of size 0 are added at the end.
2. Find the state we want to move from in the list of states, by calculating its index. this is done by using Equation 17.
3. Figure out the how many counters were removed, and from which index.
4. Find the correct index to be output. Since the input state was sorted, the index is not the same as it was, which means an index of the same value needs to be found in the proper input state.

```

103 def makeMove(self, state):
104     curState = state[:]
105     curState.sort(reverse = True)
106     nextState = []
107     move = []
108     if (len(state)<len(self.listOfStates[0])):
109         dif = len(self.listOfStates[0])-len(state)
110         for i in range(dif):
111             curState.append(0)
112     goal = 0
113     progress = self.board[0]
114     for i in range(len(self.board)):
115         heaps = len(self.board)-(i+1)
116         for j in range(progress, -1, -1):
117             if(j == curState[i]):
118                 break
119             else:
120                 if(heaps == 0):
121                     goal += j-curState[i]
122                     break
123             result = 1
124             for x in range(1, heaps + 1):
125                 result *= (j + x)
126             goal += (result)//math.factorial(heaps)
127             progress -= 1
128     nextState = self.listOfStates[np.argmax(self.Q[goal])][:]
129     numToDec = 0
130     for i in range(len(curState)-1, -1, -1):
131         if (curState[i] != nextState[i]):
```

```

132     numToDec += curState[i]-nextState[i]
133     newIndex = curState[i]
134     move.append(numToDec)
135     for i in range(len(state)):
136         if(newIndex == state[i]):
137             newIndex = i
138             break
139     move.append(newIndex)
140 return move

```

Listing 12: The *makeMove* function in the trimmed programs.

One-dimensional For this program the code might look more complex, but it is really about the same as for the trimmed programs. the only real difference lies in that instead of using numpy's *argmax* function, it needs its own version, as there isn't a matrix Q, it is an array instead, without the holes indicating legal movement. Thus the code for figuring out the highest expected value for a move needs to include T, which *argmax* doesn't do.

```

97 def makeMove(self, state):
98     curState = state[:]
99     curState.sort(reverse = True)
100    nextState = []
101    move = []
102    if (len(state)<len(self.listOfStates[0])):
103        dif = len(self.listOfStates[0])-len(state)
104        for i in range(dif):
105            curState.append(0)
106    goal = 0
107    progress = self.board[0]
108    for i in range(len(self.board)):
109        heaps = len(self.board)-(i+1)
110        for j in range(progress, -1, -1):
111            if(j == curState[i]):
112                break
113            else:
114                if(heaps == 0):
115                    goal += j-curState[i]
116                    break
117                result = 1
118                for x in range(1, heaps + 1):
119                    result *= (j + x)
120                    goal += (result)//math.factorial(heaps)
121                    progress -= 1
122    highVal = -np.inf
123    highValIndex = -1
124    for j in range(goal+1, len(self.Q)):
125        if(self.T[goal][j] == 1 and np.greater(self.Q[j], highVal)):
126            highVal = self.Q[j]
127            highValIndex = j
128    nextState = self.listOfStates[highValIndex][:]
129    numToDec = 0
130    for i in range(len(curState)-1, -1, -1):
131        if (curState[i] != nextState[i]):
132            numToDec += curState[i]-nextState[i]
133            newIndex = curState[i]

```

```

134     move.append(numToDec)
135     for i in range(len(state)):
136         if(newIndex == state[i]):
137             newIndex = i
138             break
139     move.append(newIndex)
140     return move

```

Listing 13: The *makeMove* function in the one-dimensional program.

6.2.5 Run.py

This is the program that actually runs the previous programs, and takes some inputs as to which of the programs it will run, and for which games of Nim. The program takes the following input:

1. Lower size of the Nim heaps you want to run.
2. Upper size of the Nim heaps.
3. Lower amount of Nim heaps.
4. Upper amount of Nim heaps.
5. Whether the program *ReinforcementSarsa.py* should be run. 0 for no, 1 for yes.
6. Whether the program *ReinforcementQ-learning.py* should be run.
7. Whether the program *SpaceSaveSarsa.py* should be run.
8. Whether the program *SpaceSaveQ-learning.py* should be run.
9. Whether the program *OneDSarsa.py* should be run.
10. Whether the program *SpaceSaveSarsaSpread.py* should be run.

Handling of input This is the code *Run.py* uses to handle inputs and a little other things.

```

1 import time
2 import sys
3 import ReinforcementSarsa as RS
4 import ReinforcementQlearning as RQ
5 import PerfectPlayer as PP
6 import SpaceSaveSarsa as SS
7 import SpaceSaveQlearning as SQ
8 import OneDSarsa as ODS
9 import SpaceSaveSarsaSpread as SSSS

11 if(len(sys.argv) < 11):
12     print("Error: too few arguments. Got " + str(len(sys.argv)-1) + " "
13     , Expected 10")
13 elif (len(sys.argv) > 11):

```

```

14     print("Error: too many arguments. Got, " + str(len(sys.argv)-1) +
15         ", Expected 10")
16 else:
17     #Parsing the input.
18     heapSizeLower = int(sys.argv[1])
19     heapSizeUpper = int(sys.argv[2])
20     heapAmountLower = int(sys.argv[3])
21     heapAmountUpper = int(sys.argv[4])

22     player = PP.PerfectPlayer()

23     ssstatistics = open("ssstatistics.txt", "w")

24     #writing the header of the statistics file, to show which numbers
25     #correspond to which program.
26     header = ""
27     if(int(sys.argv[5]) == 1):
28         header += "\t\t\t\t\t\tRS\t"
29     if(int(sys.argv[6]) == 1):
30         header += "\t\t\t\t\t\tRQ\t"
31     if(int(sys.argv[7]) == 1):
32         header += "\t\t\t\t\t\tTS\t"
33     if(int(sys.argv[8]) == 1):
34         header += "\t\t\t\t\t\tTQ\t"
35     if(int(sys.argv[9]) == 1):
36         header += "\t\t\t\t\t\tODS\t"
37     if(int(sys.argv[10]) == 1):
38         header += "\t\t\t\t\t\tNew Program"
39     header += "\n"
40     ssstatistics.write(header)
41

```

Listing 14: The code in *Run.py* before the programs are actually run.

Running of the programs This code is repeated 6 times, once for each program.

- Repeat for each state specified in the input:
 - Repeat for each program:
 1. Check the inputs to see if this programs should be run.
 2. Repeat 10 times
 - (a) Initialize and run the setup of the program(which also runs the training), taking the time it uses for this.
 - (b) Play against the perfect player(which uses the deterministic algorithm), for each state the program can win(that is those with a nim-sum unequal to 0) if it starts, that are less than the state the program was run with.
 3. find the average time of all the programs, and write it into the statistics file.

```

42     #The loop that defines the start-states from which the programs
43     #will be run.

```

```

43 for heapSize in range(heapSizeLower, heapSizeUpper+1):
44     for heapAmount in range(heapAmountLower, heapAmountUpper+1):
45         state = [heapSize] * heapAmount
46         print(state)
47         line = ""
48         """This repeats for each program. For each start-state they
49         are run 10 times, taking their time use, after which they are
50         set to play against the perfect player
51         for all starts it can play from, from which it can win when
52         starting. That is, the states with a nim-sum != 0. The averages
53         of the times and winrate is then
54         written into the statistics file."""
55         if(int(sys.argv[5]) == 1):
56             timeAvg = 0
57             percent = 0
58             setupTime = 0
59             trainTime = 0
60             for i in range(10):
61                 start = time.time()
62                 rs = RS.ReinforcementSarsa(state)
63                 times = rs.setup()
64                 end = time.time()
65                 timeAvg += end-start
66                 setupTime += times[0]
67                 trainTime += times[1]
68                 listOfStates = rs.getListOfStates()
69                 wins = 0
70                 losses = 0
71                 for i in range(len(listOfStates)):
72                     if(player.findNimSum(listOfStates[i]) != 0):
73                         playing = 1
74                         board = listOfStates[i][:]
75                         winBoard = [0] * len(listOfStates[i])
76                         while(playing):
77                             move = rs.makeMove(board)
78                             board[move[1]] = board[move[1]]-move[0]
79                             if(board == winBoard):
80                                 wins += 1
81                                 break
82                             move = player.makeMove(board)
83                             board[move[1]] = board[move[1]]-move[0]
84                             if(board == winBoard):
85                                 losses += 1
86                                 break
87                         percent += (wins/(wins+losses)*100)
88                         timeAvg /= 10
89                         percent /= 10
90                         setupTime /= 10
91                         trainTime /= 10
92                         line += "setup: " + str(round(setupTime, 3)) + "\t train:
93                         " + str(round(trainTime, 3)) + "\t total: " + str(round(
94                         timeAvg, 3)) + "\t" + str(int(percent)) + "%\t"
95                         if(int(sys.argv[6]) == 1):

```

Listing 15: This is the code that runs the programs, takes the time they use and creates statistics about it.

6.2.6 *Play.py*

This program lets the user play a game of Nim against one of the reinforcement programs, or get the transcript of a full game against the perfect player written to file. Its inputs are like this:

1. Should the program play agains the perfect player or the user. 0 for the perfect player, any other number for the user.
2. Which reinforcement program should be run? 0 for *reinforcementSarsa.py*, 1 for *reinforcementQ-learning.py*, 2 for *spaceSaveSarsa.py*, 3 for *spaceSaveQ-learning.py*, 4 for *oneDSarsa.py*, and any other number for *SpaceSaveSarsaSpread.py*.
3. The state from which the game starts, in the form: $n_1, n_2, n_3, \dots, n_k$.
4. Should the program make the first move, or should the user/perfect player, 1 for the program making the first move, any other number for the user/perfect player making the first move.

It works quite similar to how games between the programs and the perfect player is done in *Run.py*.

```

1 import time
2 import sys

4 import ReinforcementSarsa as RS
5 import ReinforcementQlearning as RQ
6 import PerfectPlayer as PP
7 import SpaceSaveSarsa as SS
8 import SpaceSaveQlearning as SQ
9 import OneDSarsa as ODS
10 import SpaceSaveSarsaSpread as SSSS

12 if(len(sys.argv) < 5):
13     print("Error: too few arguments. Got " + str(len(sys.argv)-1) + " "
14     , Expected 4")
14 elif (len(sys.argv) > 5):
15     print("Error: too many arguments. Got, " + str(len(sys.argv)-1) + "
16     , Expected 4")
16 else:
17     #Parsing the input
18     startState = list(map(int, sys.argv[3].split(',')))
19     trainingType = int(sys.argv[2])
20     if(trainingType == 0):
21         ml = RS.ReinforcementSarsa(startState)
22     elif(trainingType == 1):
23         ml = RQ.ReinforcementQlearning(startState)
24     elif(trainingType == 2):
25         ml = SS.SpaceSaveSarsa(startState)
26     elif(trainingType == 3):
27         ml = SQ.SpaceSaveQlearning(startState)
28     elif(trainingType == 4):
29         ml = ODS.OneDSarsa(startState)
30     else:
31         ml = SSSS.SpaceSaveSarsaSpread(startState)
```

```

32     ml.setup()
33
34     #The code for the reinforcement program playing against the
35     #deterministic program, and getting statistics from it.
36     if(int(sys.argv[1]) == 0):
37         playstatistics = open("playstatistics.txt", "w")
38         playstatistics.write(str(trainingType) + "\n")
39         player = PP.PerfectPlayer()
40         playing = 1
41         board = startState[:]
42         winBoard = [0] * len(startState)
43         #If the reinforcement algo gets the first move it gets to make
44         #a move before the loop start.
45         if(int(sys.argv[4]) == 1):
46             start = time.time()
47             move = player.makeMove(board)
48             end = time.time()
49             moveTime = end-start
50             board[move[1]] = board[move[1]]-move[0]
51             playstatistics.write("Move by Pre-algo, took " + str(move[0])
52             + " counters from heap " + str(move[1]+1) + ". New state: " +
53             str(board) + " time used: " + str(round(moveTime, 3)) + "s.\n")
54             if(board == winBoard):
55                 playstatistics.write("Pre-algo is the winner.")
56                 playing = 0
57             #Otherwise the game starts here.
58             while(playing):
59                 start = time.time()
60                 move = ml.makeMove(board)
61                 end = time.time()
62                 moveTime = end-start
63                 board[move[1]] = board[move[1]]-move[0]
64                 playstatistics.write("Move by ML, took " + str(move[0]) + "
65                 counters from heap " + str(move[1]+1) + ". New state: " + str(
66                 board) + " time used: " + str(round(moveTime, 3)) + "s.\n")
67                 if(board == winBoard):
68                     playstatistics.write("ML is the winner.")
69                     break
70                 start = time.time()
71                 move = player.makeMove(board)
72                 end = time.time()
73                 moveTime = end-start
74                 board[move[1]] = board[move[1]]-move[0]
75                 playstatistics.write("Move by Pre-algo, took " + str(move[0])
76                 + " counters from heap " + str(move[1]+1) + ". New state: " +
77                 str(board) + " time used: " + str(round(moveTime, 3)) + "s.\n")
78                 if(board == winBoard):
79                     playstatistics.write("Pre-algo is the winner.")
80                     break
81             #If the player wants to play against the Reinforcement algo.
82         else:
83             player = PP.PerfectPlayer()
84             playing = 1
85             board = startState[:]
86             winBoard = [0] * len(startState)
87             #If the player wants to start first, they make a move first.
88             if(int(sys.argv[4]) == 1):

```

```

81     move = [-1, -1]
82     print("Current board is: " + str(board) + "\n")
83     move[1] = int(input("Which heap do you want to remove from? "))
84     while((move[1] < 0 or move[1] >= len(board)) or board[move[1]] == 0):
85         move[1] = int(input("Illegal heap, choose again. "))-1
86     move[0] = int(input("How many counters do you wish to remove? "))
87     while(move[0] < 1 or move[0] > board[move[1]]):
88         move[0] = int(input("Illegal amount, choose again. "))
89     board[move[1]] = board[move[1]]-move[0]
90     if(board == winBoard):
91         print("You are the winner.")
92     playing = 0
93     #Otherwise the game starts here, with the reinforcement program
94     #making the first move.
95     while(playing):
96         move = ml.makeMove(board)
97         board[move[1]] = board[move[1]]-move[0]
98         print("Move by ML, took " + str(move[0]) + " counters from
99         heap " + str(move[1]+1) + ".")
100        if(board == winBoard):
101            print("ML is the winner.")
102            break
103        move = [-1, -1]
104        print("Current board is: " + str(board) + "\n")
105        move[1] = int(input("Which heap do you want to remove from? "))
106        while((move[1] < 0 or move[1] >= len(board) or board[move[1]] == 0)):
107            move[1] = int(input("Illegal heap, choose again. "))-1
108        move[0] = int(input("How many counters do you wish to remove? "))
109        while(move[0] < 1 or move[0] > board[move[1]]):
110            move[0] = int(input("Illegal amount, choose again. "))
111        board[move[1]] = board[move[1]]-move[0]
112        if(board == winBoard):
113            print("You are the winner.")
114            break

```

Listing 16: The code in program *Play.py*.

6.3 Programs for Supervised Learning

There are 3 programs used for the supervised learning, as well as the program *PerfectPlayer.py* which is used to ascertain success.

6.3.1 *MakeData.py*

This program makes the data used by the algorithm, for training. It is quite simple, it uses the perfect player, and for all states reachable from the one it is initialized with it writes a list of files and their correct moves to a file.

```
1 import PerfectPlayer as PP
```

```

3  class MakeData:
4      #This method just makes the list of states accessible from the
5      #start-state, and then runs the writeToFile method.
6      def setup(self, startBoard):
7          sum = startBoard[0]+1
8          board = startBoard
9          for i in range(1, len(board)):
10              sum *= board[i]+1
11          self.listOfStates = [0] * sum
12          self.listOfStates[0] = board[:]
13          for i in range(sum):
14              for x in range(len(self.listOfStates[i])):
15                  for y in range(1, self.listOfStates[i][x] + 1):
16                      tmpValue = self newIndex(x, y, board)
17                      if(self.listOfStates[i+tmpValue] == 0):
18                          tmpList = self.listOfStates[i][:]
19                          tmpList[x] -= y
20                          self.listOfStates[i+tmpValue] = tmpList[:]
21          self.writeToFile(self.listOfStates)
22
23      # This method exists for finding, given a move, how far down the
24      # list of states the new state is.
25      def newIndex(self, heapIndex, amount, board):
26          newMove = 1
27          for i in range(len(board)-1, heapIndex, -1):
28              newMove *= board[i]+1
29          newMove *= amount
30          return newMove
31
32      #This file writes each state in the list to the file, and a
33      #correct move to make when in that state.
34      def writeToFile(self, states):
35          player = PP.PerfectPlayer()
36          f = open("CorrectMoves.txt", "w")
37          for i in range(len(states)):
38              if (player.findNinSum(states[i]) != 0):
39                  tmp = [0]*len(states[0])
40                  result = player.makeMove(states[i])
41                  tmp[result[1]] = result[0]
42                  for j in range(len(states[i])):
43                      f.write(str(states[i][j]) + " ")
44                  f.write("\n")
45                  for j in range(len(states[i])):
46                      f.write(str(tmp[j]) + " ")
47                  f.write("\n")
48
49      #This method returns the list of states, which is used to test
50      #programs, by having them play from many different states.
51      def getStates(self):
52          return self.listOfStates

```

Listing 17: The program *MakeData.py* which creates the data used by the supervised learning algorithm.

6.3.2 The Multilayer Perceptron

This is the code for the MultiLayer Perceptron, in the program *mlp.py*. This code is based upon the code I made for an obligatory task in fall 2018 for the course INF4490. It has been modified to apply better to this subject, yet much of the code is the same, as both are code for Multilayer Perceptrons.

Training The training works like described earlier, except: the program only uses the sigmoid function for the hidden layer, for the output layer it uses the identity function, that is $f(x) = x$. To give an idea of the order functions call things:

1. The program is initialized.
2. The function *earlystopping* is called from the program *RunMlp.py*.
3. *earlystopping* runs the function *train* repeatedly, calling *totalError* for each time to check how well the program has been trained. If the total error stops decreasing for each repeat, the training stops.
4. *forward* is called by *train*, once for each iteration.

```
 1 import numpy as np
 2 import random
 3 from scipy.special import expit

 5 class mlp:
 6     def __init__(self, nhidden, stateLen):
 7         self.beta = 1
 8         self.eta = 0.1
 9         self.momentum = 0.0
10         self.matrix1 = []
11         self.matrix2 = []
12         self.inputAmount = stateLen
13         self.outputAmount = stateLen
14         self.hiddenAmount = nhidden
15         self.outputA = [0.0] * (self.outputAmount)
16         self.outputList = [0.0] * (self.outputAmount)
17         self.hiddenU = [0.0] * (nhidden)
18         self.hiddenList = [0.0] * (nhidden)
19         for i in range(self.inputAmount + 1):
20             self.matrix1.append([])
21             for j in range(nhidden):
22                 self.matrix1[i].append(random.random())
23         for i in range(nhidden + 1):
24             self.matrix2.append([])
25             for j in range(self.outputAmount):
26                 self.matrix2[i].append(random.random())

27     def sigmoidDerived(self, x):
28         return expit(x)*(1-expit(x))

29     def identityFunction(self, x):
30         return x
```

```

33
34     def identityDerived(self, x):
35         return 1
36
37     #this function calculates the sum of squares error of all the
38     #valid inputs and targets, and sums them, and returns that for
39     #the earlystopping to use for figuring out when to stop.
40     def totalError(self, inputs, targets):
41         accumError = 0.0
42         for i in range(len(inputs)):
43             accum = 0.0
44             self.forward(inputs[i])
45             for j in range(self.outputAmount):
46                 accum += (targets[i][j] - self.outputList[j])**2
47             accumError += accum/2
48         return accumError
49
50     #earlystopping, stops when error increases, or decreases too
51     #slowly
52     def earlystopping(self, inputs, targets, valid, validtargets):
53         stopTraining = 0
54         lowestError = self.totalError(valid, validtargets)
55         while not stopTraining:
56             self.train(inputs, targets, 10)
57             newError = self.totalError(valid, validtargets)
58             if newError > lowestError or lowestError-newError < 0.5:
59                 stopTraining = 1
60             else:
61                 lowestError = newError
62
63     #train trains the neural net for the inputted amount of
64     #iterations. One pass through the dataset is an iteration.
65     def train(self, inputs, targets, iterations=100):
66         for notUsed in range(iterations):
67             for ite in range(len(inputs)):
68                 self.forward(inputs[ite])
69                 outDelta = []
70                 for i in range(self.outputAmount):
71                     outDelta.append((self.outputList[i] - targets[ite][i]) *
72                         self.identityDerived(self.outputA[i]))
73                 hiddenDelta = []
74                 for i in range(self.hiddenAmount):
75                     accumulator = 0.0
76                     for j in range(self.outputAmount):
77                         accumulator += outDelta[j] * self.matrix2[i+1][j]
78                     hiddenDelta.append(self.sigmoidDerived(self.hiddenU[i]) *
79                         accumulator)
80                     for j in range(self.outputAmount):
81                         self.matrix2[0][j] -= self.eta * outDelta[j] * 1
82                         for i in range(self.hiddenAmount):
83                             self.matrix2[i+1][j] -= self.eta * outDelta[j] * self.
84                             hiddenList[i]
85                             for j in range(self.hiddenAmount):
86                                 self.matrix1[0][j] -= self.eta * hiddenDelta[j] * 1
87                                 for i in range(self.inputAmount):
88                                     self.matrix1[i+1][j] -= self.eta * hiddenDelta[j] *
89                                     inputs[ite][i]

```

```

83     #forward passes the inputs through the neural net to get the
84     #outputs.
85     def forward(self, inputs):
86         for i in range(len(self.hiddenList)):
87             self.hiddenU[i] = 1 * self.matrix1[0][i]
88             for j in range(self.inputAmount):
89                 self.hiddenU[i] += inputs[j] * self.matrix1[j+1][i]
90             self.hiddenList[i] = expit(self.hiddenU[i])
91             for i in range(len(self.outputList)):
92                 self.outputA[i] = 1 * self.matrix2[0][i]
93                 for j in range(len(self.hiddenList)):
94                     self.outputA[i] += self.hiddenList[j]*self.matrix2[j+1][i]
95             self.outputList[i] = self.identityFunction(self.outputA[i])

```

Listing 18: This is the part of the program *mlp.py* in which the training happens.

Making Moves The program makes moves by running the state through the neural net, and choosing the largest output. It also makes sure the move is legal, so it doesn't add counters, or remove more counters than there is in a heap. The consequences of this is that it might often make small moves when it really hasn't been trained well enough.

```

99     def makeMove(self, state):
100         self.forward(state)
101         if(state[0] < -1):
102             1/0
103         move = [0,0]
104         move[1] = np.argmax(self.outputList)
105         move[0] = int(self.outputList[move[1]])
106         if(move[0] < 1):
107             move[0] = 1
108         elif(move[0] > state[move[1]]):
109             move[0] = state[move[1]]
110         if(state[move[1]] == 0):
111             for i in range(len(state)):
112                 if(state[i] > 0):
113                     move[0] = 1
114                     move[1] = i
115         return move

```

Listing 19: The part of *mlp.py* that is used to get a move from the program.

6.3.3 *RunMlp.py*

This program is similar to *Run.py* except it runs *mlp.py* instead of any of the Reinforcement programs. Its inputs are fewer, as there is only 1 program to run, and they are:

1. Lower size of the Nim heaps you want to run.
2. Upper size of the Nim heaps.
3. Lower amount of Nim heaps.

4. Upper amount of Nim heaps.

```

1 import MakeData as MD
2 import PerfectPlayer as PP
3 import numpy as np
4 import mlp
5 import sys
6 import time

8 if(len(sys.argv) < 5):
9     print("Error: too few arguments. Got " + str(len(sys.argv)-1) + "
10      , Expected 4")
11 elif (len(sys.argv) > 5):
12     print("Error: too many arguments. Got, " + str(len(sys.argv)-1) +
13      ", Expected 4")
12 else:
13     #Parsing the input.
14     heapSizeLower = int(sys.argv[1])
15     heapSizeUpper = int(sys.argv[2])
16     heapAmountLower = int(sys.argv[3])
17     heapAmountUpper = int(sys.argv[4])
18     f = open("mlpstatistics.txt", "w")
19     filename = "CorrectMoves.txt"
20     player = PP.PerfectPlayer()
21     #The loop that defines the start-states from which the programs
22     #will be run.
22     for heapSize in range(heapSizeLower, heapSizeUpper+1):
23         for heapAmount in range(heapAmountLower, heapAmountUpper+1):
24             """The Supervised learning program is run 10 times for each
25             start-state, and plays against the perfect player each time,
26             to get a winrate for statistics. The time used is also
27             written to file as statistics. The data used is the same for a
28             given state,
29             but it is sorted differently for each time."""
30             state = [heapSize] * heapAmount
31             print(state)
32             md = MD.MakeData()
33             md.setup(state)
34             timeAvg = 0
35             percent = 0
36             for i in range(10):
37                 #sorting and separating the data into the needed categories
38
39                 data = np.loadtxt(filename)
40
41                 stateLen = len(data[0])
42
43                 boards = data[:,0]
44                 moves = data[:,1]
45
46                 order = list(range(np.shape(boards)[0]))
47                 np.random.shuffle(order)
48                 boards = boards[order,:]
49                 moves = moves[order,:]
50
51                 train = boards[::2]
52                 trainTargets = moves[::2]

```

```

50     valid = boards[1::4]
51     validTargets = moves[1::4]

53     test = boards[3::4]
54     testTargets = moves[3::4]

56     hidden = len(state)+10
57
58     start = time.time()
59     net = mlp.mlp(hidden, stateLen)
60     net.earlystopping(train, trainTargets, valid, validTargets)
61     end = time.time()
62     timeAvg += end-start
63     listOfStates = md.getStates()
64     wins = 0
65     losses = 0
66     for i in range(len(listOfStates)):
67         if(player.findNimSum(listOfStates[i]) != 0):
68             playing = 1
69             board = listOfStates[i][:]
70             winBoard = [0] * len(listOfStates[i])
71             while(playing):
72                 move = net.makeMove(board)
73                 board[move[1]] = board[move[1]]-move[0]
74                 if(board == winBoard):
75                     wins += 1
76                     break
77                 move = player.makeMove(board)
78                 board[move[1]] = board[move[1]]-move[0]
79                 if(board == winBoard):
80                     losses += 1
81                     break
82             percent += (wins/(wins+losses)*100)
83             timeAvg /= 10
84             percent /= 10
85             line = "Time used:" + str(round(timeAvg, 3)) + "s Winrate/"
86             sucessrate: " + str(int(percent)) + "% "
87             line += str(state) + "\n"
f.write(line)

```

Listing 20: The code for the program *RunMlp.py*

It shuffles the training data for each run so that it is different each time. Only the total time used is taken here, unlike for the reinforcement programs, as setup was a much larger concern in those.

6.3.4 Program With Stochastic Termination

This is the program that runs the programs with stochastic termination, and records the time used, as well as the amount of episodes it ran. The program takes the following parameters:

1. Lower size of the Nim heaps you want to run.

2. Upper size of the Nim heaps.
3. Lower amount of Nim heaps.
4. Upper amount of Nim heaps.
5. The reinforcement program it should run:

- 0 for *ReinforcementSarsa.py*.
- 1 for *ReinforcementQ-learning.py*.
- 2 for *SpaceSaveSarsa.py*.
- 3 for *SpaceSaveQ-learning.py*.
- 4 for *OneDSarsa.py*.
- Any other number for *SpaceSaveSarsaSpread*.

The program reuses much code from *Run.py*, so they are quite similar. The statistics also show the program being run, with the number input for it, as well as the desired percentage.

```

1 import time
2 import sys
3 import ReinforcementSarsa as RS
4 import ReinforcementQlearning as RQ
5 import PerfectPlayer as PP
6 import SpaceSaveSarsa as SS
7 import SpaceSaveQlearning as SQ
8 import OneDSarsa as ODS
9 import SpaceSaveSarsaSpread as SSS

11 if(len(sys.argv) < 6):
12     print("Error: too few arguments. Got " + str(len(sys.argv)-1) + " "
13         , Expected 5")
13 elif (len(sys.argv) > 6):
14     print("Error: too many arguments. Got, " + str(len(sys.argv)-1) +
15         , Expected 5")
15 else:
16     #Parsing the input.
17     heapSizeLower = int(sys.argv[1])
18     heapSizeUpper = int(sys.argv[2])
19     heapAmountLower = int(sys.argv[3])
20     heapAmountUpper = int(sys.argv[4])
21     trainingType = int(sys.argv[5])
22     player = PP.PerfectPlayer()
23     ssstatistics = open("episodesstatistics.txt", "w")
24     #The desired percentage winrate, which terminates the program.
25     goal = 90
26     ssstatistics.write("Type: " + str(trainingType) + " ,percentage "
27         , goal: " + str(goal) + "%\n")
27     for heapSize in range(heapSizeLower, heapSizeUpper+1):
28         for heapAmount in range(heapAmountLower, heapAmountUpper+1):
29             state = [heapSize] * heapAmount
30             print(state)
31             #Settings for the programs to train.

```

```

32     amountEpisodes = 10
33     epsilon = 0.1
34     mu = 0.5
35     gamma = 0.5
36     totalTime = 0
37     start = time.time()
38     #The selection of the correct program to run.
39     if(trainingType == 0):
40         ml = RS.ReinforcementSarsa(state)
41     elif(trainingType == 1):
42         ml = RQ.ReinforcementQlearning(state)
43     elif(trainingType == 2):
44         ml = SS.SpaceSaveSarsa(state)
45     elif(trainingType == 3):
46         ml = SQ.SpaceSaveQlearning(state)
47     elif(trainingType == 4):
48         ml = ODS.OneDSarsa(state)
49     else:
50         ml = SSS.SpaceSaveSarsaSpread(state)
51     ml.setup(1)
52     end = time.time()
53     listOfStates = ml.getListOfStates()
54     totalTime += end - start
55     stillTraining = 1
56     totalEpisodes = 0
57     #This loop runs until the percentage winrate is equal or
58     #higher than the desired goal.
59     while(stillTraining):
60         start = time.time()
61         ml.train(amountEpisodes, epsilon, mu, gamma)
62         end = time.time()
63         totalTime += end-start
64         wins = 0
65         losses = 0
66         """Unlike in the program Run.py the statistics aren't
67         repeated ten times. Quite understandably, this would not work
68         here,
69         because the desire is one properly trained program, not ten
70         alright ones."""
71         for i in range(len(listOfStates)):
72             if(player.findNimSum(listOfStates[i]) != 0):
73                 playing = 1
74                 board = listOfStates[i][:]
75                 winBoard = [0] * len(listOfStates[i])
76                 while(playing):
77                     move = ml.makeMove(board)
78                     board[move[1]] = board[move[1]]-move[0]
79                     if(board == winBoard):
80                         wins += 1
81                         break
82                     move = player.makeMove(board)
83                     board[move[1]] = board[move[1]]-move[0]
84                     if(board == winBoard):
85                         losses += 1
86                         break
87             totalEpisodes += amountEpisodes
88             percent = (wins//(wins+losses)*100)

```

```

85     if(percent >= goal):
86         stillTraining = 0
87     ssstatistics.write("Time used: " + str(round(totalTime, 3)) +
88     "s, Episodes used: " + str(totalEpisodes) + ", " + str(state)
89     + "\n")

```

Listing 21: The code for the program *Episodes.py*, that tries to run the reinforcement programs with stochastic termination

7 Comparing the Algorithms with Time Complexity

Usually comparing worst time-complexities is used to compare how good an algorithm is. This does not work properly with stochastic algorithms, because some of them do not have a predictable point of termination, where a deterministic algorithm does. This does not always apply, however. An algorithm that simulates a single die, say one with 6 sides, should always terminate at the same time, but it is still stochastic as the outcome will differ each time it is run. In the case of the algorithm that uses reinforcement learning, it will always run the same amount of episodes, so how long time will use has a clear upper limit but for the same output, it does not always use the same amount of time. What is important is whether an algorithm can go on forever, because if it can, there are never a worst case time use, the worst case would be infinite.

7.1 The Deterministic Algorithm

7.1.1 Time-complexity of the Nim-sum Operation

First one needs to know the time complexity of a Nim-sum operation. In many programming languages the Xor bitwise-operation will do the same, as previously mentioned. That operation uses constant time, that is $O(1)$, but that is because it is not constrained the actual size of the number. If one were to use a Turing machine instead, it would have a complexity of $O(n)$ where k is the amount of digits in the larger of the input numbers. This is simply because there has to be some comparison between the 2 numbers for each digit, with the numbers in binary, of course.

1. Repeat until one of the numbers run out of digits, starting on the smallest digit.
 - (a) Compare the current digits of the two numbers.
 - (b) If they are equal set the digit in the sum to be 0. If not set the digit in the sum to be 1.
 - (c) Move on to the next digit.
2. Now that one of the numbers has run out of digits, repeat the following until the largest number is out of digits.

- (a) Set the digit of the sum to be the current digit of the larger number.
- (b) Move on to the next digit.

This list shows a suggested algorithm to find the nim-sum, or the exclusive or comparison of 2 numbers. This repeats the three steps as many times as the lesser number has digits, and has to set the rest of the larger number's digits too, so it has the time complexity of $O(n)$ where n is the amount of digits in the larger number.

7.1.2 The Rest of the Algorithm

So for the algorithm it is required to find the nim-sum of the size of sizes of all the heaps. Since the operation only accepts 2 inputs, it will have to be repeated for $k - 1$ times, where k is the amount of heaps. Furthermore the algorithm uses the operation again to find the size of the heap which should be reached. not all heaps have a viable move, so if it is unlucky, it will have to find the nim-sum for all heaps. Thus the operation will, in the worst case, be used $(k - 1) + k$ times. The best case would be $(k - 1) + 1$ if the first heap can be reduced to get to a winning position. So then, if k is the amount of heaps, and n is the size of the biggest number, the time complexity in total is $O(k * n)$. To reiterate, this is how the algorithm works, abstractly.

1. Repeat the following until the sum includes the size of the last heap, with the current sum being defined as the size of the first heap at the first iteration.
 - Find the nim-sum of the current sum and the size of the next heap.
The result is the new current sum.
2. Starting with the first heap, do the following.
 - (a) Find the nim-sum of the total sum and the size of a heap.
 - (b) If it is less than the size of the heap, a correct move is found, which is to remove counters from the heap so that it reaches the size equal to the sum found. If the sum is larger, then repeat with the next heap.

7.1.3 Practically, Then?

In practice, the xor operation actually has a complexity of $O(1)$ on a computer, like the sum operation $+$, in which case the full algorithm is $O(k)$ as opposed to $O(k * n)$. This is a clear improvement, however it does not make much of a difference. The other algorithms use the sum operation many times, and that gets the same improvement.

7.2 The Reinforcement Learning Algorithm

This algorithm has several parts, and these will be considered separately. There are three parts. First the setup of the matrices, for legal moves and expected rewards. Second is the training itself, where the algorithm plays against itself an amount of times, to learn which moves are the best moves. Third there is the function that tells the user which move it gets the highest reward from, when given a certain game-state.

7.2.1 Setup of the Matrices

The time this takes is governed by the size of the matrix, which is the amount of possible states reachable in a game from the start-state. Let us call this m . It is clear that m increases at least exponentially when compared to the size of the input, because whenever a counter is added, m increases by at least 1, while the size of the input does not have to increase at all, since a digit might just be replaced by a bigger digit.

Size of game trees The size of a side of the matrices are all states reachable from the start states, which is just a game tree. The equation for finding the size of a certain games gametree, if all the heaps are of the same size is simply the following, or trimmed, where n is the size of the heaps, and k is the amount of heaps:

$$\frac{(n+k)!}{n! * k!} \quad (18)$$

To figure out the time complexity of this part of the algorithm, it is needed to confirm that the size of the game tree does not increase exponentially in relation to the amount of counters in the game. If it does the size of the matrices will be something like 2^{2^N} where N is the size of the input, and that won't do at all. The growth of the tree is as follows when adding new heaps, where x is the amount of heaps, and n is the arbitrary amount in each heap:

$$\frac{1}{n!} * (x+1) * (x+2) * \dots * (x+n) \quad (19)$$

The growth of the tree when adding a counter to each pile is like this:

$$\frac{1}{n!} * (x+1) * (x+2) * \dots * (x+n) \quad (20)$$

Which is the same, because the equation for finding the size of the game tree is symmetrical, after all. Do note however that this shows the growth when adding either a heap with the same numbers of counters as the other heaps, or adding 1 counter to all heaps, which is usually more than 1 counter, so adding 1 counter will usually cause a growth much less than these show. Nevertheless

this growth is not exponential, and so the size of the matrices compared to the size of the input is $O(2^N)$.

7.2.2 Training the Matrix

The training does not use the same amount of time, each time, however gauging what the worst case time use is, isn't very hard. Every episode happens a finite amount of times, and in each episode the algorithm completes a full game against itself. The longest possible game of Nim from any state is where both players remove one counter each, so the answer is the sum of the size of the heaps. the sum is of course exponentially bigger than the size of the input, since it is basically a comparison between a number and its size. So this part also has a time complexity of $O(2^N)$. It should be mentioned that the training will fail to find relevant data in most cases if it does this, however this is the closest comparison to be made.

7.2.3 Giving an Answer

Finding the move that gives the highest reward, will again be quite slow, compared to the size of the input, because when you find correct index, there needs to be a search through all the expected rewards to find the highest. That might be nearly all possible positions, which means that it will have the time complexity $O(2^N)$ and there is no way around it. Since the sides of the matrices increase exponentially in comparison to the size of the input, we are out of luck. It does not really matter how the correct index for the input state is found, as the worst solution, just searching, will only at most double the time used.

8 Comparing the Algorithms with Actual Time-use

Perhaps their actual average time use is a better measure, so some statistics is in order. For all these values the program was run 10 times, and the average of those values is what is shown, to be sure no one value is strange because of randomness. Any strangeness should be a feature of the algorithm, not just of plain luck²⁴.

²⁴ Perhaps lack of luck, being unlucky is more accurate here, as nice and accurate values in statistics would probably be better, though less interesting.

8.1 Reinforcement Learning

Here are the statistics about the programs using reinforcement-learning algorithms.

8.1.1 The naming of the programs in the graph

For smaller tables and readability the programs will be given shortform names, which are found in Table 6

Table 6: The shortform names for the different programs used in the statistics.

Name	Program
RS	Program using the untrimmed game-tree, and uses the Sarsa algorithm.
RQ	Program using the untrimmed game-tree, and uses the Q-learning algorithm.
TS	Program using the trimmed game-tree, and uses the Sarsa algorithm
TQ	Program using the trimmed game-tree, and uses the Q-learning algorithm
ODS	Program using the trimmed game-tree, uses the Sarsa algorithm, and uses the state-value function $V(s)$

8.1.2 The Settings Used by the Programs Running, and Other Relevant Information

When training, the training variables μ and γ are both set to 0.5, and ϵ is set to 0.1, except for ODS, where it is set to 0. The program runs the training through 1000 episodes. The way the winrate is found(The measure of success) is by having the program play against the perfect player(a program using the deterministic algorithm for playing Nim) from all states it knows, and it can win from. That means that for all the states in the programs list of states where the states nim-sum is not 0, it will play a game against the perfect player. The wins and losses are registered, so a percentage can be found. The program is trained, tested, and timed 10 times, and the average of the results are found.

Possible Weakness with the Measure of Success Measuring winrate this way may have a disadvantage. To win a game against the perfect player, every move needs to be a move to a state of the Nim-sum 0. If only one move is wrong, the game is lost. Since the testing is done from many states, many different games are played, but many of the same states are reached. The point is that to win any game it reaches, it must do the winning move, emptying the last heap, while many moves in larger states are just required for a small subset of games. Thus the ability to make the last move is much more important and perhaps over-valued using the current method of measuring success. On the other hand however, the move of removing from the last heap is quicker to learn, since the best move is already defined, and for any move, the closer

it is to a win, the earlier in the training it will be trained, and given a more appropriate expected reward. So the bias for those final moves is justified given their advantage in training. Indeed, that advantage might make a method where you only check if the correct moves are made from winning positions be biased the towards the final moves as well.

8.1.3 General Statistics

Table 7 contains the run-times and winrates of all the pre-mentioned programs for the sake of comparison. The supremacy of the trimmed programs should be immediatly obvious, as well as ODS doing better at harder states, but actually worse at some easier states, it is perhaps too generalized. The difference between Q-learning and Sarsa seems to be quite small, as they seem to do about as well, with some fluctuations. The biggest difference seem to be that Q-learning has a somewhat higher run-time, likely because for each move it needs to find the best move twice(once for the current action, and once for *argmax* in the assignment), while Sarsa only does so once.

Table 7: The general statistics of the different algorithms.

RS		RQ		TS		TQ		ODS		
Train	Wins	Train	Wins	Train	Wins	Train	Wins	Train	Wins	State
0.014s	100%	0.014s	100%	0.013s	100%	0.013s	100%	0.006s	100%	1 1
0.022s	100%	0.039s	100%	0.021s	100%	0.026s	100%	0.010s	100%	1 1 1
0.030s	100%	0.045s	100%	0.028s	100%	0.039s	100%	0.013s	100%	1 1 1 1
0.041s	100%	0.063s	100%	0.036s	100%	0.055s	100%	0.017s	100%	1 1 1 1 1
0.016s	100%	0.017s	100%	0.015s	100%	0.015s	100%	0.010s	100%	2 2
0.026s	97%	0.036s	98%	0.023s	98%	0.029s	97%	0.016s	100%	2 2 2
0.057s	99%	0.091s	99%	0.036s	100%	0.052s	100%	0.026s	97%	2 2 2 2
0.130s	72%	0.201s	67%	0.046s	92%	0.070s	89%	0.036s	100%	2 2 2 2 2
0.019s	100%	0.022s	100%	0.017s	100%	0.017s	100%	0.014s	100%	3 3
0.036s	87%	0.050s	85%	0.027s	90%	0.034s	93%	0.025s	100%	3 3 3
0.149s	96%	0.251s	96%	0.050s	100%	0.075s	100%	0.051s	89%	3 3 3 3
0.820s	40%	1.099s	42%	0.065s	84%	0.100s	84%	0.075s	90%	3 3 3 3 3
0.023s	100%	0.027s	100%	0.019s	100%	0.020s	100%	0.018s	100%	4 4
0.056s	67%	0.072s	57%	0.032s	84%	0.040s	66%	0.035s	98%	4 4 4
0.397s	43%	0.560s	44%	0.076s	99%	0.108s	99%	0.089s	96%	4 4 4 4
5.888s	14%	6.330s	14%	0.113s	74%	0.165s	79%	0.147s	68%	4 4 4 4 4
0.042s	100%	0.048s	100%	0.024s	100%	0.023s	100%	0.024s	94%	5 5
0.105s	62%	0.121s	46%	0.039s	66%	0.048s	63%	0.050s	94%	5 5 5
1.197s	23%	1.499s	24%	0.118s	99%	0.171s	99%	0.155s	84%	5 5 5 5
32.590s	6%	33.137s	6%	0.225s	65%	0.317s	71%	0.301s	84%	5 5 5 5 5

8.1.4 The Untrimmed Sarsa

Here are statistics on how the program running the sarsa algorithm with the untrimmed gametree. Something interesting is how the time used by the setup increases faster than that of the training.

Table 8: The statistics for applying sarsa to the untrimmed game tree.

Setup	Training	Total	Winrate	Starting State
0.000s	0.014s	0.014s	100%	1 1
0.000s	0.022s	0.022s	100%	1 1 1
0.000s	0.030s	0.030s	100%	1 1 1 1
0.001s	0.040s	0.041s	100%	1 1 1 1 1
0.000s	0.015s	0.016s	100%	2 2
0.001s	0.026s	0.026s	97%	2 2 2
0.003s	0.054s	0.057s	99%	2 2 2 2
0.026s	0.104s	0.130s	72%	2 2 2 2 2
0.000s	0.019s	0.019s	100%	3 3
0.002s	0.034s	0.036s	87%	3 3 3
0.028s	0.120s	0.149s	96%	3 3 3 3
0.455s	0.351s	0.820s	40 %	3 3 3 3 3
0.000s	0.021s	0.023s	100%	4 4
0.007s	0.049s	0.056s	67%	4 4 4
0.167s	0.226s	0.397s	43%	4 4 4 4
4.529s	1.200s	5.888s	14%	4 4 4 4 4
0.001s	0.025s	0.042s	100%	5 5
0.021s	0.084s	0.105s	62%	5 5 5
0.738s	0.438s	1.197s	23%	5 5 5 5
27.906s	3.008s	32.590s	6%	5 5 5 5 5

8.1.5 The Untrimmed Q-learning

Here are some statistics on how the program running the sarsa algorithm with the untrimmed gametree. The training time-use seems significantly higher than that for untrimmed Sarsa, while the setup time is about the same, supporting the theory of why Q-learning is slower.

Table 9: The statistics for applying Q-learning to the untrimmed game tree.

Setup	Training	Total	Winrate	Starting State
0.000s	0.014s	0.014s	100%	1 1
0.000s	0.039s	0.039s	100%	1 1 1
0.000s	0.045s	0.045s	100%	1 1 1 1
0.001s	0.062s	0.063s	100%	1 1 1 1 1
0.000s	0.017s	0.017s	100%	2 2
0.001s	0.036s	0.036s	98%	2 2 2
0.004s	0.087s	0.091s	99%	2 2 2 2
0.028s	0.173s	0.201s	67%	2 2 2 2 2
0.000s	0.021s	0.022s	100%	3 3
0.002s	0.048s	0.050s	85%	3 3 3
0.031s	0.218s	0.251s	96%	3 3 3 3
0.473s	0.611s	1.099s	42%	3 3 3 3 3
0.000s	0.025s	0.027s	100%	4 4
0.007s	0.065s	0.072s	57%	4 4 4
0.175s	0.378s	0.560s	44%	4 4 4 4
4.384s	1.805s	6.330s	14%	4 4 4 4 4
0.001s	0.032s	0.048s	100%	5 5
0.022s	0.098s	0.121s	46%	5 5 5
0.763s	0.714s	1.499s	24%	5 5 5 5
26.996s	4.696s	33.137s	6%	5 5 5 5 5

8.1.6 The Trimmed Sarsa

These are some statistics for the program using the trimmed gametree and that applies the sarsa algorithm on it. This program is clearly more useful than the untrimmed programs, being able to play larger heaps both better and faster. The time use still increases faster for setup than training though, interestingly enough. This is not strange however, since the code is quite similar, and the main difference lies in the graph used.

Table 10: The statistics for applying sarsa to the trimmed game tree.

Setup	Training	Total	Winrate	Starting State
0.000s	0.023s	0.023s	98%	2 2 2
0.000s	0.036s	0.036s	100%	2 2 2 2
0.000s	0.046s	0.046s	92%	2 2 2 2 2
0.001s	0.062s	0.063s	100%	2 2 2 2 2 2
0.001s	0.026s	0.027s	90%	3 3 3
0.001s	0.049s	0.05s	100%	3 3 3 3
0.002s	0.063s	0.065s	84%	3 3 3 3 3
0.005s	0.100s	0.106s	99%	3 3 3 3 3 3
0.001s	0.031s	0.032s	84%	4 4 4
0.004s	0.072s	0.076s	99%	4 4 4 4
0.011s	0.102s	0.113s	74%	4 4 4 4 4
0.029s	0.198s	0.228s	86%	4 4 4 4 4 4
0.003s	0.037s	0.039s	66%	5 5 5
0.011s	0.107s	0.118s	99%	5 5 5 5
0.037s	0.187s	0.225s	65%	5 5 5 5 5
0.124s	0.344s	0.472s	43%	5 5 5 5 5 5
0.006s	0.047s	0.052s	70%	6 6 6
0.027s	0.148s	0.175s	76%	6 6 6 6
0.126s	0.302s	0.433s	38%	6 6 6 6 6
0.458s	0.581s	1.051s	22%	6 6 6 6 6 6
0.009s	0.065s	0.074s	72%	7 7 7
0.067s	0.206s	0.275s	49%	7 7 7 7
0.361s	0.484s	0.856s	23%	7 7 7 7 7
1.518s	0.941s	2.499s	12%	7 7 7 7 7 7
0.016s	0.080s	0.096s	56%	8 8 8
0.128s	0.236s	0.367s	31%	8 8 8 8
0.862s	0.608s	1.492s	13%	8 8 8 8 8
4.284s	1.655s	6.056s	6%	8 8 8 8 8 8

8.1.7 The Trimmed Q-learning

These are some statistics for the program using the trimmed gametree and that applies the Q-learning algorithm on it. Trimmed Q-learning is also somewhat slower in its training, just like untrimmed Q-learning. They share the same code for training, so this was expected.

Table 11: The statistics for applying Q-learning to the trimmed game tree.

Setup	Training	Total	Winrate	Starting State
0.000s	0.029s	0.029s	97%	2 2 2
0.000s	0.051s	0.052s	100%	2 2 2 2
0.001s	0.069s	0.070s	89%	2 2 2 2 2
0.001s	0.097s	0.098s	100%	2 2 2 2 2 2
0.001s	0.034s	0.034s	93%	3 3 3
0.001s	0.074s	0.075s	100%	3 3 3 3
0.003s	0.097s	0.100s	84%	3 3 3 3 3
0.006s	0.165s	0.170s	99%	3 3 3 3 3 3
0.001s	0.038s	0.040s	66%	4 4 4
0.004s	0.105s	0.108s	99%	4 4 4 4
0.010s	0.155s	0.165s	79%	4 4 4 4 4
0.027s	0.309s	0.336s	99%	4 4 4 4 4 4
0.002s	0.045s	0.048s	63%	5 5 5
0.010s	0.161s	0.171s	99%	5 5 5 5
0.037s	0.279s	0.317s	71%	5 5 5 5 5
0.117s	0.527s	0.648s	52%	5 5 5 5 5 5
0.005s	0.056s	0.060s	53%	6 6 6
0.027s	0.234s	0.261s	88%	6 6 6 6
0.129s	0.494s	0.627s	47%	6 6 6 6 6
0.477s	0.971s	1.461s	24%	6 6 6 6 6 6
0.009s	0.071s	0.080s	51%	7 7 7
0.061s	0.309s	0.372s	56%	7 7 7 7 7
0.336s	0.698s	1.043s	26%	7 7 7 7 7 7
1.527s	1.632s	3.200s	13%	7 7 7 7 7 7 7
0.017s	0.089s	0.106s	42%	8 8 8
0.139s	0.416s	0.558s	37%	8 8 8 8
0.868s	1.076s	1.967s	15%	8 8 8 8 8
4.503s	2.843s	7.477s	7%	8 8 8 8 8 8

8.1.8 The One Dimensional Program

These are some statistics for the program using the One-dimensional array, or the State-value function, to contain the expected rewards. This program seem to generally do better than the other trimmed programs, but it seems not be able to reach 100% as easily as the other programs for some reason.

Table 12: The statistics for using the program using the one-dimensional array, and sarsa.

Setup	Training	Total	Winrate	Starting State
0.000s	0.016s	0.016s	100%	2 2 2
0.000s	0.026s	0.026s	97%	2 2 2 2
0.000s	0.035s	0.036s	100%	2 2 2 2 2
0.001s	0.051s	0.052s	97%	2 2 2 2 2 2
0.000s	0.025s	0.025s	100%	3 3 3
0.001s	0.050s	0.051s	89%	3 3 3 3
0.001s	0.074s	0.075s	90%	3 3 3 3 3
0.002s	0.124s	0.126s	96%	3 3 3 3 3 3
0.000s	0.034s	0.035s	98%	4 4 4
0.001s	0.087s	0.089s	96%	4 4 4 4
0.004s	0.143s	0.147s	68%	4 4 4 4 4
0.009s	0.304s	0.313s	98%	4 4 4 4 4 4
0.001s	0.049s	0.050s	94%	5 5 5
0.004s	0.151s	0.155s	84%	5 5 5 5
0.011s	0.291s	0.301s	84%	5 5 5 5 5
0.029s	0.722s	0.751s	81%	5 5 5 5 5 5
0.002s	0.068s	0.070s	90%	6 6 6
0.008s	0.250s	0.258s	88%	6 6 6 6
0.028s	0.610s	0.638s	64%	6 6 6 6 6
0.087s	1.385s	1.474s	46%	6 6 6 6 6 6
0.004s	0.091s	0.094s	89%	7 7 7
0.016s	0.413s	0.429s	89%	7 7 7 7
0.070s	1.088s	1.160s	58%	7 7 7 7 7
0.257s	2.529s	2.794s	32%	7 7 7 7 7 7
0.006s	0.128s	0.134s	58%	8 8 8
0.031s	0.635s	0.667s	64%	8 8 8 8
0.148s	1.656s	1.808s	44%	8 8 8 8 8
0.655s	4.147s	4.826s	23%	8 8 8 8 8 8

8.1.9 Phenomena: Total Time Use Is More than the Sum of the Training and Setup Time.

A phenomena found in the statistics for all of the different programs is the fact that at times the total time used is more than the training and setup times together. This seems strange, because it seems that they should be the same, and they mostly are. The Table 13 shows some examples of this occurring, from the two programs using the untrimmed game tree. It does happen with the other algorithms too, however.

Table 13: The lines which show a total time use that is more than the sum of the setup and train times.

Setup	Training	Total	Winrate	Starting State
RS				
0.455s	0.351s	0.820s	40 %	3 3 3 3 3
0.000s	0.021s	0.023s	100%	4 4
0.167s	0.226s	0.397s	43%	4 4 4 4
4.529s	1.200s	5.888s	14%	4 4 4 4 4
0.001s	0.025s	0.042s	100%	5 5
0.738s	0.438s	1.197	23%	5 5 5 5
27.906s	3.008s	32.590s	6%	5 5 5 5 5
RQ				
0.473s	0.611s	1.099s	42%	3 3 3 3 3
4.384s	1.805s	6.330s	14%	4 4 4 4 4
0.001s	0.032s	0.048s	100%	5 5
0.763s	0.714s	1.499s	24%	5 5 5 5
26.996s	4.696s	33.137s	6%	5 5 5 5 5

Hypothesis I believe this happens because the times train and setup do not measure in total the part of code that total does, it runs a little more code than the other 2 in total. See Listing 22 where the the total time is taken in the file *Run.py*. The setup and training time is in Listing 23 and the initialization, run only within the total time is in Listing 24. Based upon all that it seems reasonable that the computer might just have used a little extra time to initialize when running, skewing the numbers a bit.

```
57 start = time.time()
58 rs = RS.ReinforcementSarsa(state)
59 times = rs.setup()
60 end = time.time()
```

Listing 22: Total time is taken here. From file *Run.py*

```
23 def setup(self):
24     timeStart = time.time()
25     for i in range(self.sum):
```

```

26
46 timeEnd = time.time()
47 setupTime = timeEnd-timeStart
48 if(type == 0):
49     epsilon = 0
50     mu = 0.5
51     gamma = 0.5
52     amountEpisodes = 1000
53     timeStart = time.time()
54     self.train(amountEpisodes, epsilon, mu, gamma)
55     timeEnd = time.time()
56     trainTime = timeEnd-timeStart

58 times = []
59 times.append(setupTime)
60 times.append(trainTime)
61
62 return times
63 return -1

```

Listing 23: Setup and training time is taken here. From file *ReinforcementSarsa.py*, but is the same for all the similar programs.

```

8 def __init__(self, startBoard):
9     self.sum = startBoard[0]+1
10    self.board = startBoard
11    for i in range(1, len(self.board)):
12        self.sum *= self.board[i]+1
13    self.listOfStates = [0] * self.sum
14    self.listOfStates[0] = self.board[:]
15    self.T = []
16    self.Q = []

```

Listing 24: The initialization of the file, which is part of the what is run in the total time only. From file *ReinforcementSarsa.py*, but is the same for all the similar programs.

8.1.10 Phenomena: Irregular Decrease in Success-rate

A phenomena to be observed here, is the fact that the success rate does not decrease in a regular way, when given larger states. It seems to actually be so that states with an even number of heaps has a higher success rate, than states with an odd number of heaps. That is, the algorithm trains better when it starts with a state of an even amount of heaps. The difference between these states is clear, and that is that a state with an even amount of heaps of the same size is a losing state, it has a Nim-sum of 0, while a state with an odd number heaps of the same size is a winning state, its Nim-sum is not 0. Table 14.

Table 14: The strange phenomena, success rate is higher where there are an even number of heaps. From running TS.

Setup	Training	Total	Winrate	Starting State
0.000s	0.019s	0.019s	100%	4 4
0.001s	0.031s	0.032s	84%	4 4 4
0.004s	0.072s	0.076s	99%	4 4 4 4
0.011s	0.102s	0.113s	74%	4 4 4 4 4
0.029s	0.198s	0.228s	86%	4 4 4 4 4 4

Hypothesis I believe that because the even number of heaps is a losing position, the algorithm would be more exploretative, because its first move would mostly change each time, as the policy in use is ϵ -greedy. This means it tries out more positions, which might be beneficial. On the odd number of heaps it might however be too exploitative, and not explore all possible positions, as perhaps it should. If ϵ was increased, making it more explorative, the phenomena should in this case disappear, because its advantage would be nullified. With ϵ set to 1, where in the training every move is made by picking a random move, the expected result would be that the winrate should decrease properly. This would be because any extra exploration for some of the starting states would be nullified, as the policy would just be random moves anyway. The results were:

Results Table 15 shows that the winrate is now regular again, so I conclude that the effect was caused by some states causing more exploration, which is beneficial for the winrate. This makes sense, as the algorithms for reinforcement learning really find the best way to get from the start to the end, while in the case of solving Nim training for all possible states is desirable. Based upon this phenomena a better solution might be found. If the goal is more exploration with the goal being training from all states, I suggest that rather than just starting in the starting state, the training should start games from all possible positions, which might show results similar to those states with an even number of heaps, for all states.

Table 15: Statistics for trimmed sarsa, with ϵ set to 1.

Setup	Training	Total	Winrate	Starting State
0.000s	0.027s	0.027s	39%	5 5
0.003s	0.057s	0.060s	30%	5 5 5
0.011s	0.132s	0.143s	18%	5 5 5 5
0.039s	0.287s	0.327s	11%	5 5 5 5 5
0.001s	0.033s	0.034s	42%	6 6
0.005s	0.078s	0.083s	23%	6 6 6
0.028s	0.207s	0.235s	12%	6 6 6 6
0.129s	0.567s	0.698s	6%	6 6 6 6 6

8.1.11 Training From All States

The program is the similar to the trimmed programs, and uses the sarsa algorithm. The only difference is in the method for the training, of course, since only the training is different in theory. The new method is shown in Listing 25.

```

144 def train(self, amountEpisodes, epsilon, mu, gamma):
145     curStartState = 0
146     for neverUsed in range(amountEpisodes):
147         curState = curStartState
148         playing = 1
149         while(playing):
150             curAction = self.epsilonGreedyChoice(curState, epsilon)
151             curReward = self.Q[curState][curAction]
152             nextState = curAction
153             if(nextState == len(self.listOfStates)-1):
154                 break
155             nextAction = self.epsilonGreedyChoice(nextState, epsilon)
156             self.Q[curState][curAction] -= mu * (curReward + gamma*self.Q
157             [nextState][nextAction] - self.Q[curState][curAction])
158             curState = nextState
159             if (nextAction == self.sum-1):
160                 playing = 0
161             curStartState += 1
162             if(curStartState >= len(self.listOfStates)-2):
163                 curStartState = 0

```

Listing 25: The method *train* in the program that trains from all the states; *SpaceSaveSarsaSpread.py*

Table 16 contains the results for this program compared to the ordinary trimmed program using sarsa. As can be seen the new program does better in most cases, with the exception of some cases where there are an equal amount of heaps. This is expected, as the new program trains less per episode, as it starts from lower states where the game can be beaten with fewer moves. As for the timeuse, the training should be quicker for the new program, and it is.

Table 16: Comparing the Algorithm that trains from all states with Trimmed Sarsa.

Program	Training	Total	Winrate	Starting State
New Program TS	0.041s 0.052s	0.047s 0.058s	98% 73%	6 6 6
New Program TS	0.90s 0.170s	0.122s 0.202s	70% 74%	6 6 6 6
New Program TS	0.204s 0.341s	0.350s 0.482s	41% 38%	6 6 6 6 6
New Program TS	0.052s 0.071s	0.063s 0.081s	88% 73%	7 7 7
New Program TS	0.123s 0.218s	0.196s 0.294s	49% 50%	7 7 7 7
New Program TS	0.293s 0.542s	0.686s 0.973s	26% 22%	7 7 7 7 7
New Program TS	0.071s 0.109s	0.094s 0.131s	70% 60%	8 8 8
New Program TS	0.183s 0.303s	0.350s 0.467s	33% 31%	8 8 8 8
New Program TS	0.491s 0.876s	1.537s 1.957s	16% 14%	8 8 8 8 8
New Program TS	0.066s 0.105s	0.093s 0.132s	57% 55%	9 9 9
New Program TS	0.191s 0.315s	0.450s 0.579s	24% 21%	9 9 9 9
New Program TS	0.594s 0.987s	2.527s 3.047s	10% 8%	9 9 9 9 9

8.1.12 Difference in Time-use, Comparing the Algorithms

It is quite clear that the algorithms using the trimmed game trees uses less time, in fact, as can be seen in Table 17. They are also more successful, that is they have a higher winrate against the perfect player. This is to be expected after all, since a lesser game-tree means that the matrices are faster to set up, searching for the best move is also faster, as there are less to pick from, and there is a smaller Q to train, so finding the proper policy is easier.

Table 17: The timeuse and winrate for the different algorithms when given the gamestate: 5 5 5 5 5.

Program	Setup	Training	Total	Winrate	Starting State
RS	27.906s	3.008s	32.590s	6%	5 5 5 5 5
RQ	26.996s	4.696s	33.137s	6%	5 5 5 5 5
TS	0.037s	0.187s	0.225s	65%	5 5 5 5 5
TQ	0.037s	0.279s	0.317s	71%	5 5 5 5 5
ODS	0.011s	0.291s	0.301s	84%	5 5 5 5 5

It is clear that the growth of the setup time is greater than the growth of the training time for the programs that use the untrimmed game-tree. The same is true for the programs that use the trimmed game-tree. Table 18 Shows 2 examples of this. This is not surprising, given that the parts of the programs that use the most time, is going to be those that have to iterate over the matrices somehow, because the matrices are clearly the largest datastructures in the programs. Since the setup have to go through the whole matrix twice, it should always use the most time as the training will only go through one column the matrix an amount of time, which is at the highest the amount of counters in the startstate times the amount of episodes. For these statistics it always trains with 1000 episodes. As long as a matrix is used this will be the case. Thus we look towards the last program, the one that does away with the matrix altogether, and uses an array instead.

Table 18: How the timeuse develops for greater start-states.

Program	Setup	Training	Total	Winrate	Starting State
RS	0.001s	0.040s	0.041s	100%	1 1 1 1 1
RS	0.026s	0.104s	0.130s	72%	2 2 2 2 2
RS	0.455s	0.351s	0.820s	40 %	3 3 3 3 3
RS	4.529s	1.200s	5.888s	14%	4 4 4 4 4
RS	27.906s	3.008s	32.590s	6%	5 5 5 5 5
TS	0.124s	0.344s	0.472s	43%	5 5 5 5 5 5
TS	0.458s	0.581s	1.051s	22%	6 6 6 6 6 6
TS	1.518s	0.941s	2.499s	12%	7 7 7 7 7 7
TS	4.284s	1.655s	6.056s	6%	8 8 8 8 8 8
TS	11.957s	2.808s	15.306s	4%	9 9 9 9 9 9
TS	32.736s	4.666s	39.09s	2%	10 10 10 10 10 10
ODS	0.009s	0.304s	0.313s	98%	4 4 4 4 4 4
ODS	0.029s	0.722s	0.751s	81%	5 5 5 5 5 5
ODS	0.087s	1.385s	1.474s	46%	6 6 6 6 6 6
ODS	0.257s	2.529s	2.794s	32%	7 7 7 7 7 7
ODS	0.655s	4.147s	4.826s	23%	8 8 8 8 8 8

The last selection of statistics in Table 18 refers to the program that uses a one dimensional array to train. It is incidentally clear how much faster the trimmed programs are, so even though only certain states can be trained from, it is still better. If the input state was changed within the program to be the state of all heaps of the size of the biggest input heap, it would perhaps train much more of the game than necessary, but it would still in most cases be far quicker than the untrimmed programs.

8.1.13 Stochastic Termination

The program *Episodes.py* runs a program until it has been trained to reach a predefined winrate. This table of results are presented to argue against the use of stochastic termination. Table 19 are the result of running *Episodes.py* twice with program TS, with a desired winrate of 90%.

Table 19: The results of stochastic termination of TS

Time used	Episodes used	startstate
0.001s	10	2 2
0.000s	10	
0.046s	2010	2 2 2
0.004s	80	
0.001s	30	2 2 2 2
0.001s	20	
0.001s	20	3 3
0.001s	20	
1.149s	44210	3 3 3
0.069s	3510	
0.007s	90	3 3 3 3
0.012s	130	

This shows how unreliable stochastic termination might be, and how time-consuming training the program actually is, when one desires a high winrate. the training required for 3 3 3, for example, seems hard to predict, and from this it seems obvious the problem will just intensify. The reason that 3 3 3 3 is so much easier would be the phenomena described earlier.

8.1.14 Playing the Game

Thus far all the time-use has been based upon the time-use of the setup and the training. There is one part of the program that has not been explored yet, which is the time-use of the move-selection, that is, the time used when the program finds the best move it can take(based upon its statistics).

Table 20: Time-use of moves

Program	Move	Time-use	Current State
Deterministic	RS	5 counters from heap 4	0.001s 5 5 5 5 5
	1 counter from heap 1	0.000s	5 5 5 0 5
	3 counters from heap 2	0.002s	4 5 5 0 5
Deterministic	TS	4 counters from heap 1	0.000s 8 8 8 8 8 8
	4 counters from heap 1	0.000s	4 8 8 8 8 8 8
	5 counters from heap 2	0.000s	0 8 8 8 8 8 8
ODS	ODS	6 counters from heap 1	0.001s 8 8 8 8 8 8
	Deterministic	2 counters from heap 1	0.000s 2 8 8 8 8 8 8
	ODS	4 counters from heap 2	0.002s 0 8 8 8 8 8 8

Table 20 shows the timeuse for making moves of some selected programs. The states were chosen for being rather large, so that the limits of the programs might be tested. It only shows small sections of complete games, as the rest of the games had as small time-uses. The conclusion here is that making moves is very quick and simple in comparison to the time used for setup and training. That is to be expected as in the training the program might find moves several times, and as such should be slower. The timeuse of the deterministic algorithm is very small, which I predicted when I considered the theoretical time-use of the algorithms.

8.2 Supervised Learning

Separating training from setup doesn't make as much sense for this program, so time-use is just shown in total for the whole program. Table 21 shows that it generally does worse than the untrimmed programs, both in success and time-use. It also does not show the phenomena of higher winrates for states with an equal amount of heaps.

Table 21: Statistics of a run of the program for playing nim with Supervised learning.

Time-use	Winrate	Starting state
0.017s	68%	2 2
0.082s	45%	2 2 2
0.447s	41%	2 2 2 2
1.371s	30%	2 2 2 2 2
0.025s	61%	3 3
0.260s	44%	3 3 3
1.443s	33%	3 3 3 3
9.087s	16%	3 3 3 3 3
0.051s	57%	4 4
0.574s	29%	4 4 4
2.444s	14%	4 4 4 4
26.476s	6%	4 4 4 4 4
0.140s	62%	5 5
0.871s	24%	5 5 5
6.608s	9%	5 5 5 5
48.948s	3%	5 5 5 5 5

9 Conclusion

Solving Nim with machine learning is evidently possible. I am confident in saying that it will never be better than the deterministic solution. This seems quite obvious, as any machine learning solution is by definition stochastic, the supervised solution is undecidable, and the reinforcement solution is based upon the game-tree, which is too big, even when trimmed. The results show the same, so it seems a fair conclusion. That was not surprising however, because the deterministic solution is very good, and this problem is not within the area which machine learning is considered best. What to take from this is rather the troubles faced by the machine when trying to solve Nim, that is the problem of achieving general success when playing. That is that the amount of training is based upon the size of the game-tree, and the tree gets large, fast. In cases where the training graph remains fairly small, and does not increase too much, playing well from all positions from the graph should be feasible, especially of methods which are more explorative are used.

The Importance What is explored here might be interesting when approaching other problems with similar traits, and similar solutions might apply. For example, trimming the graph is important(But that might as well apply in all cases of reinforcement learning), training from all possible starts might help. I suspect that training for all states is more interesting when preparing for a human opponent, as their unpredictability makes being versatile more important.

9.1 Where Now?

There are some extra solutions and statistics that might be interesting, for an expanded understanding of the problem.

- Trying to train from all positions could be applied to ODS, which already has some advantages, to see if that union will create an even more successful program.
- Making a program that just checks how well the reinforcement programs do if the only goal is to play as well as possible from the startstate might be interesting, as that is closer to what reinforcement learning is considered good at.

10 The Code in Full

Here are the programs used in this paper, in full. They are all written in the language Python.

```

1  class PerfectPlayer:
2      #This method applies the deterministic algorithm for making
3      #perfect moves in Nim, and returns the index of the move and the
4      #amount of counters removed.
5      def makeMove(self, state):
6          nimSum = self.findNimSum(state)
7          for i in range(len(state)):
8              remove = nimSum^state[i]
9              if(remove < state[i]):
10                  result = [state[i] - remove, i]
11                  return result
12          for i in range(len(state)):
13              if(state[i] != 0):
14                  result = [1, i]
15                  return result
16          return -1
17
18      #Since finding nimsums is relevant in the program Run.py this
19      #program finds the nimsum of the input state.
20      def findNimSum(self, state):
21          nimSum = 0
22          for i in range(len(state)):
23              nimSum = nimSum^state[i]
24          return nimSum

```

Listing 26: This is the code in the file *PerfectPlayer.py*.

```

1  import time
2  import sys
3  import ReinforcementSarsa as RS
4  import ReinforcementQlearning as RQ
5  import PerfectPlayer as PP
6  import SpaceSaveSarsa as SS
7  import SpaceSaveQlearning as SQ
8  import OneDSarsa as ODS
9  import SpaceSaveSarsaSpread as SSSS

11 if(len(sys.argv) < 11):
12     print("Error: too few arguments. Got " + str(len(sys.argv)-1) + "
13         , Expected 10")
14 elif (len(sys.argv) > 11):
15     print("Error: too many arguments. Got, " + str(len(sys.argv)-1) +
16         ", Expected 10")
17 else:
18     #Parsing the input.
19     heapSizeLower = int(sys.argv[1])
20     heapSizeUpper = int(sys.argv[2])
21     heapAmountLower = int(sys.argv[3])
22     heapAmountUpper = int(sys.argv[4])

player = PP.PerfectPlayer()

```

```

24 ssstatistics = open("ssstatistics.txt", "w")
25
26 #writing the header of the statistics file, to show which numbers
27     #correspond to which program.
28 header = ""
29 if(int(sys.argv[5]) == 1):
30     header += "\t\t\t\t\t\tRS\t"
31 if(int(sys.argv[6]) == 1):
32     header += "\t\t\t\t\t\tRQ\t"
33 if(int(sys.argv[7]) == 1):
34     header += "\t\t\t\t\t\tTS\t"
35 if(int(sys.argv[8]) == 1):
36     header += "\t\t\t\t\t\tTQ\t"
37 if(int(sys.argv[9]) == 1):
38     header += "\t\t\t\t\t\tODS\t"
39 if(int(sys.argv[10]) == 1):
40     header += "\t\t\t\t\t\tNew Program"
41 header += "\n"
42 ssstatistics.write(header)
43 #The loop that defines the start-states from which the programs
44     #will be run.
45 for heapSize in range(heapSizeLower, heapSizeUpper+1):
46     for heapAmount in range(heapAmountLower, heapAmountUpper+1):
47         state = [heapSize] * heapAmount
48         print(state)
49         line = ""
50         """This repeats for each program. For each start-state they
51             are run 10 times, taking their time use, after which they are
52                 set to play against the perfect player
53                     for all starts it can play from, from which it can win when
54                         starting. That is, the states with a nim-sum != 0. The averages
55                             of the times and winrate is then
56                                 written into the statistics file."""
57         if(int(sys.argv[5]) == 1):
58             timeAvg = 0
59             percent = 0
60             setupTime = 0
61             trainTime = 0
62             for i in range(10):
63                 start = time.time()
64                 rs = RS.ReinforcementSarsa(state)
65                 times = rs.setup()
66                 end = time.time()
67                 timeAvg += end-start
68                 setupTime += times[0]
69                 trainTime += times[1]
70                 listOfStates = rs.getListOfStates()
71                 wins = 0
72                 losses = 0
73                 for i in range(len(listOfStates)):
74                     if(player.findNimSum(listOfStates[i]) != 0):
75                         playing = 1
76                         board = listOfStates[i][:]
77                         winBoard = [0] * len(listOfStates[i])
78                         while(playing):
79                             move = rs.makeMove(board)
80                             board[move[1]] = board[move[1]]-move[0]

```

```

75         if(board == winBoard):
76             wins += 1
77             break
78         move = player.makeMove(board)
79         board[move[1]] = board[move[1]]-move[0]
80         if(board == winBoard):
81             losses += 1
82             break
83         percent += (wins/(wins+losses)*100)
84         timeAvg /= 10
85         percent /= 10
86         setupTime /= 10
87         trainTime /= 10
88         line += "setup: " + str(round(setupTime, 3)) + "\t train:
" + str(round(trainTime, 3)) + "\t total: " + str(round(
timeAvg, 3)) + "\t" + str(int(percent)) + "%\t"
89         if(int(sys.argv[6]) == 1):
90             timeAvg = 0
91             percent = 0
92             setupTime = 0
93             trainTime = 0
94             for i in range(10):
95                 start = time.time()
96                 rq = RQ.ReinforcementQlearning(state)
97                 times = rq.setup()
98                 end = time.time()
99                 timeAvg += end-start
100                setupTime += times[0]
101                trainTime += times[1]
102                listOfStates = rq.getListOfStates()
103                wins = 0
104                losses = 0
105                for i in range(len(listOfStates)):
106                    if(player.findNimSum(listOfStates[i]) != 0):
107                        playing = 1
108                        board = listOfStates[i][:]
109                        winBoard = [0] * len(listOfStates[i])
110                        while(playing):
111                            move = rq.makeMove(board)
112                            board[move[1]] = board[move[1]]-move[0]
113                            if(board == winBoard):
114                                wins += 1
115                                break
116                            move = player.makeMove(board)
117                            board[move[1]] = board[move[1]]-move[0]
118                            if(board == winBoard):
119                                losses += 1
120                                break
121                            percent += (wins/(wins+losses)*100)
122                            timeAvg /= 10
123                            percent /= 10
124                            setupTime /= 10
125                            trainTime /= 10
126                            line += "setup: " + str(round(setupTime, 3)) + "\t train:
" + str(round(trainTime, 3)) + "\t total: " + str(round(
timeAvg, 3)) + "\t" + str(int(percent)) + "%\t"
127                            if(int(sys.argv[7]) == 1):

```

```

128     timeAvg = 0
129     percent = 0
130     setupTime = 0
131     trainTime = 0
132     for i in range(10):
133         start = time.time()
134         sss = SS.SpaceSaveSarsa(state)
135         times = sss.setup()
136         end = time.time()
137         timeAvg += end-start
138         setupTime += times[0]
139         trainTime += times[1]
140         listOfStates = sss.getListOfStates()
141         wins = 0
142         losses = 0
143         for i in range(len(listOfStates)):
144             if(player.findNimSum(listOfStates[i]) != 0):
145                 playing = 1
146                 board = listOfStates[i][:]
147                 winBoard = [0] * len(listOfStates[i])
148                 while(playing):
149                     move = sss.makeMove(board)
150                     board[move[1]] = board[move[1]]-move[0]
151                     if(board == winBoard):
152                         wins += 1
153                         break
154                     move = player.makeMove(board)
155                     board[move[1]] = board[move[1]]-move[0]
156                     if(board == winBoard):
157                         losses += 1
158                         break
159                     percent += (wins/(wins+losses)*100)
160             timeAvg /= 10
161             percent /= 10
162             setupTime /= 10
163             trainTime /= 10
164             line += "setup: " + str(round(setupTime, 3)) + "\t train:
" + str(round(trainTime, 3)) + "\t total: " + str(round(
timeAvg, 3)) + "\t" + str(int(percent)) + "%\t"
165             if(int(sys.argv[8]) == 1):
166                 timeAvg = 0
167                 percent = 0
168                 setupTime = 0
169                 trainTime = 0
170                 for i in range(10):
171                     start = time.time()
172                     ssq = SQ.SpaceSaveQlearning(state)
173                     times = ssq.setup()
174                     end = time.time()
175                     timeAvg += end-start
176                     setupTime += times[0]
177                     trainTime += times[1]
178                     listOfStates = ssq.getListOfStates()
179                     wins = 0
180                     losses = 0
181                     for i in range(len(listOfStates)):
182                         if(player.findNimSum(listOfStates[i]) != 0):

```

```

183         playing = 1
184         board = listOfStates[i][:]
185         winBoard = [0] * len(listOfStates[i])
186         while(playing):
187             move = ssq.makeMove(board)
188             board[move[1]] = board[move[1]]-move[0]
189             if(board == winBoard):
190                 wins += 1
191                 break
192             move = player.makeMove(board)
193             board[move[1]] = board[move[1]]-move[0]
194             if(board == winBoard):
195                 losses += 1
196                 break
197             percent += (wins/(wins+losses)*100)
198             timeAvg /= 10
199             percent /= 10
200             setupTime /= 10
201             trainTime /= 10
202             line += "setup: " + str(round(setupTime, 3)) + "\t train:"
203             line += str(round(trainTime, 3)) + "\t total: " + str(round(
204                 timeAvg, 3)) + "\t" + str(int(percent)) + "%\t"
205             if(int(sys.argv[9]) == 1):
206                 timeAvg = 0
207                 percent = 0
208                 setupTime = 0
209                 trainTime = 0
210                 for i in range(10):
211                     start = time.time()
212                     ods = ODS.OneDSarsa(state)
213                     times = ods.setup()
214                     end = time.time()
215                     timeAvg += end-start
216                     setupTime += times[0]
217                     trainTime += times[1]
218                     listOfStates = ods.getListOfStates()
219                     wins = 0
220                     losses = 0
221                     for i in range(len(listOfStates)):
222                         if(player.findNimSum(listOfStates[i]) != 0):
223                             playing = 1
224                             board = listOfStates[i][:]
225                             winBoard = [0] * len(listOfStates[i])
226                             while(playing):
227                                 move = ods.makeMove(board)
228                                 board[move[1]] = board[move[1]]-move[0]
229                                 if(board == winBoard):
230                                     wins += 1
231                                     break
232                                 move = player.makeMove(board)
233                                 board[move[1]] = board[move[1]]-move[0]
234                                 if(board == winBoard):
235                                     losses += 1
236                                     break
237             percent += (wins/(wins+losses)*100)
238             timeAvg /= 10
239             percent /= 10

```

```

238     setupTime /= 10
239     trainTime /= 10
240     line += "setup: " + str(round(setupTime, 3)) + "\t train:
241     " + str(round(trainTime, 3)) + "\t total: " + str(round(
242         timeAvg, 3)) + "\t" + str(int(percent)) + "%\t"
243     if(int(sys.argv[10]) == 1):
244         timeAvg = 0
245         percent = 0
246         setupTime = 0
247         trainTime = 0
248         for i in range(10):
249             start = time.time()
250             ssss = SSSS.SpaceSaveSarsaSpread(state)
251             times = ssss.setup()
252             end = time.time()
253             timeAvg += end-start
254             setupTime += times[0]
255             trainTime += times[1]
256             listOfStates = ssss.getListOfStates()
257             wins = 0
258             losses = 0
259             for i in range(len(listOfStates)):
260                 if(player.findNimSum(listOfStates[i]) != 0):
261                     playing = 1
262                     board = listOfStates[i][:]
263                     winBoard = [0] * len(listOfStates[i])
264                     while(playing):
265                         move = ssss.makeMove(board)
266                         board[move[1]] = board[move[1]]-move[0]
267                         if(board == winBoard):
268                             wins += 1
269                             break
270                         move = player.makeMove(board)
271                         board[move[1]] = board[move[1]]-move[0]
272                         if(board == winBoard):
273                             losses += 1
274                             break
275                         percent += (wins/(wins+losses)*100)
276             timeAvg /= 10
277             percent /= 10
278             setupTime /= 10
279             trainTime /= 10
280             line += "setup: " + str(round(setupTime, 3)) + "\t train:
281             " + str(round(trainTime, 3)) + "\t total: " + str(round(
282                 timeAvg, 3)) + "\t" + str(int(percent)) + "%\t"
283             line += str(state) + "\n"
284             ssstatistics.write(line)

```

Listing 27: This is the code in the file *Run.py*.

```

1 import time
2 import sys

4 import ReinforcementSarsa as RS
5 import ReinforcementQlearning as RQ
6 import PerfectPlayer as PP
7 import SpaceSaveSarsa as SS
8 import SpaceSaveQlearning as SQ

```

```

9 import OneDSarsa as ODS
10 import SpaceSaveSarsaSpread as SSSS

12 if(len(sys.argv) < 5):
13     print("Error: too few arguments. Got " + str(len(sys.argv)-1) + "
14         , Expected 4")
14 elif (len(sys.argv) > 5):
15     print("Error: too many arguments. Got, " + str(len(sys.argv)-1) +
16         ", Expected 4")
16 else:
17     #Parsing the input
18     startState = list(map(int, sys.argv[3].split(',')))
19     trainingType = int(sys.argv[2])
20     if(trainingType == 0):
21         ml = RS.ReinforcementSarsa(startState)
22     elif(trainingType == 1):
23         ml = RQ.ReinforcementQlearning(startState)
24     elif(trainingType == 2):
25         ml = SS.SpaceSaveSarsa(startState)
26     elif(trainingType == 3):
27         ml = SQ.SpaceSaveQlearning(startState)
28     elif(trainingType == 4):
29         ml = ODS.OneDSarsa(startState)
30     else:
31         ml = SSSS.SpaceSaveSarsaSpread(startState)
32     ml.setup()

34 #The code for the reinforcement program playing against the
34 deterministic program, and getting statistics from it.
35 if(int(sys.argv[1]) == 0):
36     playstatistics = open("playstatistics.txt", "w")
37     playstatistics.write(str(trainingType) + "\n")
38     player = PP.PerfectPlayer()
39     playing = 1
40     board = startState[:]
41     winBoard = [0] * len(startState)
42     #If the reinforcement algo gets the first move it gets to make
42     #a move before the loop start.
43     if(int(sys.argv[4]) == 1):
44         start = time.time()
45         move = player.makeMove(board)
46         end = time.time()
47         moveTime = end-start
48         board[move[1]] = board[move[1]]-move[0]
49         playstatistics.write("Move by Pre-algo, took " + str(move[0])
49         + " counters from heap " + str(move[1]+1) + ". New state: " +
50         str(board) + " time used: " + str(round(moveTime, 3)) + "s.\n")
51         if(board == winBoard):
52             playstatistics.write("Pre-algo is the winner.")
53             playing = 0
53     #Otherwise the game starts here.
54     while(playing):
55         start = time.time()
56         move = ml.makeMove(board)
57         end = time.time()
58         moveTime = end-start
59         board[move[1]] = board[move[1]]-move[0]

```

```

60     playstatistics.write("Move by ML, took " + str(move[0]) + " "
61     counters from heap " + str(move[1]+1) + ". New state: " + str(
62     board) + " time used: " + str(round(moveTime, 3)) + "s.\n")
63     if(board == winBoard):
64         playstatistics.write("ML is the winner.")
65         break
66     start = time.time()
67     move = player.makeMove(board)
68     end = time.time()
69     moveTime = end-start
70     board[move[1]] = board[move[1]]-move[0]
71     playstatistics.write("Move by Pre-algo, took " + str(move[0])
72     + " counters from heap " + str(move[1]+1) + ". New state: " +
73     str(board) + " time used: " + str(round(moveTime, 3)) + "s.\n")
74     if(board == winBoard):
75         playstatistics.write("Pre-algo is the winner.")
76         break
77     #If the player wants to play against the Reinforcement algo.
78     else:
79         player = PP.PerfectPlayer()
80         playing = 1
81         board = startState[:]
82         winBoard = [0] * len(startState)
83         #If the player wants to start first, they make a move first.
84         if(int(sys.argv[4]) == 1):
85             move = [-1, -1]
86             print("Current board is: " + str(board) + "\n")
87             move[1] = int(input("Which heap do you want to remove from? "
88            ))-1
89             while((move[1] < 0 or move[1] >= len(board)) or board[move
90             [1]] == 0):
91                 move[1] = int(input("Illegal heap, choose again. "))-1
92                 move[0] = int(input("How many counters do you wish to remove?
93             "))
94                 while(move[0] < 1 or move[0] > board[move[1]]):
95                     move[0] = int(input("Illegal amount, choose again. "))
96                     board[move[1]] = board[move[1]]-move[0]
97                     if(board == winBoard):
98                         print("You are the winner.")
99                         playing = 0
100                    #Otherwise the game starts here, with the reinforcement program
101                    #making the first move.
102                    while(playing):
103                        move = ml.makeMove(board)
104                        board[move[1]] = board[move[1]]-move[0]
105                        print("Move by ML, took " + str(move[0]) + " counters from
heap " + str(move[1]+1) + ".")
106                        if(board == winBoard):
107                            print("ML is the winner.")
108                            break
109                        move = [-1, -1]
110                        print("Current board is: " + str(board) + "\n")
111                        move[1] = int(input("Which heap do you want to remove from? "
112                        ))-1
113                        while((move[1] < 0 or move[1] >= len(board) or board[move[1]]
114                         == 0)):
115                            move[1] = int(input("Illegal heap, choose again. "))-1

```

```

106     move[0] = int(input("How many counters do you wish to remove?
107     "))
108     while(move[0] < 1 or move[0] > board[move[1]]):
109         move[0] = int(input("Illegal amount, choose again. "))
110         board[move[1]] = board[move[1]]-move[0]
111         if(board == winBoard):
112             print("You are the winner.")
113             break

```

Listing 28: This is the code in the file *Play.py*.

```

1 import numpy as np
2 import random
3 import time

5 class ReinforcementSarsa:
6
7     #The initialization takes the starting board, and finds the
8     #amount of states in its untrimmed game-tree, which is called
9     #self.sum. It also initializes T and Q.
10    def __init__(self, startBoard):
11        self.sum = startBoard[0]+1
12        self.board = startBoard
13        for i in range(1, len(self.board)):
14            self.sum *= self.board[i]+1
15        self.listOfStates = [0] * self.sum
16        self.listOfStates[0] = self.board[:]
17        self.T = []
18        self.Q = []
19
20        """This method is for setting up the matrices T and Q, as well as
21        creating the list of states. For each state it makes all
22        possible moves it can,
23        and for each of those, it adds them to their spot in listOfStates
24        , if the state wasnt already there, using the method newIndex
25        to find the states' index.
26        All those moves are also added to T, which searching is the main
27        purpose. Then Q is made, by setting all positions in T with 1,
28        to have small random values in Q.
29        The other positions in Q gets the value of negative infinity, -np
30        .inf, so that they are always smaller than the expected reward
31        of a move.
32        Lastly the method sets up the inputs for the training, and runs
33        it as well. Lastly the method takes the time used both by the
34        setup and the training."""
35    def setup(self, type):
36        timeStart = time.time()
37        for i in range(self.sum):
38            self.T.append([0] * self.sum)
39            for x in range(len(self.listOfStates[i])):
40                for y in range(1, self.listOfStates[i][x] + 1):
41                    tmpValue = self.newIndex(x, y)
42                    self.T[i][i+tmpValue] = 1
43                    if(self.listOfStates[i+tmpValue] == 0):
44                        tmpList = self.listOfStates[i][:]
45                        tmpList[x] -= y
46                        self.listOfStates[i+tmpValue] = tmpList[:]

```

```

36     self.Q.append([])
37     for j in range(self.sum-1):
38         if(self.T[i][j] == 1):
39             self.Q[i].append(random.random())
40         else:
41             self.Q[i].append(-np.inf)
42         if (self.T[i][self.sum-1] == 1):
43             self.Q[i].append(100.0)
44         else:
45             self.Q[i].append(-np.inf)
46     timeEnd = time.time()
47     setupTime = timeEnd-timeStart
48     if(type == 0):
49         epsilon = 0
50         mu = 0.5
51         gamma = 0.5
52         amountEpisodes = 1000
53         timeStart = time.time()
54         self.train(amountEpisodes, epsilon, mu, gamma)
55         timeEnd = time.time()
56         trainTime = timeEnd-timeStart
57
58     times = []
59     times.append(setupTime)
60     times.append(trainTime)
61
62     return times
63     return -1
64
65 # This method exists for finding, given a move, how far down the
66 # list of states the new state is.
67 def newIndex(self, heapIndex, amount):
68     newMove = 1
69     for i in range(len(self.board)-1, heapIndex, -1):
70         newMove *= self.board[i]+1
71     newMove *= amount
72     return newMove
73
74 """This is the method that finds the best state to make a move to
75 , from the input state. It does this by searching through list
76 to find the current state,
77 and then using np.argmax to find the following state with the
78 highest expected reward. Then it finds the heap where the move
79 was made,
80 to find exactly what to return, that is the amount of counters
81 removed and the index of the heap removed from."""
82 def makeMove(self, state):
83     curState = state[:]
84     nextState = []
85     move = []
86     if (len(state)<len(self.listOfStates[0])):
87         dif = len(self.listOfStates[0])-len(state)
88         for i in range(dif):
89             curState.append(0)
90     for i in range(len(self.listOfStates)):
91         if (curState == self.listOfStates[i]):
92             nextState = self.listOfStates[np.argmax(self.Q[i])]
```

```

87     for i in range(len(curState)):
88         if (curState[i] != nextState[i]):
89             move.append(curState[i] - nextState[i])
90             move.append(i)
91     return move
92
93     """This method is for training the matrix Q, by using the Sarsa
94     algorithm. The main difference is that it removes the
95     calculated reward, instead of adding it.
96     It trains by playing against itself, until it wins, and then
97     starts again. each game is an episode."""
98     def train(self, amountEpisodes, epsilon, mu, gamma):
99         for neverUsed in range(amountEpisodes):
100             curState = 0
101             curAction = self.epsilonGreedyChoice(curState, epsilon)
102             playing = 1
103             while(playing):
104                 curReward = self.Q[curState][curAction]
105                 nextState = curAction
106                 nextAction = self.epsilonGreedyChoice(nextState, epsilon)
107                 self.Q[curState][curAction] -= mu * (curReward + gamma*self
108                 .Q[nextState][nextAction] - self.Q[curState][curAction])
109                 curState = nextState
110                 curAction = nextAction
111                 if (curAction == self.sum-1):
112                     playing = 0
113
114             #This is an implementation of the epsilon greedy policy, by the
115             #use of numpy.argmax when finding the best move, or just random
116             #choice of the possible moves otherwise.
117             def epsilonGreedyChoice(self, state, epsilon):
118                 if(random.random()<epsilon):
119                     tmpList = []
120                     for i in range(state, self.sum):
121                         if(self.T[state][i] == 1):
122                             tmpList.append(i)
123                     return tmpList[random.randrange(0, len(tmpList))]
124                 else:
125                     return np.argmax(self.Q[state])
126
127             #This method returns the list of states, which is used to test
128             #the program, by having it play from many different states.
129             def getListOfStates(self):
130                 return self.listOfStates

```

Listing 29: This is the code in the file *ReinforcementSarsa.py*.

```

1 import numpy as np
2 import random
3 import time
4
5 class ReinforcementQlearning:
6
7     #The initialization takes the starting board, and finds the
8     #amount of states in its untrimmed game-tree, which is called
9     #self.sum. It also initializes T and Q.
10    def __init__(self, startBoard):
11        self.sum = startBoard[0]+1

```

```

10     self.board = startBoard
11     for i in range(1, len(self.board)):
12         self.sum *= self.board[i]+1
13     self.listOfStates = [0] * self.sum
14     self.listOfStates[0] = self.board[:]
15     self.T = []
16     self.Q = []
17
18     """This method is for setting up the matrices T and Q, as well as
19     creating the list of states. For each state it makes all
20     possible moves it can,
21     and for each of those, it adds them to their spot in listOfStates
22     , if the state wasnt already there, using the method newIndex
23     to find the states' index.
24     All those moves are also added to T, which searching is the main
25     purpose. Then Q is made, by setting all positions in T with 1,
26     to have small random values in Q.
27     The other positions in Q gets the value of negative infinity, -np
28     .inf, so that they are always smaller than the expected reward
29     of a move.
30     Lastly the method sets up the inputs for the training, and runs
31     it as well. Lastly the method takes the time used both by the
32     setup and the training."""
33
34     def setup(self, type):
35         timeStart = time.time()
36         for i in range(self.sum):
37             self.T.append([0] * self.sum)
38             for x in range(len(self.listOfStates[i])):
39                 for y in range(1, self.listOfStates[i][x] + 1):
40                     tmpValue = self newIndex(x, y)
41                     self.T[i][i+tmpValue] = 1
42                     if(self.listOfStates[i+tmpValue] == 0):
43                         tmpList = self.listOfStates[i][:]
44                         tmpList[x] -= y
45                         self.listOfStates[i+tmpValue] = tmpList[:]
46
47             self.Q.append([])
48             for j in range(self.sum-1):
49                 if(self.T[i][j] == 1):
50                     self.Q[i].append(random.random())
51                 else:
52                     self.Q[i].append(-np.inf)
53                 if (self.T[i][self.sum-1] == 1):
54                     self.Q[i].append(100.0)
55                 else:
56                     self.Q[i].append(-np.inf)
57         timeEnd = time.time()
58         setupTime = timeEnd-timeStart
59         if(type == 0):
60             epsilon = 0
61             mu = 0.5
62             gamma = 0.5
63             amountEpisodes = 1000
64             timeStart = time.time()
65             self.train(amountEpisodes, epsilon, mu, gamma)
66             timeEnd = time.time()
67             trainTime = timeEnd-timeStart

```

```

57
58     times = []
59     times.append(setupTime)
60     times.append(trainTime)
61
62     return times
63     return -1
64
65 # This method exists for finding, given a move, how far down the
66 # list of states the new state is.
67 def newIndex(self, heapIndex, amount):
68     newMove = 1
69     for i in range(len(self.board)-1, heapIndex, -1):
70         newMove *= self.board[i]+1
71     newMove *= amount
72     return newMove
73
74 """This is the method that finds the best state to make a move to
75 , from the input state. It does this by searching through list
76 to find the current state,
77 and then using np.argmax to find the following state with the
78 highest expected reward. Then it finds the heap where the move
79 was made,
80 to find exactly what to return, that is the amount of counters
81 removed and the index of the heap removed from."""
82 def makeMove(self, state):
83     curState = state[:]
84     nextState = []
85     move = []
86     if (len(state)<len(self.listOfStates[0])):
87         dif = len(self.listOfStates[0])-len(state)
88         for i in range(dif):
89             curState.append(0)
90     for i in range(len(self.listOfStates)):
91         if (curState == self.listOfStates[i]):
92             nextState = self.listOfStates[np.argmax(self.Q[i])]
93     for i in range(len(curState)):
94         if (curState[i] != nextState[i]):
95             move.append(curState[i] - nextState[i])
96         move.append(i)
97     return move
98
99 """This method is for training the matrix Q, by using the Q-
100 learning algorithm. The main difference is that it removes the
101 calculated reward, instead of adding it.
102 It trains by playing against itself, until it wins, and then
103 starts again. each game is an episode."""
104 def train(self, amountEpisodes, epsilon, mu, gamma):
105     for neverUsed in range(amountEpisodes):
106         curState = 0
107         playing = 1
108         while(playing):
109             curAction = self.epsilonGreedyChoice(curState, epsilon)
110             curReward = self.Q[curState][curAction]
111             nextState = curAction
112             nextAction = np.argmax(self.Q[nextState])
113             self.Q[curState][curAction] -= mu * (curReward + gamma*self

```

```

105     .Q[nextState][nextAction] - self.Q[curState][curAction])
106     curState = nextState
107     if (nextAction == self.sum-1):
108         playing = 0
109
110     #This is an implementation of the epsilon greedy policy, by the
111     #use of numpy.argmax when finding the best move, or just random
112     #choice of the possible moves otherwise.
113     def epsilonGreedyChoice(self, state, epsilon):
114         if(random.random()<epsilon):
115             tmpList = []
116             for i in range(state, self.sum):
117                 if(self.T[state][i] == 1):
118                     tmpList.append(i)
119             return tmpList[random.randrange(0, len(tmpList))]
120         else:
121             return np.argmax(self.Q[state])
122
123     #This method returns the list of states, which is used to test
124     #the program, by having it play from many different states.
125     def getListOfStates(self):
126         return self.listOfStates

```

Listing 30: This is the code in the file *ReinforcementQlearning.py*.

```

1 import numpy as np
2 import random
3 import math
4 import time
5
6 class SpaceSaveSarsa:
7
8     #The initialization takes the starting board, and finds the
9     #amount of states in its untrimmed game-tree, which is called
10    #self.sum. It also initializes T and Q.
11    def __init__(self, startBoard):
12        self.board = startBoard
13        self.sum = self.getSum()
14        self.board.sort(reverse = True)
15        self.listOfStates = []
16        self.listOfStates.append(self.board[:])
17        self.T = []
18        self.Q = []
19
20    #This method is for finding the amount of states reachable from
21    #the start-state, for the trimmed game-tree.
22    def getSum(self):
23        return int(math.factorial(len(self.board)+ self.board[0])/(math
24        .factorial(len(self.board))*math.factorial(self.board[0])))
25
26    #This method is for taking a state and a move in the state, and
27    #returning the resulting state in the correct order.
28    def newState(self, heapIndex, amount, state):
29        curState = state[:]
30        curAmount = amount
31        for i in range(heapIndex, len(state)-1):
32            dif = curState[i] - curState[i+1]
33            if (dif >= curAmount):

```

```

29         curState[i] -= curAmount
30         curAmount = 0
31         break
32     else:
33         curState[i] -= dif
34         curAmount -= dif
35     if (curAmount > 0):
36         curState[len(curState)-1] -= curAmount
37     return curState

38 """This method is for setting up the matrices T and Q, as well as
39     creating the list of states. For each state it makes all
40     possible moves it can,
41     and for each of those, it adds them to their spot in listOfStates
42     . The function newState is used to get the resulting state from
43     a given move.
44     the states are added at the end of the list of states in this
45     case, unless they are already in the list.
46     All those moves are also added to T, which is the main purpose.
47     Then Q is made, by setting all positions in T with 1, to have
48     small random values in Q.
49     The other positions in Q gets the value of negative infinity, -np
50     .inf, so that they are always smaller than the expected reward
51     of a move.
52     Lastly the method sets up the inputs for the training, and runs
53     it as well. Lastly the method takes the time used both by the
54     setup and the training."""
55 def setup(self, type = 0):
56     timeStart = time.time()
57     emptyBoard = [0] * len(self.listOfStates[0])
58     for i in range(self.sum):
59         self.T.append([0] * self.sum)
60         for x in range(len(self.listOfStates[i])-1, -1, -1):
61             atEnd = 0
62             curIndex = i+1
63             for y in range(1, self.listOfStates[i][x]+1):
64                 newState = self.newState(x, y, self.listOfStates[i])
65                 if(atEnd):
66                     self.T[i][curIndex] = 1
67                     self.listOfStates.append(newState)
68                     curIndex += 1
69                 else:
70                     atEnd = 1
71                     for j in range(curIndex, len(self.listOfStates)):
72                         if (newState == self.listOfStates[j]):
73                             self.T[i][j] = 1
74                             curIndex = j+1
75                             atEnd = 0
76                             break
77                         if (atEnd):
78                             self.T[i][len(self.listOfStates)] = 1
79                             self.listOfStates.append(newState)
80                             curIndex = len(self.listOfStates)
81             self.Q.append([])
82             for j in range(self.sum-1):
83                 if(self.T[i][j] == 1):
84                     self.Q[i].append(random.random())

```

```

75     else:
76         self.Q[i].append(-np.inf)
77     if (self.T[i][self.sum-1] == 1):
78         self.Q[i].append(100.0)
79     else:
80         self.Q[i].append(-np.inf)
81 timeEnd = time.time()
82 setupTime = timeEnd-timeStart
83 if(type == 0):
84     epsilon = 0.1
85     mu = 0.5
86     gamma = 0.5
87     amountEpisodes = 1000
88     timeStart = time.time()
89     self.train(amountEpisodes, epsilon, mu, gamma)
90     timeEnd = time.time()
91     trainTime = timeEnd-timeStart
92
93     times = []
94     times.append(setupTime)
95     times.append(trainTime)
96
97     return times
98 return -1

99 """
100     This is the method for making a move, it finds the state(
101         which is sorted first) in the list of states, by calculating
102         its index. Then the best move is found by use
103         of the function np.argmax. Lastly it identifies the index of the
104         heap in which the move is made, and it finds a heap of the same
105         size
106     in the input state."""
107 def makeMove(self, state):
108     curState = state[:]
109     curState.sort(reverse = True)
110     nextState = []
111     move = []
112     if (len(state)<len(self.listOfStates[0])):
113         dif = len(self.listOfStates[0])-len(state)
114         for i in range(dif):
115             curState.append(0)
116     goal = 0
117     progress = self.board[0]
118     for i in range(len(self.board)):
119         heaps = len(self.board)-(i+1)
120         for j in range(progress, -1, -1):
121             if(j == curState[i]):
122                 break
123             else:
124                 if(heaps == 0):
125                     goal += j-curState[i]
126                     break
127                 result = 1
128                 for x in range(1, heaps + 1):
129                     result *= (j + x)
130                 goal += (result)//math.factorial(heaps)
131                 progress -= 1

```

```

128     nextState = self.listOfStates[np.argmax(self.Q[goal])][:]
129     numToDec = 0
130     for i in range(len(curState)-1, -1, -1):
131         if (curState[i] != nextState[i]):
132             numToDec += curState[i]-nextState[i]
133             newIndex = curState[i]
134     move.append(numToDec)
135     for i in range(len(state)):
136         if(newIndex == state[i]):
137             newIndex = i
138             break
139     move.append(newIndex)
140     return move

142 """This method is for training the matrix Q, by using the Sarsa
143 algorithm. The main difference is that it removes the
144 calculated reward, instead of adding it.
145 It trains by playing against itself, until it wins, and then
146 starts again. each game is an episode."""
147 def train(self, amountEpisodes, epsilon, mu, gamma):
148     for neverUsed in range(amountEpisodes):
149         curState = 0
150         curAction = self.epsilonGreedyChoice(curState, epsilon)
151         playing = 1
152         while(playing):
153             curReward = self.Q[curState][curAction]
154             nextState = curAction
155             nextAction = self.epsilonGreedyChoice(nextState, epsilon)
156             self.Q[curState][curAction] -= mu * (curReward + gamma*self
157             .Q[nextState][nextAction] - self.Q[curState][curAction])
158             curState = nextState
159             curAction = nextAction
160             if (curAction == self.sum-1):
161                 playing = 0

162 #This is an implementation of the epsilon greedy policy, by the
163 #use of numpy.argmax when finding the best move, or just random
164 #choice of the possible moves otherwise.
165 def epsilonGreedyChoice(self, state, epsilon):
166     if(random.random()<epsilon):
167         tmpList = []
168         for i in range(state, self.sum):
169             if(self.T[state][i] == 1):
170                 tmpList.append(i)
171         return tmpList[random.randrange(0, len(tmpList))]
172     else:
173         return np.argmax(self.Q[state])

174 #This method returns the list of states, which is used to test
175 #the program, by having it play from many different states.
176 def getListOfStates(self):
177     return self.listOfStates

```

Listing 31: This is the code in the file *SpaceSaveSarsa.py*.

```

1 import numpy as np
2 import random
3 import math

```

```

4 import time
6
6 class SpaceSaveQlearning:
8
8     #The initialization takes the starting board, and finds the
9         amount of states in its untrimmed game-tree, which is called
10        self.sum. It also initializes T and Q.
11
11 def __init__(self, startBoard):
12     self.board = startBoard
13     self.sum = self.getSum()
14     self.board.sort(reverse = True)
15     self.listOfStates = []
16     self.listOfStates.append(self.board[:])
17     self.T = []
18     self.Q = []
19
20     # This method is for finding the amount of states reachable from
21         the start-state, for the trimmed game-tree.
22
22 def getSum(self):
23     return int(math.factorial(len(self.board)+ self.board[0])/(math
24         .factorial(len(self.board))*math.factorial(self.board[0])))
25
26     #This method is for taking a state and a move in the state, and
27         returning the resulting state in the correct order.
28
28 def newState(self, heapIndex, amount, state):
29     curState = state[:]
30     curAmount = amount
31
32     for i in range(heapIndex, len(state)-1):
33         dif = curState[i] - curState[i+1]
34         if (dif >= curAmount):
35             curState[i] -= curAmount
36             curAmount = 0
37             break
38         else:
39             curState[i] -= dif
40             curAmount -= dif
41         if (curAmount > 0):
42             curState[len(curState)-1] -= curAmount
43     return curState
44
45     """This method is for setting up the matrices T and Q, as well as
46         creating the list of states. For each state it makes all
47         possible moves it can,
48     and for each of those, it adds them to their spot in listOfStates
49         . The function newState is used to get the resulting state from
50         a given move.
51     the states are added at the end of the list of states in this
52         case, unless they are already in the list.
53     All those moves are also added to T, which is the main purpose.
54         Then Q is made, by setting all positions in T with 1, to have
55         small random values in Q.
56     The other positions in Q gets the value of negative infinity, -np
57         .inf, so that they are always smaller than the expected reward
58         of a move.
59     Lastly the method sets up the inputs for the training, and runs
60         it as well. Lastly the method takes the time used both by the
61         setup and the training."""

```

```

45     def setup(self, type):
46         timeStart = time.time()
47         emptyBoard = [0] * len(self.listOfStates[0])
48         for i in range(self.sum):
49             self.T.append([0] * self.sum)
50             for x in range(len(self.listOfStates[i])-1, -1, -1):
51                 atEnd = 0
52                 curIndex = i+1
53                 for y in range(1, self.listOfStates[i][x]+1):
54                     newState = self.newState(x, y, self.listOfStates[i])
55                     if(atEnd):
56                         self.T[i][curIndex] = 1
57                         self.listOfStates.append(newState)
58                         curIndex += 1
59                     else:
60                         atEnd = 1
61                         for j in range(curIndex, len(self.listOfStates)):
62                             if (newState == self.listOfStates[j]):
63                                 self.T[i][j] = 1
64                                 curIndex = j+1
65                                 atEnd = 0
66                                 break
67                             if (atEnd):
68                                 self.T[i][len(self.listOfStates)] = 1
69                                 self.listOfStates.append(newState)
70                                 curIndex = len(self.listOfStates)
71             self.Q.append([])
72             for j in range(self.sum-1):
73                 if(self.T[i][j] == 1):
74                     self.Q[i].append(random.random())
75                 else:
76                     self.Q[i].append(-np.inf)
77                 if (self.T[i][self.sum-1] == 1):
78                     self.Q[i].append(100.0)
79                 else:
80                     self.Q[i].append(-np.inf)
81             timeEnd = time.time()
82             setupTime = timeEnd-timeStart
83             if(type == 0):
84                 epsilon = 0
85                 mu = 0.5
86                 gamma = 0.5
87                 amountEpisodes = 1000
88                 timeStart = time.time()
89                 self.train(amountEpisodes, epsilon, mu, gamma)
90                 timeEnd = time.time()
91                 trainTime = timeEnd-timeStart
92
93             times = []
94             times.append(setupTime)
95             times.append(trainTime)
96
97             return times
98         return -1
99
100    """ This is the method for making a move, it finds the state(
101        which is sorted first) in the list of states, by calculating

```

```

    its index. Then the best move is found by use
101   of the function np.argmax. Lastly it identifies the index of the
102     heap in which the move is made, and it finds a heap of the same
103       size
104     in the input state."""
105
106   def makeMove(self, state):
107     curState = state[:]
108     curState.sort(reverse = True)
109     nextState = []
110     move = []
111     if (len(state)<len(self.listOfStates[0])):
112       dif = len(self.listOfStates[0])-len(state)
113       for i in range(dif):
114         curState.append(0)
115     goal = 0
116     progress = self.board[0]
117     for i in range(len(self.board)):
118       heaps = len(self.board)-(i+1)
119       for j in range(progress, -1, -1):
120         if(j == curState[i]):
121           break
122         else:
123           if(heaps == 0):
124             goal += j-curState[i]
125             break
126           result = 1
127           for x in range(1, heaps + 1):
128             result *= (j + x)
129           goal += (result)//math.factorial(heaps)
130           progress -= 1
131     nextState = self.listOfStates[np.argmax(self.Q[goal])][:]
132     numToDec = 0
133     for i in range(len(curState)-1, -1, -1):
134       if (curState[i] != nextState[i]):
135         numToDec += curState[i]-nextState[i]
136         newIndex = curState[i]
137       move.append(numToDec)
138       for i in range(len(state)):
139         if(newIndex == state[i]):
140           newIndex = i
141           break
142       move.append(newIndex)
143     return move
144
145   """This method is for training the matrix Q, by using the Q-
146     learning algorithm. The main difference is that it removes the
147       calculated reward, instead of adding it.
148   It trains by playing against itself, until it wins, and then
149     starts again. each game is an episode."""
150
151   def train(self, amountEpisodes, epsilon, mu, gamma):
152     for neverUsed in range(amountEpisodes):
153       curState = 0
154       playing = 1
155       while(playing):
156         curAction = self.epsilonGreedyChoice(curState, epsilon)
157         curReward = self.Q[curState][curAction]
158         nextState = curAction

```

```

152     nextAction = np.argmax(self.Q[nextState])
153     self.Q[curState][curAction] -= mu * (curReward + gamma*self
154     .Q[nextState][nextAction] - self.Q[curState][curAction])
155     curState = nextState
156     if (nextAction == self.sum-1):
157         playing = 0
158
#This is an implementation of the epsilon greedy policy, by the
#use of numpy.argmax when finding the best move, or just random
#choice of the possible moves otherwise.
159 def epsilonGreedyChoice(self, state, epsilon):
160     if(random.random()<epsilon):
161         tmpList = []
162         for i in range(state, self.sum):
163             if(self.T[state][i] == 1):
164                 tmpList.append(i)
165         return tmpList[random.randrange(0, len(tmpList))]
166     else:
167         return np.argmax(self.Q[state])
168
#This method returns the list of states, which is used to test
#the program, by having it play from many different states.
169 def getListOfStates(self):
170     return self.listOfStates
171

```

Listing 32: This is the code in the file *SpaceSaveQlearning.py*.

```

1 import numpy as np
2 import random
3 import math
4 import time
5
6 class OneDSarsa:
7
#The initialization takes the starting board, and finds the
#amount of states in its untrimmed game-tree, which is called
#self.sum. It also initializes T and Q.
8     def __init__(self, startBoard):
9         self.board = startBoard
10        self.sum = self.getSum()
11        self.board.sort(reverse = True)
12        self.listOfStates = []
13        self.listOfStates.append(self.board[:])
14        self.T = []
15        self.Q = []
16
#This method is for finding the amount of states reachable from
#the start-state, for the trimmed game-tree.
17     def getSum(self):
18         return int(math.factorial(len(self.board)+ self.board[0])/(math
19             .factorial(len(self.board))*math.factorial(self.board[0])))
20
#This method is for taking a state and a move in the state, and
#returning the resulting state in the correct order.
21     def newState(self, heapIndex, amount, state):
22         curState = state[:]
23         curAmount = amount
24         for i in range(heapIndex, len(state)-1):
25

```

```

27     dif = curState[i] - curState[i+1]
28     if (dif >= curAmount):
29         curState[i] -= curAmount
30         curAmount = 0
31         break
32     else:
33         curState[i] -= dif
34         curAmount -= dif
35     if (curAmount > 0):
36         curState[len(curState)-1] -= curAmount
37     return curState
38
39 """This method is for setting up the matrices T and Q, as well as
40     creating the list of states. For each state it makes all
41     possible moves it can,
42     and for each of those, it adds them to their spot in listOfStates
43     . The function newState is used to get the resulting state from
44     a given move.
45     the states are added at the end of the list of states in this
46     case, unless they are already in the list.
47     All those moves are also added to T, which is the main purpose.
48     Then Q is made, and since each state's index in it is reachable
49     all fields
50     are filled with small random numbers.
51     Lastly the method sets up the inputs for the training, and runs
52     it as well. Lastly the method takes the time used both by the
53     setup and the training."""
54 def setup(self, type = 0):
55     timeStart = time.time()
56     emptyBoard = [0] * len(self.listOfStates[0])
57     for i in range(self.sum):
58         self.T.append([0] * self.sum)
59         for x in range(len(self.listOfStates[i])-1, -1, -1):
60             atEnd = 0
61             curIndex = i+1
62             for y in range(1, self.listOfStates[i][x]+1):
63                 newState = self.newState(x, y, self.listOfStates[i])
64                 if(atEnd):
65                     self.T[i][curIndex] = 1
66                     self.listOfStates.append(newState)
67                     curIndex += 1
68                 else:
69                     atEnd = 1
70                     for j in range(curIndex, len(self.listOfStates)):
71                         if (newState == self.listOfStates[j]):
72                             self.T[i][j] = 1
73                             curIndex = j+1
74                             atEnd = 0
75                             break
76                         if (atEnd):
77                             self.T[i][len(self.listOfStates)] = 1
78                             self.listOfStates.append(newState)
79                             curIndex = len(self.listOfStates)
80             if(i == (self.sum-1)):
81                 self.Q.append(100.0)
82             else:
83                 self.Q.append(random.random())

```

```

75     timeEnd = time.time()
76     setupTime = timeEnd-timeStart
77     if(type == 0):
78         epsilon = 0
79         mu = 0.5
80         gamma = 0.5
81         amountEpisodes = 1000
82         timeStart = time.time()
83         self.train(amountEpisodes, epsilon, mu, gamma)
84         timeEnd = time.time()
85         trainTime = timeEnd-timeStart
86
87         times = []
88         times.append(setupTime)
89         times.append(trainTime)
90
91     return times
92     return -1
93
94     """ This is the method for making a move, it finds the state(
95         which is sorted first) in the list of states, by calculating
96         its index. Then the best move is found,
97         By checking which of the accessible positions are the best.
98         Lastly it identifies the index of the heap in which the move is
99             made, and it finds a heap of the same size
100            in the input state."""
101    def makeMove(self, state):
102        curState = state[:]
103        curState.sort(reverse = True)
104        nextState = []
105        move = []
106        if (len(state)<len(self.listOfStates[0])):
107            dif = len(self.listOfStates[0])-len(state)
108            for i in range(dif):
109                curState.append(0)
110        goal = 0
111        progress = self.board[0]
112        for i in range(len(self.board)):
113            heaps = len(self.board)- (i+1)
114            for j in range(progress, -1, -1):
115                if(j == curState[i]):
116                    break
117                else:
118                    if(heaps == 0):
119                        goal += j-curState[i]
120                        break
121                    result = 1
122                    for x in range(1, heaps + 1):
123                        result *= (j + x)
124                    goal += (result)//math.factorial(heaps)
125                    progress -= 1
126        highVal = -np.inf
127        highValIndex = -1
128        for j in range(goal+1, len(self.Q)):
129            if(self.T[goal][j] == 1 and np.greater(self.Q[j], highVal)):
130                highVal = self.Q[j]
131                highValIndex = j

```

```

128     nextState = self.listOfStates[highValIndex][:]
129     numToDec = 0
130     for i in range(len(curState)-1, -1, -1):
131         if (curState[i] != nextState[i]):
132             numToDec += curState[i]-nextState[i]
133             newIndex = curState[i]
134     move.append(numToDec)
135     for i in range(len(state)):
136         if(newIndex == state[i]):
137             newIndex = i
138             break
139     move.append(newIndex)
140     return move
141
142 """This method is for training the matrix Q, by using the Sarsa
143 algorithm. The main difference is that it removes the
144 calculated reward, instead of adding it.
145 It trains by playing against itself, until it wins, and then
146 starts again. each game is an episode."""
147 def train(self, amountEpisodes, epsilon, mu, gamma):
148     for neverUsed in range(amountEpisodes):
149         curState = 0
150         curAction = self.epsilonGreedyChoice(curState, epsilon)
151         playing = 1
152         while(playing):
153             curReward = self.Q[curAction]
154             nextState = curAction
155             nextAction = self.epsilonGreedyChoice(nextState, epsilon)
156             self.Q[curAction] -= mu * (curReward + gamma*self.Q[
157             nextAction] - self.Q[curAction])
158             curState = nextState
159             curAction = nextAction
160             if (curAction == self.sum-1):
161                 playing = 0
162
163 #This is an implementation of the epsilon greedy policy, by the
164 #use of numpy.argmax when finding the best move, or just random
165 #choice of the possible moves otherwise.
166 """This is an implementation of the epsilon greedy policy,
167 however it finds the best move without using np.argmax, as Q
168 doesn't correspond to T in the same way here,
169 as it is only a 1-D array."""
170 def epsilonGreedyChoice(self, state, epsilon):
171     if(random.random()<epsilon):
172         tmpList = []
173         for i in range(state, self.sum):
174             if(self.T[state][i] == 1):
175                 tmpList.append(i)
176         return tmpList[random.randrange(0, len(tmpList))]
177     else:
178         highVal = -np.inf
179         highValIndex = -1
180         for i in range(state+1, len(self.Q)):
181             if(self.T[state][i] == 1 and np.greater(self.Q[i], highVal)
182             ):
183                 highVal = self.Q[i]
184                 highValIndex = i

```

```

176     return highValIndex
178
179     #This method returns the list of states, which is used to test
180     #the program, by having it play from many different states.
181     def getListOfStates(self):
182         return self.listOfStates

```

Listing 33: This is the code in the file *OneDSarsa.py*.

```

1 import MakeData as MD
2 import PerfectPlayer as PP
3 import numpy as np
4 import mlp
5 import sys
6 import time
7
8 if(len(sys.argv) < 5):
9     print("Error: too few arguments. Got " + str(len(sys.argv)-1) + " "
10        , Expected 4")
11 elif (len(sys.argv) > 5):
12     print("Error: too many arguments. Got, " + str(len(sys.argv)-1) + " "
13        , Expected 4")
14 else:
15     #Parsing the input.
16     heapSizeLower = int(sys.argv[1])
17     heapSizeUpper = int(sys.argv[2])
18     heapAmountLower = int(sys.argv[3])
19     heapAmountUpper = int(sys.argv[4])
20     f = open("mlpstatistics.txt", "w")
21     filename = "CorrectMoves.txt"
22     player = PP.PerfectPlayer()
23     #The loop that defines the start-states from which the programs
24     #will be run.
25     for heapSize in range(heapSizeLower, heapSizeUpper+1):
26         for heapAmount in range(heapAmountLower, heapAmountUpper+1):
27             """The Supervised learning program is run 10 times for each
28             start-state, and plays against the perfect player each time,
29             to get a winrate for statistics. The time used is also
30             written to file as statistics. The data used is the same for a
31             given state,
32             but it is sorted differently for each time."""
33             state = [heapSize] * heapAmount
34             print(state)
35             md = MD.MakeData()
36             md.setup(state)
37             timeAvg = 0
38             percent = 0
39             for i in range(10):
40                 #sorting and separating the data into the needed categories
41
42                 data = np.loadtxt(filename)
43
44                 stateLen = len(data[0])
45
46                 boards = data[:,::2]
47                 moves = data[:,1::2]
48
49                 order = list(range(np.shape(boards)[0]))

```

```

43     np.random.shuffle(order)
44     boards = boards[order,:]
45     moves = moves[order,:]

46     train = boards[:,2]
47     trainTargets = moves[:,2]

48     valid = boards[:,4]
49     validTargets = moves[:,4]

50     test = boards[:,6]
51     testTargets = moves[:,6]

52     hidden = len(state)+10

53
54
55     start = time.time()
56     net = mlp.mlp(hidden, stateLen)
57     net.earlystopping(train, trainTargets, valid, validTargets)
58     end = time.time()
59     timeAvg += end-start
60     listOfStates = md.getStates()
61     wins = 0
62     losses = 0
63
64     for i in range(len(listOfStates)):
65         if(player.findNimSum(listOfStates[i]) != 0):
66             playing = 1
67             board = listOfStates[i][:]
68             winBoard = [0] * len(listOfStates[i])
69             while(playing):
70                 move = net.makeMove(board)
71                 board[move[1]] = board[move[1]]-move[0]
72                 if(board == winBoard):
73                     wins += 1
74                     break
75                 move = player.makeMove(board)
76                 board[move[1]] = board[move[1]]-move[0]
77                 if(board == winBoard):
78                     losses += 1
79                     break
80             percent += (wins/(wins+losses)*100)
81             timeAvg /= 10
82             percent /= 10
83             line = "Time used:" + str(round(timeAvg, 3)) + "s Winrate/"
84             sucessrate: " + str(int(percent)) + "%"
85             line += str(state) + "\n"
86             f.write(line)
87

```

Listing 34: This is the code in the file *RunMlp.py*.

```

1 import PerfectPlayer as PP

3 class MakeData:
4     #This method just makes the list of states accessible from the
5     #start-state, and then runs the writeToFile method.
6     def setup(self, startBoard):
7         sum = startBoard[0]+1
8         board = startBoard
9         for i in range(1, len(board)):

```

```

9     sum *= board[i]+1
10    self.listOfStates = [0] * sum
11    self.listOfStates[0] = board[:]
12    for i in range(sum):
13        for x in range(len(self.listOfStates[i])):
14            for y in range(1, self.listOfStates[i][x] + 1):
15                tmpValue = self newIndex(x, y, board)
16                if(self.listOfStates[i+tmpValue] == 0):
17                    tmpList = self.listOfStates[i][:]
18                    tmpList[x] -= y
19                    self.listOfStates[i+tmpValue] = tmpList[:]
20    self.writeToFile(self.listOfStates)
21
22 # This method exists for finding, given a move, how far down the
23 # list of states the new state is.
24 def newIndex(self, heapIndex, amount, board):
25     newMove = 1
26     for i in range(len(board)-1, heapIndex, -1):
27         newMove *= board[i]+1
28     newMove *= amount
29     return newMove
30
31 #This file writes each state in the list to the file, and a
32 #correct move to make when in that state.
33 def writeToFile(self, states):
34     player = PP.PerfectPlayer()
35     f = open("CorrectMoves.txt", "w")
36     for i in range(len(states)):
37         if (player.findNimSum(states[i]) != 0):
38             tmp = [0]*len(states[0])
39             result = player.makeMove(states[i])
40             tmp[result[1]] = result[0]
41             for j in range(len(states[i])):
42                 f.write(str(states[i][j]) + " ")
43             f.write("\n")
44             for j in range(len(states[i])):
45                 f.write(str(tmp[j]) + " ")
46             f.write("\n")
47
48 #This method returns the list of states, which is used to test
49 #programs, by having them play from many different states.
50 def getStates(self):
51     return self.listOfStates

```

Listing 35: This is the code in the file *MakeData.py*.

```

1 import numpy as np
2 import random
3 from scipy.special import expit
4
5 class mlp:
6     def __init__(self, nhidden, stateLen):
7         self.beta = 1
8         self.eta = 0.1
9         self.momentum = 0.0
10        self.matrix1 = []
11        self.matrix2 = []
12        self.inputAmount = stateLen

```

```

13     self.outputAmount = stateLen
14     self.hiddenAmount = nhidden
15     self.outputA = [0.0] * (self.outputAmount)
16     self.outputList = [0.0] * (self.outputAmount)
17     self.hiddenU = [0.0] * (nhidden)
18     self.hiddenList = [0.0] * (nhidden)
19     for i in range(self.inputAmount + 1):
20         self.matrix1.append([])
21         for j in range(nhidden):
22             self.matrix1[i].append(random.random())
23     for i in range(nhidden + 1):
24         self.matrix2.append([])
25         for j in range(self.outputAmount):
26             self.matrix2[i].append(random.random())
27
28 def sigmoidDerived(self, x):
29     return expit(x)*(1-expit(x))
30
31 def identityFunction(self, x):
32     return x
33
34 def identityDerived(self, x):
35     return 1
36
37 #this function calculates the sum of squares error of all the
38 #valid inputs and targets, and sums them, and returns that for
39 #the earlystopping to use for figuring out when to stop.
40 def totalError(self, inputs, targets):
41     accumError = 0.0
42     for i in range(len(inputs)):
43         accum = 0.0
44         self.forward(inputs[i])
45         for j in range(self.outputAmount):
46             accum += (targets[i][j] - self.outputList[j])**2
47         accumError += accum/2
48     return accumError
49
50 #earlystopping, stops when error increases, or decreases too
51 #slowly
52 def earlystopping(self, inputs, targets, valid, validtargets):
53     stopTraining = 0
54     lowestError = self.totalError(valid, validtargets)
55     while not stopTraining:
56         self.train(inputs, targets, 10)
57         newError = self.totalError(valid, validtargets)
58         if newError > lowestError or lowestError-newError < 0.5:
59             stopTraining = 1
60         else:
61             lowestError = newError
62
63 #train trains the neural net for the inputted amount of
64 #iterations. One pass through the dataset is an iteration.
65 def train(self, inputs, targets, iterations=100):
66     for notUsed in range(iterations):
67         for ite in range(len(inputs)):
68             self.forward(inputs[ite])
69             outDelta = []

```

```

66     for i in range(self.outputAmount):
67         outDelta.append((self.outputList[i] - targets[ite][i]) *
68             self.identityDerived(self.outputA[i]))
69         hiddenDelta = []
70         for i in range(self.hiddenAmount):
71             accumulator = 0.0
72             for j in range(self.outputAmount):
73                 accumulator += outDelta[j] * self.matrix2[i+1][j]
74             hiddenDelta.append(self.sigmoidDerived(self.hiddenU[i]) *
75                 accumulator)
76             for j in range(self.outputAmount):
77                 self.matrix2[0][j] -= self.eta * outDelta[j] * 1
78             for i in range(self.hiddenAmount):
79                 self.matrix2[i+1][j] -= self.eta * outDelta[j] * self.
80             hiddenList[i]
81             for j in range(self.hiddenAmount):
82                 self.matrix1[0][j] -= self.eta * hiddenDelta[j] * 1
83             for i in range(self.inputAmount):
84                 self.matrix1[i+1][j] -= self.eta * hiddenDelta[j] *
85             inputs[ite][i]

86     #forward passes the inputs through the neural net to get the
87     #outputs.
88     def forward(self, inputs):
89         for i in range(len(self.hiddenList)):
90             self.hiddenU[i] = 1*self.matrix1[0][i]
91             for j in range(self.inputAmount):
92                 self.hiddenU[i] += inputs[j] * self.matrix1[j+1][i]
93             self.hiddenList[i] = expit(self.hiddenU[i])
94         for i in range(len(self.outputList)):
95             self.outputA[i] = 1*self.matrix2[0][i]
96             for j in range(len(self.hiddenList)):
97                 self.outputA[i] += self.hiddenList[j]*self.matrix2[j+1][i]
98             self.outputList[i] = self.identityFunction(self.outputA[i])

99     """To find the best move, the input is run forward through the
100    neural net. The highest output is considered the chosen move.
101    If the chosen output is illegal it is trimmed, to make it legal,
102    either making it bigger or smaller. if the index is of a hap of
103    0,
104    a removal of 1 in some non-empty heap is chosen."""
105    def makeMove(self, state):
106        self.forward(state)
107        if(state[0] < -1):
108            1/0
109        move = [0,0]
110        move[1] = np.argmax(self.outputList)
111        move[0] = int(self.outputList[move[1]])
112        if(move[0] < 1):
113            move[0] = 1
114        elif(move[0] > state[move[1]]):
115            move[0] = state[move[1]]
116        if(state[move[1]] == 0):
117            for i in range(len(state)):
118                if(state[i] > 0):
119                    move[0] = 1
120                    move[1] = i

```

```
115     return move
```

Listing 36: This is the code in the file *mlp.py*.

```
1 import numpy as np
2 import random
3 import math
4 import time

6 class SpaceSaveSarsaSpread:

8     #The initialization takes the starting board, and finds the
9     #amount of states in its untrimmed game-tree, which is called
10    self.sum. It also initializes T and Q.
11    def __init__(self, startBoard):
12        self.board = startBoard
13        self.sum = self.getSum()
14        self.board.sort(reverse = True)
15        self.listOfStates = []
16        self.listOfStates.append(self.board[:])
17        self.T = []
18        self.Q = []

19    #This method is for finding the amount of states reachable from
20    #the start-state, for the trimmed game-tree.
21    def getSum(self):
22        return int(math.factorial(len(self.board)+ self.board[0])/(math
23            .factorial(len(self.board))*math.factorial(self.board[0])))

24    #This method is for taking a state and a move in the state, and
25    #returning the resulting state in the correct order.
26    def newState(self, heapIndex, amount, state):
27        curState = state[:]
28        curAmount = amount
29        for i in range(heapIndex, len(state)-1):
30            dif = curState[i] - curState[i+1]
31            if (dif >= curAmount):
32                curState[i] -= curAmount
33                curAmount = 0
34                break
35            else:
36                curState[i] -= dif
37                curAmount -= dif
38            if (curAmount > 0):
39                curState[len(curState)-1] -= curAmount
40        return curState

41    """This method is for setting up the matrices T and Q, as well as
42    creating the list of states. For each state it makes all
43    possible moves it can,
44    and for each of those, it adds them to their spot in listOfStates
45    . The function newState is used to get the resulting state from
46    a given move.
47    the states are added at the end of the list of states in this
48    case, unless they are already in the list.
49    All those moves are also added to T, which is the main purpose.
50    Then Q is made, by setting all positions in T with 1, to have
51    small random values in Q.
```

```

43     The other positions in Q gets the value of negative infinity, -np
44         .inf, so that they are always smaller than the expected reward
45         of a move.
46     Lastly the method sets up the inputs for the training, and runs
47         it as well. Lastly the method takes the time used both by the
48         setup and the training."""
49
50     def setup(self, type = 0):
51         timeStart = time.time()
52         emptyBoard = [0] * len(self.listOfStates[0])
53         for i in range(self.sum):
54             self.T.append([0] * self.sum)
55             for x in range(len(self.listOfStates[i])-1, -1, -1):
56                 atEnd = 0
57                 curIndex = i+1
58                 for y in range(1, self.listOfStates[i][x]+1):
59                     newState = self.newState(x, y, self.listOfStates[i])
60                     if(atEnd):
61                         self.T[i][curIndex] = 1
62                         self.listOfStates.append(newState)
63                         curIndex += 1
64                     else:
65                         atEnd = 1
66                         for j in range(curIndex, len(self.listOfStates)):
67                             if (newState == self.listOfStates[j]):
68                                 self.T[i][j] = 1
69                                 curIndex = j+1
70                                 atEnd = 0
71                                 break
72                             if (atEnd):
73                                 self.T[i][len(self.listOfStates)] = 1
74                                 self.listOfStates.append(newState)
75                                 curIndex = len(self.listOfStates)
76                                 self.Q.append([])
77                                 for j in range(self.sum-1):
78                                     if(self.T[i][j] == 1):
79                                         self.Q[i].append(random.random())
80                                     else:
81                                         self.Q[i].append(-np.inf)
82                                     if (self.T[i][self.sum-1] == 1):
83                                         self.Q[i].append(100.0)
84                                     else:
85                                         self.Q[i].append(-np.inf)
86             timeEnd = time.time()
87             setupTime = timeEnd-timeStart
88             if(type == 0):
89                 epsilon = 0
90                 mu = 0.5
91                 gamma = 0.5
92                 amountEpisodes = 1000
93                 timeStart = time.time()
94                 self.train(amountEpisodes, epsilon, mu, gamma)
95                 timeEnd = time.time()
96                 trainTime = timeEnd-timeStart
97
98                 times = []
99                 times.append(setupTime)
100                times.append(trainTime)

```

```

96
97         return times
98     return -1

99
100    """ This is the method for making a move, it finds the state(
101        which is sorted first) in the list of states, by calculating
102        its index. Then the best move is found by use
103        of the function np.argmax. Lastly it identifies the index of the
104        heap in which the move is made, and it finds a heap of the same
105        size
106    in the input state."""
107
108    def makeMove(self, state):
109        curState = state[:]
110        curState.sort(reverse = True)
111        nextState = []
112        move = []
113
114        if (len(state)<len(self.listOfStates[0])):
115            dif = len(self.listOfStates[0])-len(state)
116            for i in range(dif):
117                curState.append(0)
118
119        goal = 0
120        progress = self.board[0]
121
122        for i in range(len(self.board)):
123            heaps = len(self.board)- (i+1)
124            for j in range(progress, -1, -1):
125                if(j == curState[i]):
126                    break
127                else:
128                    if(heaps == 0):
129                        goal += j-curState[i]
130                        break
131                    result = 1
132                    for x in range(1, heaps + 1):
133                        result *= (j + x)
134                    goal += (result)//math.factorial(heaps)
135                    progress -= 1
136
137        nextState = self.listOfStates[np.argmax(self.Q[goal])][:]
138        numToDec = 0
139
140        for i in range(len(curState)-1, -1, -1):
141            if (curState[i] != nextState[i]):
142                numToDec += curState[i]-nextState[i]
143                newIndex = curState[i]
144
145                move.append(numToDec)
146
147                for i in range(len(state)):
148                    if(newIndex == state[i]):
149                        newIndex = i
150                        break
151
152                move.append(newIndex)
153
154    return move

155
156    """This method is for training the matrix Q, by using the Sarsa
157        algorithm, as well as training not just from the start states,
158        but from all states(except the win-state)"""
159
160    def train(self, amountEpisodes, epsilon, mu, gamma):
161        curStartState = 0
162
163        for neverUsed in range(amountEpisodes):
164            curState = curStartState

```

```

148     curAction = self.epsilonGreedyChoice(curState, epsilon)
149     playing = 1
150     while(playing):
151         curReward = self.Q[curState][curAction]
152         nextState = curAction
153         if(nextState == len(self.listOfStates)-1):
154             break
155         nextAction = self.epsilonGreedyChoice(nextState, epsilon)
156         self.Q[curState][curAction] -= mu * (curReward + gamma*self
157         .Q[nextState][nextAction] - self.Q[curState][curAction])
158         curState = nextState
159         curAction = nextAction
160         if (nextAction == self.sum-1):
161             playing = 0
162         curStartState += 1
163         if(curStartState >= len(self.listOfStates)-2):
164             curStartState = 0
165
166         #This is an implementation of the epsilon greedy policy, by the
167         #use of numpy.argmax when finding the best move, or just random
168         #choice of the possible moves otherwise.
169     def epsilonGreedyChoice(self, state, epsilon):
170         if(random.random()<epsilon):
171             tmpList = []
172             for i in range(state, self.sum):
173                 if(self.T[state][i] == 1):
174                     tmpList.append(i)
175             return tmpList[random.randrange(0, len(tmpList))]
176         else:
177             return np.argmax(self.Q[state])
178
179         #This method returns the list of states, which is used to test
180         #the program, by having it play from many different states.
181     def getListOfStates(self):
182         return self.listOfStates

```

Listing 37: This is the code in the file *SpaceSaveSarsaSpread.py*.

References

- Bouton, Charles L. “Nim, A Game with a Complete Mathematical Theory”. In: *Annals of Mathematics* 3.1/4 (1901), pp. 35–39. ISSN: 0003486X. URL: <http://www.jstor.org/stable/1967631>.
- Conway, John H., Elwyn R. Berlekamp, and Richard K. Guy. *Winning Ways for Your Mathematical Plays*. Second. Vol. 1. CRC Press, 2001.
- Marsland, Stephen. *Machine learning: an algorithmic perspective*. Second. CRC Press/Taylor & Francis Group, 2015.
- The Ferranti NIMROD Digital Computer*. URL: <https://www.goodeveca.net/nimrod>.

List of Figures

1	A model of a McCulloch Pitts neuron.	10
2	A model of a single layer Perceptron. Note that only the darker nodes are actual neurons, the lighter ones are there to show how the inputs are given to all neurons.	11
3	A model of a MultiLayer Perceptron.	12
4	First graph is the graph for the game tree of a game of Nim, starting with the state 2 2. Notice how many times state 0 0 (blue) occurs. State 1 1 (red) occurs twice as well. The second graph is the resulting graph if all equal state-nodes in the game tree is combined. Notice how the states 0 0, and 1 1 occur just once in this graph.	17
5	Graph to the left is the untrimmed graph, the one to the right the trimmed one. See how the trimmed one has almost been cut in two? It is much smaller, another improvement to the graph. . .	19
6	A simple model of a Multilayer Perceptron.	22
7	The equation for finding the index of a state, in relation to the start state, for a specific move. n is the size of the heaps in the starting state, x is the amount of counters being removed, and k is the amount of piles right if the pile being removed, which would be $\text{length}(\text{state}) - k'$ where k' is the pile being removed from.	33

List of Tables

1	A Nim-sum example	15
2	This table shows the decision process for each digit in the nim-sum operation.	15
3	Nim-sum of the state 3 5 7	15
4	Nim-sum of Table 3s Nim-sum and the first heap in the game-state. .	16
5	The graph as a matrix	18
6	The shortform names for the different programs used in the statistics.	53
7	The general statistics of the different algorithms.	54
8	The statistics for applying sarsa to the untrimmed game tree. . .	55
9	The statistics for applying Q-learning to the untrimmed game tree. .	56
10	The statistics for applying sarsa to the trimmed game tree. . . .	57
11	The statistics for applying Q-learning to the trimmed game tree. .	58
12	The statistics for using the program using the one-dimensional array, and sarsa.	59
13	The lines which show a total time use that is more than the sum of the setup and train times.	60
14	The strange phenomena, success rate is higher where there are an even number of heaps. From running TS.	62

15	Statistics for trimmed sarsa, with ϵ set to 1	63
16	Comparing the Algorithm that trains from all states with Trimmed Sarsa.	64
17	The timeuse and winrate for the different algorithms when given the gamestate: 5 5 5 5	65
18	How the timeuse develops for greater start-states	65
19	The results of stochastic termination of TS	66
20	Time-use of moves	67
21	Statistics of a run of the program for playing nim with Supervised learning.	68

Listings

1	Algorithm for finding the correct move from a certain state in Nim	23
2	Initialization of the programs using the untrimmed game-tree . .	24
3	Initialization of the programs using the trimmed game-tree . . .	25
4	A function for applying Equation 15 to find the amount of states reachable from the start state.	25
5	The setup function for the untrimmed programs.	26
6	The setup function for the trimmed programs.	28
7	The part of the setup function which is different in the One-dimensional program.	29
8	The code for the Q-learning algorithm. Note that linenumbers are from the untrimmed program.	30
9	The code for the Sarsa algorithm. Note that linenumbers are from the untrimmed program.	30
10	The code for the Sarsa algorithm, in the One-dimensional program	31
11	The <i>makeMove</i> function in the untrimmed programs.	32
12	The <i>makeMove</i> function in the trimmed programs.	33
13	The <i>makeMove</i> function in the one-dimensional program.	34
14	The code in <i>Run.py</i> before the programs are actually run.	35
15	This is the code that runs the programs, takes the time they use and creates statistics about it.	36
16	The code in program <i>Play.py</i>	38
17	The program <i>MakeData.py</i> which creates the data used by the supervised learning algorithm.	40
18	This is the part of the program <i>mlp.py</i> in which the training happens.	42
19	The part of <i>mlp.py</i> that is used to get a move from the program.	44
20	The code for the program <i>RunMlp.py</i>	45
21	The code for the program <i>Episodes.py</i> , that tries to run the reinforcement programs with stochastic termination	47
22	Total time is taken here. From file <i>Run.py</i>	60
23	Setup and training time is taken here. From file <i>ReinforcementSarsa.py</i> , but is the same for all the similar programs.	60

24	The initialization of the file, which is part of the what is run in the total time only. From file <i>ReinforcementSarsa.py</i> , but is the same for all the similar programs.	61
25	The method <i>train</i> in the program that trains from all the states; <i>SpaceSaveSarsaSpread.py</i>	63
26	This is the code in the file <i>PerfectPlayer.py</i>	70
27	This is the code in the file <i>Run.py</i>	70
28	This is the code in the file <i>Play.py</i>	75
29	This is the code in the file <i>ReinforcementSarsa.py</i>	78
30	This is the code in the file <i>ReinforcementQlearning.py</i>	80
31	This is the code in the file <i>SpaceSaveSarsa.py</i>	83
32	This is the code in the file <i>SpaceSaveQlearning.py</i>	86
33	This is the code in the file <i>OneDSarsa.py</i>	90
34	This is the code in the file <i>RunMlp.py</i>	94
35	This is the code in the file <i>MakeData.py</i>	95
36	This is the code in the file <i>mlp.py</i>	96
37	This is the code in the file <i>SpaceSaveSarsaSpread.py</i>	99