

Reinforcement Learning for Blackjack

Saqib A Kakvi
Supervisor: Dr Marco Gillies

April 30, 2009

Abstract

This report explores the development of an Artificial Intelligence system for an already existing framework of card games, called SKCards. This system was developed by the author as a way of learning Java. The old Artificial intelligence is highly flawed.

Reinforcement Learning was chosen as the method to be employed. Reinforcement Learning attempts to teach a computer certain actions, given certain states, based on past experience and numerical rewards gained. The agent either assigns values to states, or actions in states.

This will initially be developed for Blackjack, with possible extensions to other games. Blackjack is one of the simpler games and the only game in the SKCards package which needs an Artificial Intelligence agent. All the other games are single player.

Acknowledgements

I am hugely grateful to my project supervisor, Dr. Marco Gillies, and all the members of staff, for all the support and assistance they gave me. I would also like to thank my colleagues, whose help was invaluable. In particular I would like to thank David Jarret, Rohan Sherrard, Daniel Smith, Benoy Sen, Jusna Begum, Arron Donnely & Josephene Amirthanayagam. And a very special thank you to Nidhi Shah and Fatima Nassir for their help in proofreading this report.

Contents

1	Introduction	1
1.1	The problem I am solving	1
1.1.1	Reinforcement Learning Agent	1
1.1.2	SKCards	1
1.2	The layout of the report	1
2	Background	3
2.1	Blackjack	3
2.1.1	Overview	3
2.1.2	Game play	3
2.2	Reinforcement Learning	5
2.2.1	Introduction	5
2.2.2	Rewards	5
2.2.3	Policies	6
2.2.4	Agent Behaviour	6
2.3	Blackjack in Reinforcement Learning	7
3	Project Description	9
3.1	Technologies involved	9
3.2	System Design	9
4	System Development	15
4.1	Prototyping	15
4.1.1	Prototype 1 - Next card higher/lower	15
4.1.2	Prototype 2 - GridWorld	15
4.2	Final System	16
4.2.1	AI	16
4.2.2	SKCards	20
5	System Evaluation	21
5.1	System testing	21
5.2	Results	22
5.2.1	Table	22
5.2.2	Graphs	23
5.3	Result Analysis	30
5.4	Possible Future Improvements	30
5.4.1	Card Counting	30
5.4.2	Probabilities	31
5.4.3	Busting	31
5.4.4	SKCards	31

6	Summary and Conclusions	33
6.1	Summary of the Project	33
6.1.1	Goals	33
6.1.2	Outcomes	33
6.2	Conclusion	33
7	Project Self-Evaluation	35
7.1	Introduction	35
7.2	Project outcomes	35
7.3	Project planning	35
7.4	Self-evaluation	35
A	Card High/Low prototype	37
A.1	Source Code	37
A.2	Results	42
B	GridWorld prototype	59
B.1	Source Code	59
B.2	Results	67
C	Blackjack AI	83
C.1	Source Code	83
C.1.1	Face	83
C.1.2	Suit	83
C.1.3	Card	84
C.1.4	Pack	89
C.1.5	Blackjack	93
D	Final System	109
D.1	Source Code	109
D.1.1	Face	109
D.1.2	Suit	109
D.1.3	Card	109
D.1.4	Pack	109
D.1.5	Shoe	109
D.1.6	BJPlayer	110
D.1.7	BJUI	114
D.1.8	BJRules	119
D.1.9	RLPlayer	122
E	Sample Blackjack Game	125

List of Figures

3.1	Old System Architecture	11
3.2	New System Architecture	12
3.3	Use Cases	13
3.4	Single Episode	14
4.1	The frequency of the extended actions	17
5.1	Q-Learning Policies	25
5.2	Win Percentages	26
5.3	Loss Percentages	27
5.4	Bust Percentages	28
5.5	Net Win Percentages	29
A.1	Values for Joker	43
A.2	Values for Ace	44
A.3	Values for 2	45
A.4	Values for 3	46
A.5	Values for 4	47
A.6	Values for 5	48
A.7	Values for 6	49
A.8	Values for 7	50
A.9	Values for 8	51
A.10	Values for 9	52
A.11	Values for 10	53
A.12	Values for Jack	54
A.13	Values for Queen	55
A.14	Values for King	56
B.1	Episode 1	69
B.2	Episode 25	70
B.3	Episode 50	71
B.4	Episode 100	72
B.5	Episode 150	73
B.6	Episode 200	74
B.7	Episode 250	75
B.8	Episode 500	76
B.9	Episode 750	77
B.10	Episode 1000	78
B.11	Episode 1500	79
B.12	Episode 2000	80
B.13	Episode 2500	81

List of Tables

5.1	Final Results	23
A.1	Final Policy	57
B.1	Predicted values for the GridWorld problem	68

Chapter 1

Introduction

1.1 The problem I am solving

This project has two goals. The major goal is the development of a Reinforcement Learning agent for Blackjack, using parts of an existing system. The second and less important goal is to improve the design of the existing system, so as to comply with software engineering norms. These goals may be solved independently and then combined to make the final system.

1.1.1 Reinforcement Learning Agent

When the initial system was built, it lacked a proper Artificial Intelligence Agent for the player to play against. This project aims to build an AI agent using Reinforcement Learning for Blackjack. This agent can be the dealer, or it can be another player in the game. Reinforcement Learning is very suited for the task of Blackjack due to the easily represented states. And the agent could learn to match a user's style of play, over time.

1.1.2 SKCards

SKCards was built as a framework for all card games. It provides the basic Cards, Packs of Cards and a standard interface. The concept is that any programmer could create their own card game and easily integrate it into SKCards. This, would give us an infinitely extendable suite of card games. This was created as an experiment with the Java programming language.

1.2 The layout of the report

This report will aim to break down the problem and present it in several smaller pieces. Each piece will logically follow from the next and will try to build on top of the last. In this way this report should give a smooth transition from theory to implementation and finally the results and conclusions.

The next chapter will discuss the separate parts of the project, namely, Blackjack and Reinforcement Learning. Then it discusses the two together. The following chapter will discuss the design aspects of the system. After that the development of the system will be looked at. This will detail the implementation of the system. Following that will be the evaluation of the system. This will detail the experiments, the results and an analysis of these results. Finally, the last chapter will cover the self-evaluation.

Chapter 2

Background

2.1 Blackjack

2.1.1 Overview

Blackjack is a very simple game. The object is to get as close to, but not over 21 points in a maximum of 5 cards. This is achieved by assigning values to each of the cards in the deck, irrespective of suit, as listed below:

2 - 10 \Rightarrow Face value (i.e. a 2 is 2 points, a 3 is 3 points)

Courts (Jack, Queen, King) \Rightarrow 10 points

A \Rightarrow 1 or 11 points. The player decides if he wants the ace “high” (11) or “low” (1). This can be changed in a situation where a “high” ace would cause a player to be bust.

2.1.2 Game play

The player is dealt 2 cards face down. The dealer then deals themselves 2 cards, 1st face up then face down. The player places a bet on his hand and turns over his cards. The player may then take any of several actions, listed below. These have been divided into core and extended rules, where the core rules make for the basis of the game, and the extended rules give additional gameplay.

Core Rules

HIT: The player may ask for another card from the pack to be added to their hand, this is known as a “hit” or “twist” (the term “hit” shall be used henceforth). The dealer “hits” the player’s hand. The player’s score is increased by the point value of the card added. If the player’s score exceeds 21, they are “bust” and loses on that hand. If not, they can continue “hitting” until they are bust or have 5 cards in their hand.

STAND: If the player is satisfied with their hand they may “stand”. This means that he is done with the current hand. The player then moves on to play any other hands they have (see Split in section 2.1.2). After “standing” on all their hands, the player ends their turn.

Extended rules

These rules are extensions of the rules in the previous section. They provide further functionality, however they are mostly based on the core rules.

BLACKJACK: The quickest way to win. A blackjack is an Ace and a Court Card (King, Queen or Jack) bringing the player up to a total of 21, known as a “natural”. The player receives 2.5 times the money he was supposed to receive. This situation does not apply if the player has an Ace and a Ten. This rule was introduced in casinos to increase popularity of the game. Originally it only applied to an Ace and the Jack of Clubs, the “Black Jack”, which is where the game derives its name from.

INSURANCE: If the dealer's exposed card is an ace, threatening a Blackjack, the player may take out "insurance" against a Blackjack. They place another bet on the table. If the dealer has Blackjack, they get a 2-to-1 payout. If the dealer doesn't have a Blackjack, the player loses that money. As the name states, this serves as the player's backup, should they lose to a dealer blackjack.

SPLIT: If both the cards have the same face value (two Queens, two Fives, two Sevens...) the player may opt to "split" their hand into two separate hands. The cards are separated and one card is added to each hand to make two cards in each hand. The player then places the same bet on the second hand. A player may have up to 5 hands.

DOUBLE DOWN: If the total points of the player's cards is greater than or equal to 11 the player may opt to "double down". In this case the player doubles their bet and the player's hand is hit, after which they are forced to stand, provided the player is not bust. This rule is basically a gamble by the player that they will need no more than 1 Card to reach a favourable score. This increased risk is why there is an increased reward.

After the player is finished, the next player plays in the same manner. After all players have finished playing, play moves to the dealer, who plays in a similar manner. After the dealer plays, the winner is the hand which is closest to 21. It must be noticed that the dealer generally has some restrictions imposed on them to make it fair on the players.

Dealer restrictions normally come in the form of a hitting threshold, a score above which they can no longer hit. It must be noted that a dealer can hit on high counting hands that are "soft" i.e. they have an Ace counting high. The standard policy in casinos tends to be a threshold score of 17.

Rule variations

The rules listed above are mostly in effect in the United States of America. There are some variations, specifically on the British version of the game, known as Pontoon. The details of this game are described in Parlett (1994). Although the basic principles of the game are the same, there are several rule variations. These variations are described below.

In this game, the role of dealer rotates between the players. The other players are called punters. The dealer also places a stake and plays as normal. Listed below are some of the other rule variations:

BUY: This is a variation on the hit rule. Under this rule, the player may request for a card to be dealt face-down to them. This allows the player to keep their score secret. To do this the player must pay an amount greater than what they paid for the last card, where applicable, and a less than their total stake. The fifth card of the hand may not be bought.

5-CARD TRICK: If a player has a hand of 5 cards and has not gone bust, then they automatically win. They receive a payment equivalent double to their stake.

PONTOON: In Pontoon Tens, Jacks, Queens & Kings are called tenths. A regular Pontoon is any tenth and an Ace. A Royal Pontoon is three sevens. The payouts for a Pontoon and Royal Pontoon are 2 times and 3 times respectively. Royal Pontoons can only be claimed by a punter. If the dealer holds a pontoon, then each punter pays the dealer double their stake. If the punter holds a Pontoon, they only pay an amount equal to their stake.

BANKER'S PLAY: The banker plays in a similar way as the punters, however there are some restrictions on their play. The banker is not allowed to split pairs. The players may only beat the dealer if they hold a score of 16 or higher. After the dealer has stood and is not bust, at a score of $s < 21$, they make a call of "pay $s + 1$ ". Any players holding hands scored at $s + 1$ or greater must reveal their hands and claim payment. A Royal Pontoon is unbeatable.

A sample game which illustrates most of the gameplay and the basic and extended rules can be found in Appendix E.

2.2 Reinforcement Learning

2.2.1 Introduction

Reinforcement learning is an unsupervised Artificial Intelligence learning method in which the agent teaches itself. The central principle is based on rewards for a sequence of actions taken or sequence of state transitions. Given a state S , the agent will take an action A , based on its knowledge of the environment. This will lead to a change in the environment, and eventually a reward, either positive or negative, will be given. The agent will then take this into consideration and accordingly adjust its knowledge. It must be noted that rewards are not given after every action, but may be after a series of actions.

The agent has a preference to each possible action for a given state. This is based on what is known as the action's value. This is a numeric value attached to each action A , given a state S . Every reward will adjust the action's value and hence make it more or less desirable. There are several ways to value action and the choice is mainly arbitrary. It depends mainly on the developers understanding of the system and their preference as to how the agent should behave.

Reinforcement Learning agents can act based on 3 things: values of states, values of actions or values of actions in states. The choice of which method to employ is left to the developer. This choice is made based on the developer's understanding of the environment and his ability to represent it accurately. It must be noted that an agent can learn in more than one way from the same environment. Action valuing is only applicable in very trivial tasks and hence is not used as much as state or state-action valuing.

The agent learns based on 2 things, namely the current reward and the previous rewards for A , S , or A & S . To simplify calculations and storage, we simply take the average of the previous rewards and update it. As the rewards are arbitrary, although mostly people use +1, -1 and 0, the agent will not know how good a reward is. By comparing it to the average of the past rewards, it can judge how high or low the reward was. A "high" reward would increase the value of the action and a "low" reward would reduce it.

How exactly rewards are allocated is a bit of a mathematical quandary. Some methods allow rewards to be "passed back" to previous actions. This follows directly from the fact that rewards are allocated after a series of actions. This is so a low valued action which can be followed by a high valued action seems more attractive. If this were not done then the action values would not be balanced and would only reflect immediate rewards.

2.2.2 Rewards

The basic value method is based on discounted expected future returns. As mentioned above, when an agent receives a reward, it passes back some of that reward. However this is not always back all the way and may just be to the previous state. By doing this, when the agent next encounters the newly rewarded state, it passes back to the previous state and so on and so forth. Based on these values, and the probability of the states occurring, we calculate the expected reward. So in effect we are summing the rewards of the given possible states over their probability distribution, from $0 \rightarrow \infty$.

This is not always possible to do, due the large problem space. Instead whenever a reward is received, all previous states or state-action pairs are updated by a fraction of that reward. This fraction is a constant call the discount rate and is denoted by γ . It follows that a state $S1$ n steps from state $S2$, which receives a reward of r , will get a reward of $\gamma^n \times r$, where $0 < \gamma < 1$.

It is very important, especially for stochastic tasks, to minimise the effect noise will have on the learning algorithm. To do this, we use what is call the step-size parameter, denoted by α . This means that only a small portion of any reward will be passed back, mitigating any effects of noise in the data. valid data

will be observed several times, which will eventually add up to its true value. The value of the step-size parameter is $0 < \alpha < 1$.

2.2.3 Policies

By assigning rewards to states or state-action pairs, the agent learns what is known as a policy, denoted by P . The policy is the set of actions the agent will take in a given state. The aim of Reinforcement Learning is to learn a policy which is optimal or as close to the optimal policy as possible, denoted by P^* . For some problems, such as the grid world, the optimal policy is defined. However for stochastic and/or non-stationary problems, this is much harder and some other statistic must be used to determine if the agent has learnt the policy it was required to.

Most Reinforcement Learning techniques have been mathematically proven to converge on P^* for the limit $n \rightarrow \infty$, where n is the number of runs. However as it is not always possible to learn the optimal policy in a reasonable amount of time, so the focus on is to learn a policy that is as close as possible to P^* , or not exactly P^* .

Reinforcement learning is a very powerful Artificial Intelligence tool. This very closely mimics the way humans learn by association. For every set of actions taken, there is a reward and this is associated to those actions, relative to past experience. However not every action is rewarded independently, but instead the entire series of actions receives a reward, which is mathematically calculated based on the whole reward. Hence it is not as powerful as human association, but it is none the less, very close to it.

What makes it even more powerful is that it is an unsupervised technique. So the agent can run through thousands of plays on its own and learn some policy P . With supervised learning, the user would be required to tell the agent something after every play. This way it saves time and completely eliminates any human error in the learning process. Assuming the implementation is done correctly, as stated above, it will, in the limit $n \rightarrow \infty$, converge on P^* .

2.2.4 Agent Behaviour

Agents can behave in several ways. The primary question with Reinforcement Learning is "exploitation vs. exploration." An agent can either exploit what it believes to be the best action, or it can explore new options, or ideally it can do some of both. The real trick is striking the correct balance between the two factors. Exploration is best for learning, as the agent will cover as much of the state space as possible, whereas exploitation is very useful after the agent has learnt a good policy.

The simplest exploitative and most intuitive behaviour is called a greedy agent. This agent will always choose the action with the highest reward value, known as exploitation. In this way it is guaranteed large immediate rewards, but possibly lower rewards in the long run. I believe given a small number of actions needed to reach a goal, this is the best approach. But it must still be said that some lowly valued actions may lead to much higher rewards than greedy actions.

In slight contrast to a greedy agent, we have what is known as an ε -greedy agent, which is more explorative. This agent explores other actions with lower reward values with a small probability ε . All non-greedy actions have an equal chance of being chosen. This allows the agent to explore more alternatives and possibly find a better action than the current greedy action. However in my opinion, this is not a very balanced way of having the agent behave as it may take a very poor action as likely as the next best action.

A more explorative agent behaviour is called a softmax selection agent. It takes all actions with a probability proportionate to their value. It creates a probability distribution according to the formula below:

$$\frac{e^{Q_t(a) \setminus \tau}}{\sum_{b=1}^n e^{Q_t(b) \setminus \tau}}$$

(Sutton and Barto, 1998)

This ensures that the agent picks better actions more often than not. This is, in my opinion, very close to human behaviour and can be perceived as the agent being able to take risks, as it will not always take the greediest action, but it tends to take good actions rather than bad. Also this strikes a fine balance between exploitation and exploration. The parameter τ is called the temperature and decides how much the agent explores. The values of the temperature are bounded as $0 < \tau < 1$

High temperatures cause the actions to be all (nearly) equiprobable. Low temperatures cause a greater difference in selection probability for actions that differ in their value estimates. In the limit as $\tau \rightarrow 0$, softmax action selection becomes the same as greedy action selection. (Sutton and Barto, 1998)

Most agents start off with initial values all set to 0. However based on human knowledge of the system, optimistic initial values can be set. As they systems interacts with its environment, it will update the values based on formula defined by the user. This gives the agent a statistical bias. However Sutton and Barto (1998) make the point that it provides “an easy way to supply some prior knowledge about what level of rewards can be expected.”

2.3 Blackjack in Reinforcement Learning

Taking the existing Card/Pack system, which works very well, I intend to rebuild the functionality of the Player and Main classes. I will re-implement the basic GUI I had before, but the Player class will now be redefined. Each game will have 2 types of players, a human player, an improved version of the previous system, and a reinforcement learning agent. This will be the main concentration of my project.

Blackjack is a fairly simple game with minimal mathematical complexity both in terms of rules and betting. Also a relatively small number of actions are required to reach an outcome and hence a reward. Both these factors, make Blackjack suitable as a reinforcement learning problem.

After having made a successful agent for blackjack, it should be possible to abstract the basic functionality of the agent and hence create AIPlayer, which would integrate with my old framework and make an infinitely extensible suite of card games, with AI playing capabilities.

Blackjack is what is known as a Markov Decision Process (MDP). Blackjack is said to have the Markov property because, given any state, you can easily trace back previous states, as well as the actions leading to those states. This is the case with a large majority of Reinforcement Learning problems.

The basic state representation will be the sum of the value of the cards the player has. This is the most important information, as far as the game is concerned. It would be wasteful to have a separate state for each combination, as they are irrelevant, a 6 and 7 is the same as a 5 and 8.

However there are some complexities involved in this. The most prominent case is the fact that an Ace can count high (11) or low (1). Also deriving from this is the fact that when a high ace would cause a player to be bust, it is automatically switched to a low ace. Another complexity would be enabling the agent to split. Considering the state is only the sum, we can miss on a split, as a 5 and 5 is treated the same as a 2 and 8. Also the double down rule would be a minor complication, but this would not be due to the state representation, but because of the inherent complexity of the rule itself.

Chapter 3

Project Description

3.1 Technologies involved

The system will be developed using the Java programming language, as the existing system is written in Java. From this fact, it follows that any additions to the system must be in Java. To change to any other programming language would require rebuilding the entire system, after it has been refactored. Due to the time constraints, it was decided to continue using Java.

In terms of Reinforcement Learning, the algorithms provided in Sutton and Barto (1998) are all in pseudo-code. Hence we can take these algorithms and implement them in any given programming language. As mentioned above, the existing system was developed in Java, hence the learning algorithms will all be implemented in Java to be added to the system.

3.2 System Design

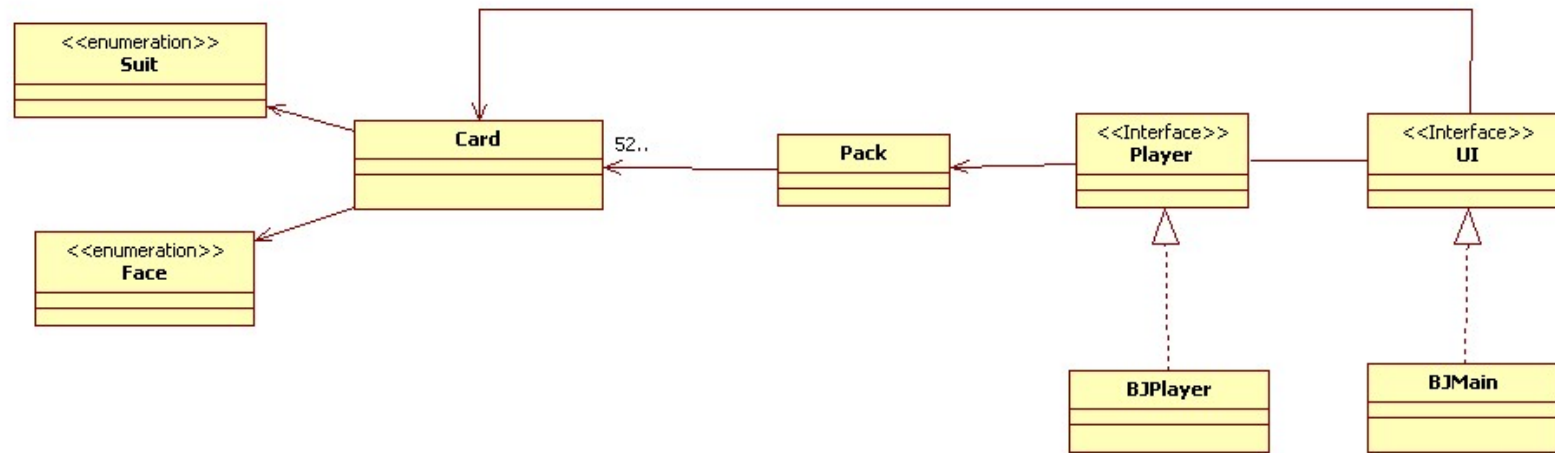
The old system is based on 3 main components. Firstly the Card/Pack system, which is currently the only fully operational part. This may still be modified to improve on the current functionality. The second part is Player interface, and its subclasses, which for the representation for the Player. Thirdly we have the Main interface, its subclasses and SKMain which provide the GUI elements. SKMain is the Main Menu of the program, but does not implement Main. This can be illustrated in the class diagram (see figure 3.1).

This system is not very stable, as the dependencies and interactions between Main and Player are very complicated and tend to cause unexpected errors. The new system uses a less densely connected system to ensure a more streamlined system and easier communication between objects. The interface “GameRules” will be introduced to act as a conduit for communication. This is illustrated in the class diagram (see figure 3.2).

There are two basic Use Cases for the system, which a Player playing the game. This case is implemented in the old system. However this case is specific for a Human Player. The second Use Case is the Reinforcement Learning agent actually learning how to play the game. This case actually includes the first Use Case. This can be illustrated by the Use Case Diagram (see figure 3.3).

Based on the reading and research, it has become apparent that Blackjack is an episodic task. The agent will learnt something each time it plays an episode, which in this case is a game of Blackjack. For Blackjack, there can be no more than 2 states between the initial state and the final state. This is because a Player may not have more than 5 cards in his hand. The initial state is 2 cards. Between these states we have 3 and 4 cards as the two possible intermediate states. The agent will operate as illustrated in the sequence diagram (see figure 3.4).

However it must be noted that in some problems, there may be a reward at an intermediate state. This is not true for Blackjack, hence this was not included in the diagram.



1

Figure 3.1: Old System Architecture

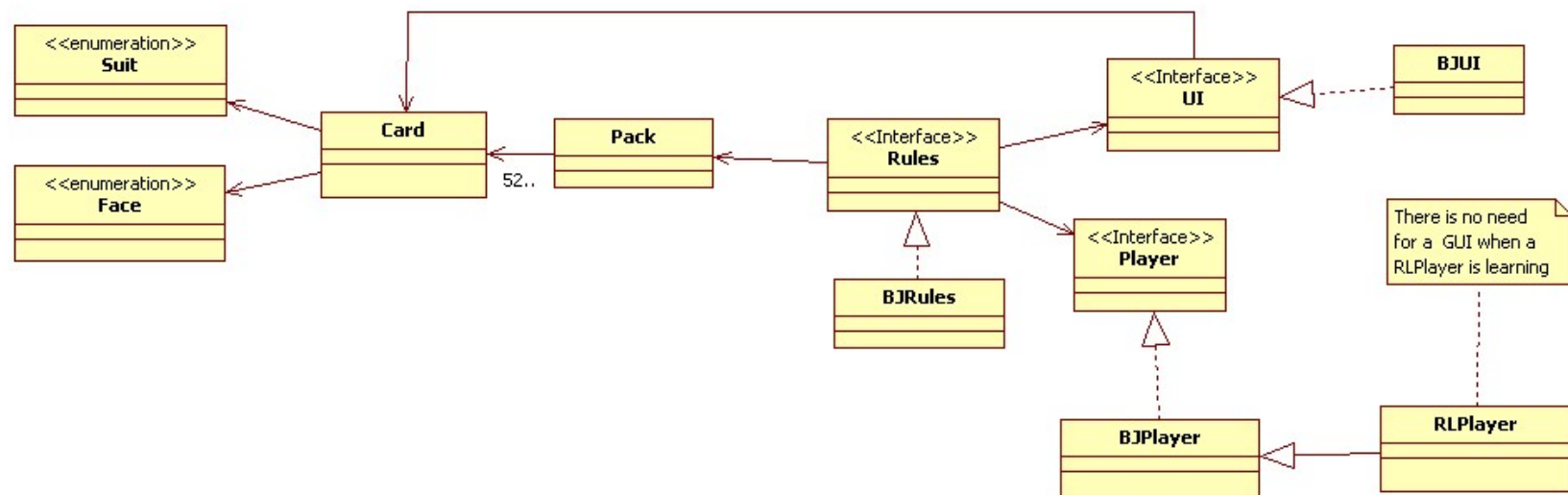


Figure 3.2: New System Architecture

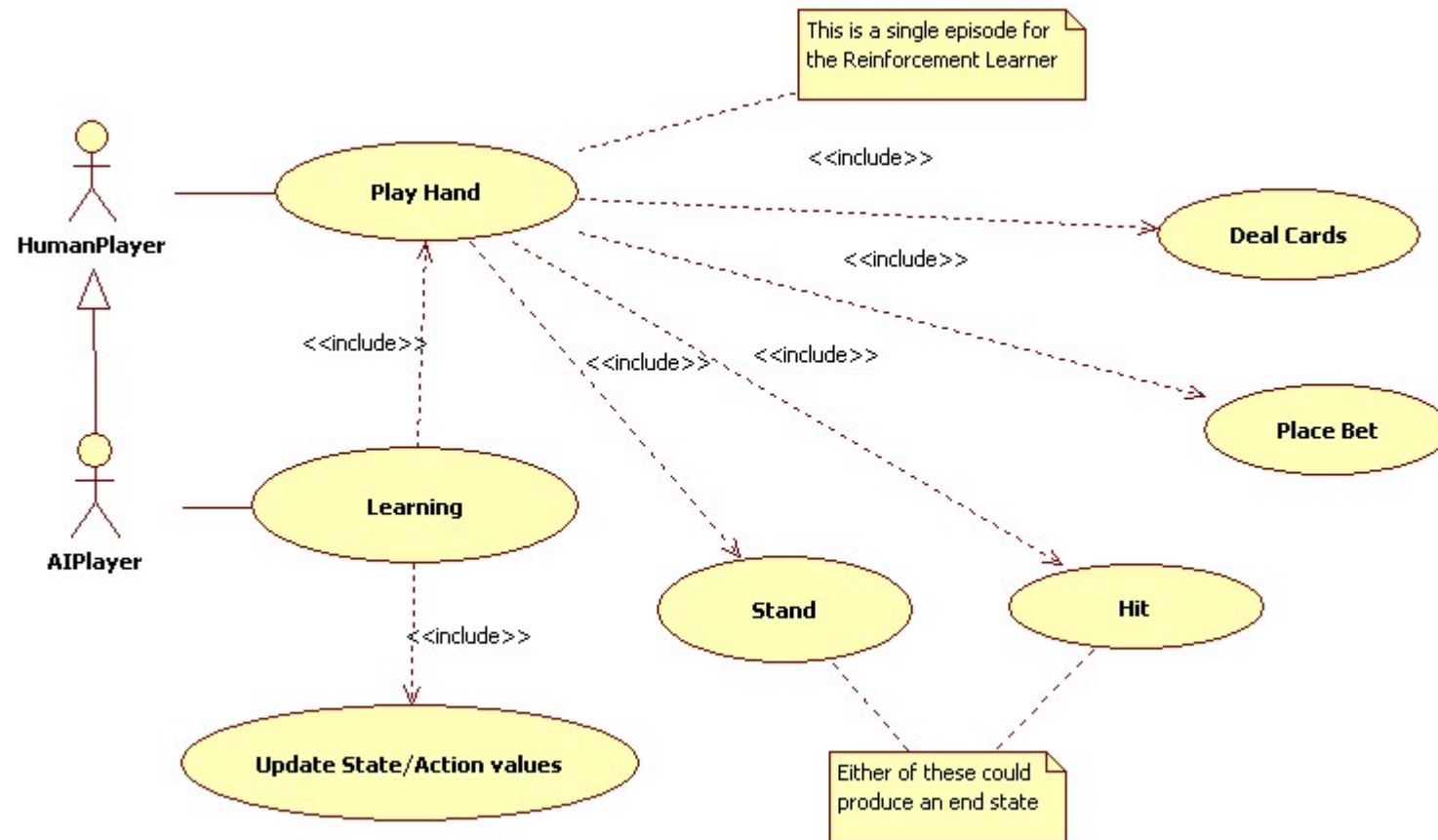


Figure 3.3: Use Cases

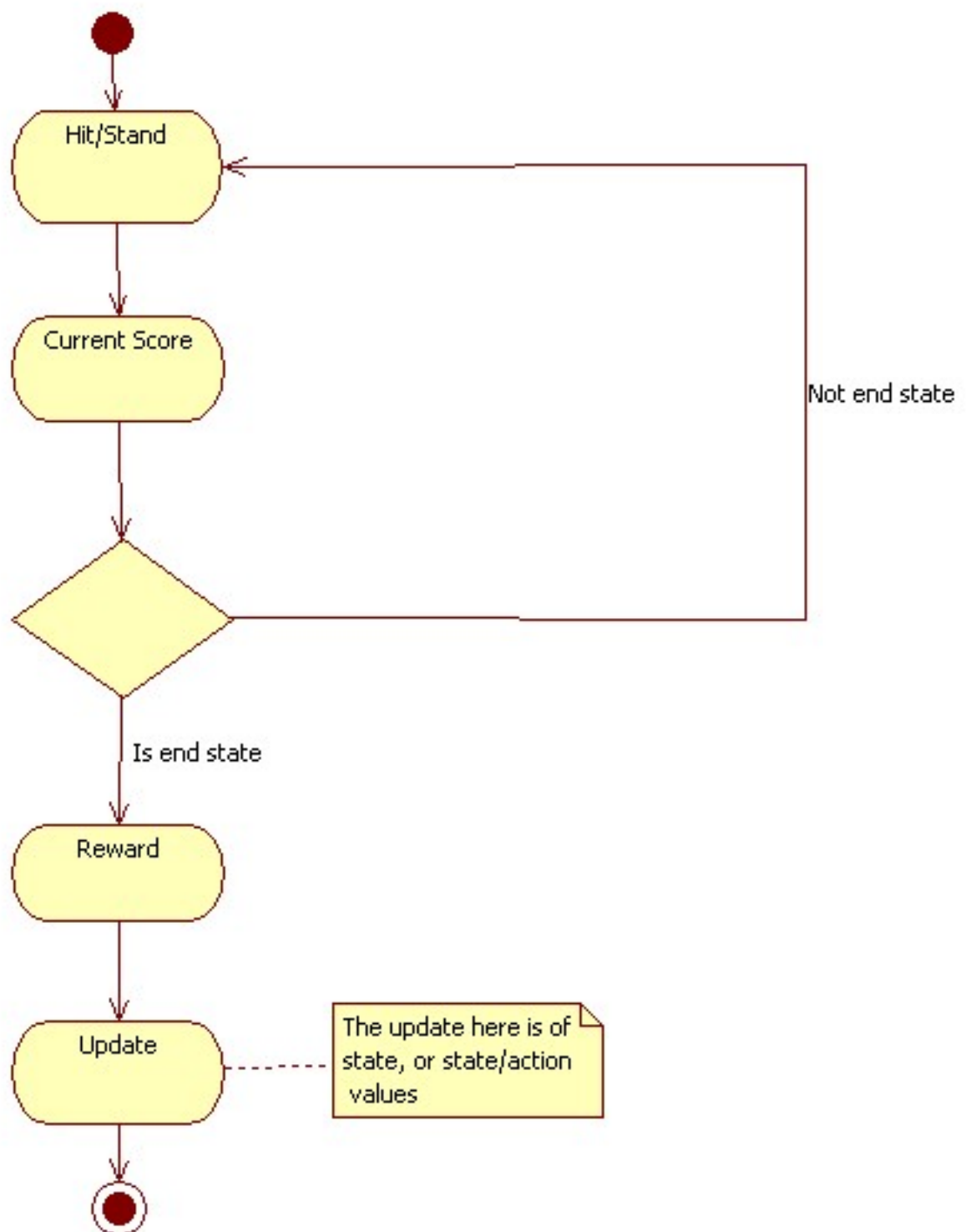


Figure 3.4: Single Episode

Chapter 4

System Development

4.1 Prototyping

To better understand Reinforcement Learning, two prototypes were built before the final system was built. These were built chosen in such a way that they shared some features with the final system, thus enabling code reuse. Both of these prototypes are discussed below.

4.1.1 Prototype 1 - Next card higher/lower

This system trains an agent to guess if the next Card is higher or lower than the current card. This was chosen so as to see how the Card/Pack system would be used in an agent. No proper Reinforcement Learning algorithm was used, but merely a sum of all past rewards.

The agent used the `points()` method in the Card class to determine if a card was higher or lower. This gives a Joker 0 points, Ace 1 point, number cards face value and Jack, Queen and King get 11, 12, and 13 respectively. The agent received a reward of +1 for guessing correctly, and -1 for guessing incorrectly. The source code for this prototype can be found in Appendix A.1.

The prototype did not take into account when two cards were equal and just counted this as an incorrect guess. This situation does not occur very often and was not considered for the prototype as it was simply a proof of concept as opposed to a proper system. Results from this prototype can be found in Appendix A.2.

This prototype learns a policy which is fairly consistent with what is expected. However as this does not use a proper Reinforcement Learning method, it is not guaranteed to converge on P^* . Despite this, it uses the softmax probability as described in section 2.2.3.

From this prototype, code that calculated softmax probabilities was made. This was carried over to the next prototype and the final system. Also it showed how to correctly use the Card/Pack system in a Reinforcement Learning Agent. Parts of this code was carried over to the final system.

4.1.2 Prototype 2 - GridWorld

This prototype solves what is known as the GridWorld problem. This problem has an agent in a grid-based environment. Each square is either a goal, a penalty or empty. The agent begins in a random square and must navigate its way to the goal(s) and avoid any penalty(ies). If the agent moves to a penalty or goal square, the episode ends. It receives a reward of +100 for a goal and -100 for a penalty. The source code for this prototype can be found in Appendix B.1.

For this prototype, the Monte Carlo method was used (see Sutton and Barto (1998) chapter 5). This is mathematically proven to converge on P^* , which this system does. For an easier graphical representation, states were valued instead of state/action pairs. As stated before either approach is valid and can be chosen based on the users requirements. Some results of this can be found in Appendix B.2.

From the previous prototype, code that calculated softmax probabilities was used. This prototype created the code for reward functions and value calculation. The code was modified, as the Monte Carlo method was abandoned in favour of Q-Learning for the final system. However, the learning parameters were carried over.

4.2 Final System

4.2.1 AI

As the Reinforcement Learning Agent was the primary objective of the project, it was built first. It was decided that the agent should be built using the Temporal Difference Q-Learning Algorithm. Although Sutton and Barto (1998) cover Blackjack under the MonteCarlo methods, Temporal Difference is a more powerful method in general.

For this agent, only hit and stand were considered as valid actions. The rest of the actions (as listed in 2.1.2) were not considered for the system, as it was thought that the agent would not learn anything of substance. These reasons are enumerated below:

- **BLACKJACK:** This would not help the agent learn anything except that a starting pair adding up to 21 is greatly desired. Hence the logical action would be to stand. The system will already learn it, as it will always get a negative reward for hitting at 21. Although the increased payout was added into the system, however the agent does not learn anything from this.
- **INSURANCE:** Taking out an insurance would not be beneficial in anyway to an AI agents policy, as the payout of this is based on probability. Also this occurs so infrequently, that it was not considered worth adding in.
- **SPLIT:** Splitting hands would only give the agent more hands to learn on. Although it would lead to more learning, it would not reduce the learning time enough to justify its inclusion. Also due to the complexity of being able to split many times it was not included in the final system.
- **DOUBLE DOWN:** Doubling down would lead to a hit and then stand. The agent would effectively learn a policy for hitting at a score s_n and standing at score s_{n+1} . This is would not lead to any additional learning and may also impede the learning of the agent. As a separate action, it would not add anything useful, as it is just a hit and a stand.

For further clarification, a graph indicating how frequently each action may have been taken is included below. As can be seen, all the actions occur less than 10% of the time. This would not have lead to any significant learning. Although double down occurs over 70% of the time, as stated above it would not add anything meaningful to the learning.

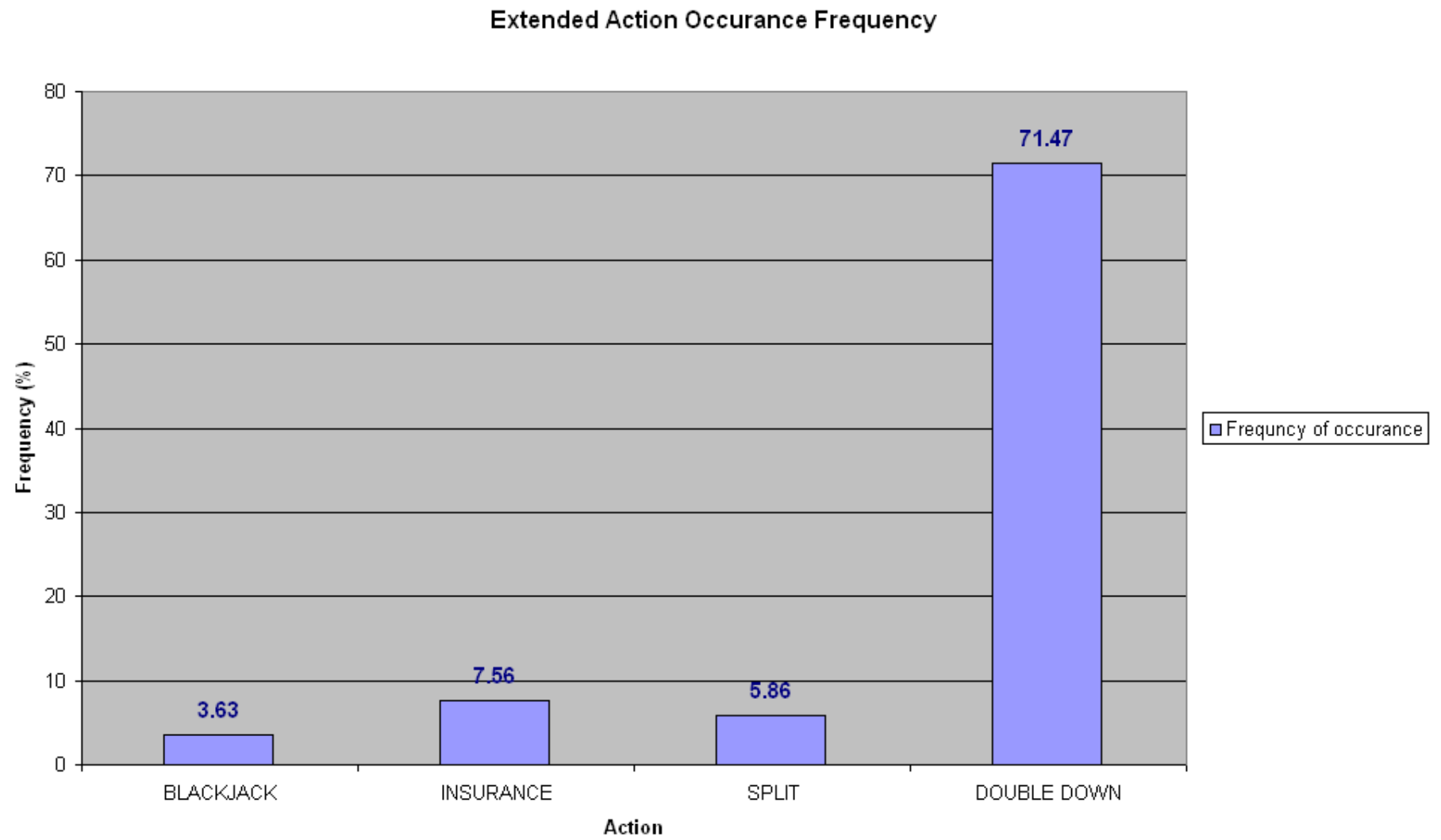


Figure 4.1: The frequency of the extended actions

```

1 double hit = Math.exp(val[HIT][score - 1]/TAU);
2 double stand = Math.exp(val[STAND][score - 1]/TAU);
3 double sum = hit + stand;
4 double hprob = hit / sum;
5 double sprob = stand / sum;

```

The first step was to build the graphical user interface. This was built with the same design as the previous system. This was to ensure that the maximum of code could be re-used in the new system, so as to avoid any complications of porting the code from one interface to another.

After the interface was built, the underlying game mechanics were built. For the AI, only hit and stand were considered for actions. The reasons for this have already been discussed above. To facilitate this, a class called BPlayer was created as an inner class, to facilitate easy access. This was used instead of using accessor methods to manipulate the object for simplicity and quick implementation.

The BPlayer was created with the basic functionality required for hit and stand. The BPlayer had no real functionality, apart from the Reinforcement Learning, as most of the gameplay is done in the Blackjack object. After this was established, the Reinforcement Learning methods were added. The player has an exploration method, which uses softmax, a greedy method and a fixed method, for the dealer.

Finally the Reinforcement Learning agent was developed. This mostly used code from the prototype, with modifications to fit into the system. The completion of the Reinforcement Learning agent makes both system equal in functionality. Although it must be noted that new AI works, where as the old AI does not.

To achieve this, we need a state representation. For Blackjack, this is simple as it just the score of the player minus 1. The actions were stored as numbers. The values of the state action pairs were stored in a 2-dimensional array called val. The value for action A in state (score) S is given by $\text{val}[A][S-1]$. The states are represented as $S-1$ as Java uses 0 as the first array index and $n - 1$ as the n^{th} index, hence the $S-1$.

The softmax method was fairly simple to implement as there are only 2 possible actions, hence only two probabilities. These were calculated using the following code:

After these probabilities are calculated, we update the last state-action pair, if any. For the case where it is the first action, *lastaxn* is set to DEAL(-1). Then a random r number is generated to choose the action. After the action is chosen, we set the values of *lastaxn* and *lastscore*. This is achieved with the code given below:

The above code is an implementation of the general formula:

$$Q_{(s_t, a_t)} = Q_{(s_t, a_t)} + \alpha(r_{t+1} + \gamma \max_a Q_{(s_{t+1}, a)} - Q_{(s_t, a_t)})$$

(Sutton and Barto, 1998)

Here the reward r_{t+1} is 0, as no reward is given.

After the learning method, a greedy method was devised. This is for use when the agent has learnt a policy $P \approx P^*$. As stated previously, after exploring, the agent should exploit what it learnt. This was implemented using the following code:

```

1  double maxval = Math.max(val[STAND][score - 1], val[HIT][score - 1]);
2  if(lastaxn != DEAL){
3      double rd = A * ((G * maxval) - val[lastaxn][lastscore - 1]);
4      val[lastaxn][lastscore - 1] += rd;
5  }
6  double r = Math.random();
7  lastscore = score;
8  if(r > hprob){
9      stand();
10     lastaxn = STAND;
11 }
12 else{
13     hit();
14     lastaxn = HIT;
15 }

1  if(val[HIT][score - 1] > val[STAND][score - 1]){
2     hit();
3 }
4 else{
5     stand();
6 }

```

Finally the reward function was created. This gives a reward to the agent, based on the game outcome. After the dealer stands, the outcome is calculated. There are 5 possible outcomes. Based on the outcome, the player was judged to have achieved one of the 4 results and based on that it was given a reward. Listed below are the 5 outcomes and their end state and reward in brackets:

1. Player Busts (Bust -10)
2. Dealer Busts (Win +10)
3. Player and Dealer bust (Bust -10)
4. Player Wins (Win -10)
5. Player Loses (Lose -10)
6. Push (Push +5)

This was achieved using the code below:

```

1  if(player.score > 21 && dealer.score <= 21){
2      //player bust
3      player.reward(-20);
4  }
5  else if(player.score <= 21 && dealer.score > 21){
6      //dealer bust
7      player.funds += bet * 2;
8      player.reward(10);
9  }
10 else if(player.score > 21 && dealer.score > 21){
11     //dealer and player bust
12     player.reward(-10);

```

```

13 }
14 else if(player.score > dealer.score){
15     //player wins
16     player.funds += bet * 2;;
17     player.reward(10);
18 }
19 else if(player.score < dealer.score){
20     //player loses
21     player.reward(-10);
22 }
23 else{
24     //push
25     player.funds += bet;
26     player.reward(5);
27 }

```

The reward method uses the same formula as employed by the learning method, except the reward is explicitly passed to it as can be seen above. The code for this method is given below:

```

1 private void reward(double r){
2     double maxval = Math.max(val[STAND][lastscore-1], val[HIT][lastscore-1]);
3     if(lastaxn != DEAL){
4         double rd = A * (r + (G * maxval) - val[lastaxn][lastscore-1]);
5         val[lastaxn][lastscore-1] += rd;
6     }
7 }

```

It must be noted that the final class has more code than listed above, but this code was for testing purposes only and does not interfere with the learning algorithm. For the final code see Appendix C.1.5.

4.2.2 SKCards

The first feature that was added to the SKCards was the Shoe class (see Appendix D.1.5). This class is used to represent what is known as a dealer's shoe. This is in effect a Pack, which is not fully used. A cut card is inserted at random in the Pack. Once this cut card is dealt, the Shoe is discarded and a new one is brought into play. This is used mainly in casinos for Blackjack to counter card counting. This was implemented for completeness of the system. The dealer does not distinguish between "soft" and "hard" scores. This is as dealer play was not the objective of the experiments.

Then the Player interface was reworked, to fit in the new system. This included eliminating all relations to the Main interface, which was renamed to UI. Simultaneously, the UI interface was modified to eliminate all relations to the Player interface. The Rules abstract class was then created and the two interfaces were joined using this. With this the framework was completed.

After completing the framework, the existing BJPlayer and BJUI (formerly BJMain) were built up. Where possible, code from the original system was used, in its original form or modified. Where it was not possible, code from the prototype was used, or it was made from scratch. This now makes the Old system, equivalent to the old system in functionality, barring the AI.

Chapter 5

System Evaluation

5.1 System testing

There are 4 basic tests that have been devised, all of which depend on the agent being able to react to a dealer policy effectively. As previously stated, we wish to come as close as possible to P^* , which in this case would be a policy better than its opponent (the dealer).

Each test uses different values of rewards to achieve a different purpose. The first test uses the normal values as listed below:

- Win = +10
- Push = +5
- Lose = -10
- Bust = -10

The subsequent tests involve doubling the reward for win, lose and bust and keeping the rest at their normal values. The push value was not modified as its low frequency of occurrence would mean its effect is minimal. These values were tested against 7 dealer policies, namely dealer hits below 11-17. As the dealer plays a more aggressive strategy, the agent should also adjust its policy, within the parameters provided.

This range of policies is that the two boundary values represent significant policies. The lower bound (11) represent the policy a player must play to avoid going bust. The upper bound (17) represents the policy that is employed by several casinos, making it the de facto standard. All other policies were included for completion.

The constants in the learning equation were kept the same as listed below:

- Softmax temperature $\tau = 0.5$
- Step size parameter $\alpha = 0.6$
- Discount rate $\gamma = 0.75$

5.2 Results

5.2.1 Table

This table provides the results of running the agent with different parameters. The data and therefore the results are highly susceptible to noise, due to the random element in the system. However separate data sets with the same parameters tend show a consistency of approximately $\pm 2.5\%$, giving a total error range of 5%. This error can be attributed to noise in the data, as well as the random factor within the game. Due to the nature of card games, even an optimal action can lead to a sub-optimal result, causing interference in the data. This is mitigated to a certain extent by the Law of Large numbers.

DEALER POLICY	WIN REWARD	LOSE REWARD	BUST REWARD	WIN(%)	LOSE(%)	BUST(%)	NET WINS(%)
11	10	-10	-10	44.62	17.6	30.66	-3.64
12	10	-10	-10	43.38	19.88	29	-5.5
13	10	-10	-10	37.58	21.98	29.06	-8.96
14	10	-10	-10	34.76	22.28	29.34	-10.76
15	10	-10	-10	27.86	23.08	30.22	-12.46
16	10	-10	-10	26.2	23.9	29.26	-11.74
17	10	-10	-10	19.08	17.4	38.74	-12.86
11	20	-10	-10	45.98	23.1	22.02	0.86
12	20	-10	-10	42.76	19.54	29.82	-6.6
13	20	-10	-10	37.62	20.96	29.12	-7.92
14	20	-10	-10	35.02	21.92	29.66	-11.04
15	20	-10	-10	28.14	24.06	28.72	-12.26
16	20	-10	-10	26.12	23.88	29.88	-12.82
17	20	-10	-10	17.92	25.46	29.54	-13.18
11	10	-10	-20	43.02	47.22	0	-4.2
12	10	-10	-20	37.1	52.5	0	-15.4
13	10	-10	-20	30.06	56.22	0	-22.02
14	10	-10	-20	27.04	57.42	0	-24.4
15	10	-10	-20	19.48	58.88	0	-26.26
16	10	-10	-20	18.3	58.26	0	-25.72
17	10	-10	-20	12.44	54.34	0	-16.76
11	10	-20	-10	29.26	3.92	60.08	-34.74
12	10	-20	-10	27.52	5.48	59.56	-37.52
13	10	-20	-10	25.28	5.02	59.86	-35.44
14	10	-20	-10	25.28	5.02	59.32	-32.9
15	10	-20	-10	20.66	5.86	60.84	-32.28
16	10	-20	-10	20.46	5.8	59.86	-31.12
17	10	-20	-10	16.14	6.04	59.02	-25.7

Table 5.1: Final Results

5.2.2 Graphs

All the data from table 5.1 is illustrated graphically in the graphs below. Each graphs show the comparison between the 4 separate testbeds. For each run the following values were calculated:

1. Final Policy
2. Winning Percentage
3. Losing Percentage
4. Bust Percentage
5. Net Winning Percentage

The final policy is the highest score at which a greedy agent will hit and will stand on any higher score, based on the values learnt. It is expected that as the dealer policy increases, the agent will also increase its policy to try and stay ahead of it. It is also expected as agent's reward values are adjusted, its policy will also be adjusted.

The winning percentage is the percentage of games in which neither dealer nor agent has gone bust and the player has won the game. Although when the dealer goes bust, it is considered that the agent has won, this was not considered as this situation arises from poor play from the dealer and not necessarily good play from the agent.

Similarly, the losing percentage is the percentage of the games in which the dealer has won. Again the case of the player going bust was not considered under here as it can be attributed to poor play by the agent and not necessarily good play from the dealer.

However, busting is an essential part of the game and must not be ignored. To this effect, the percentage of busts was recorded. This number is simply the percentage of times the player has been bust during the course of the learning.

When considering all of these figures, we must remember that there is a large random factor that affects the game and therefore the results. Although the step-size parameter can smooth the learning process, there is really no way to compensate for the noise in the final results data.

To consolidate all these figures, the net winning percentage was also calculated. This can be found using the formula below:

$$\text{net wins} = (\text{wins} + \text{dealer busts}) - (\text{losses} + \text{player bust})$$



Figure 5.1: Q-Learning Policies

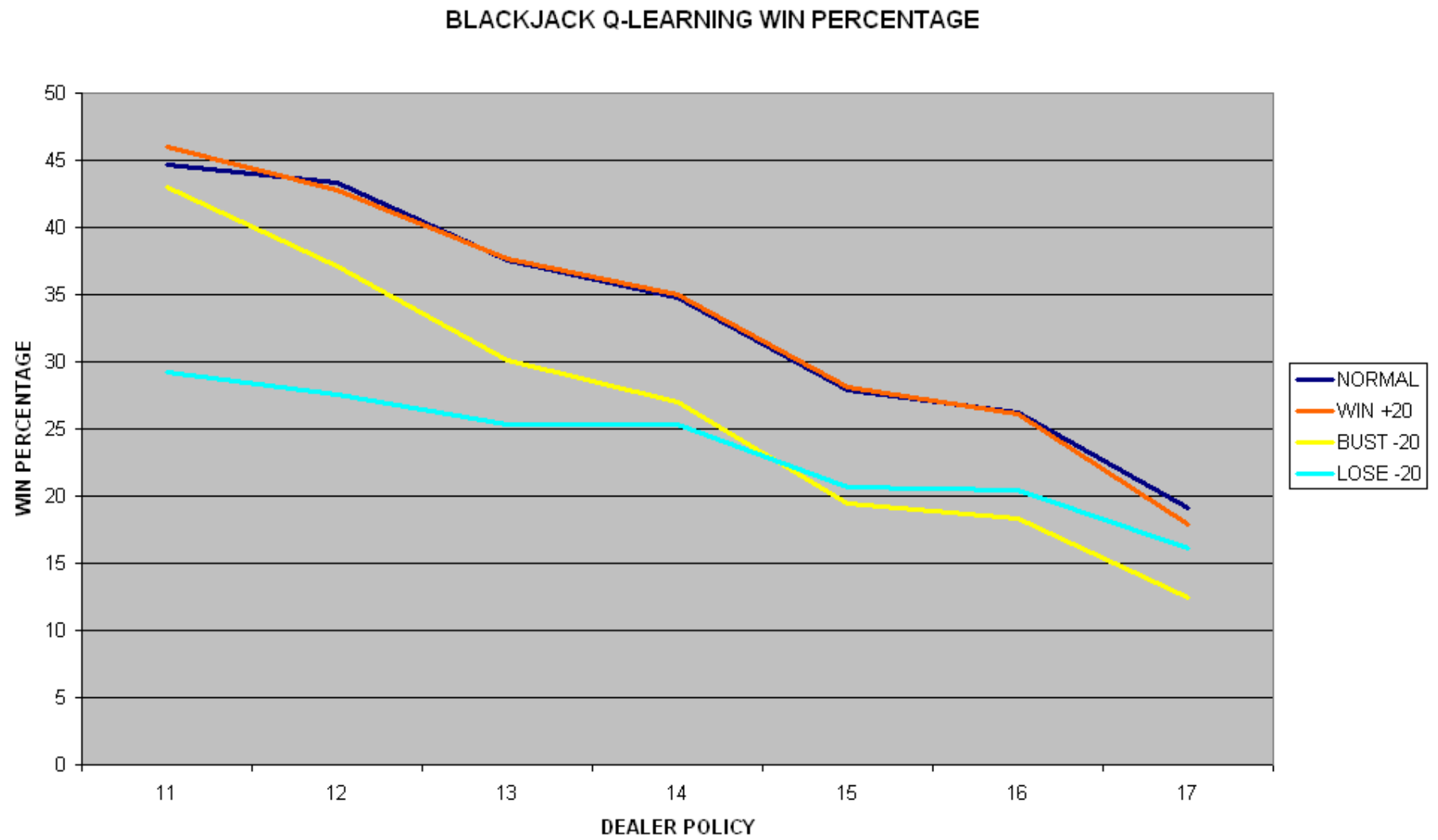


Figure 5.2: Win Percentages

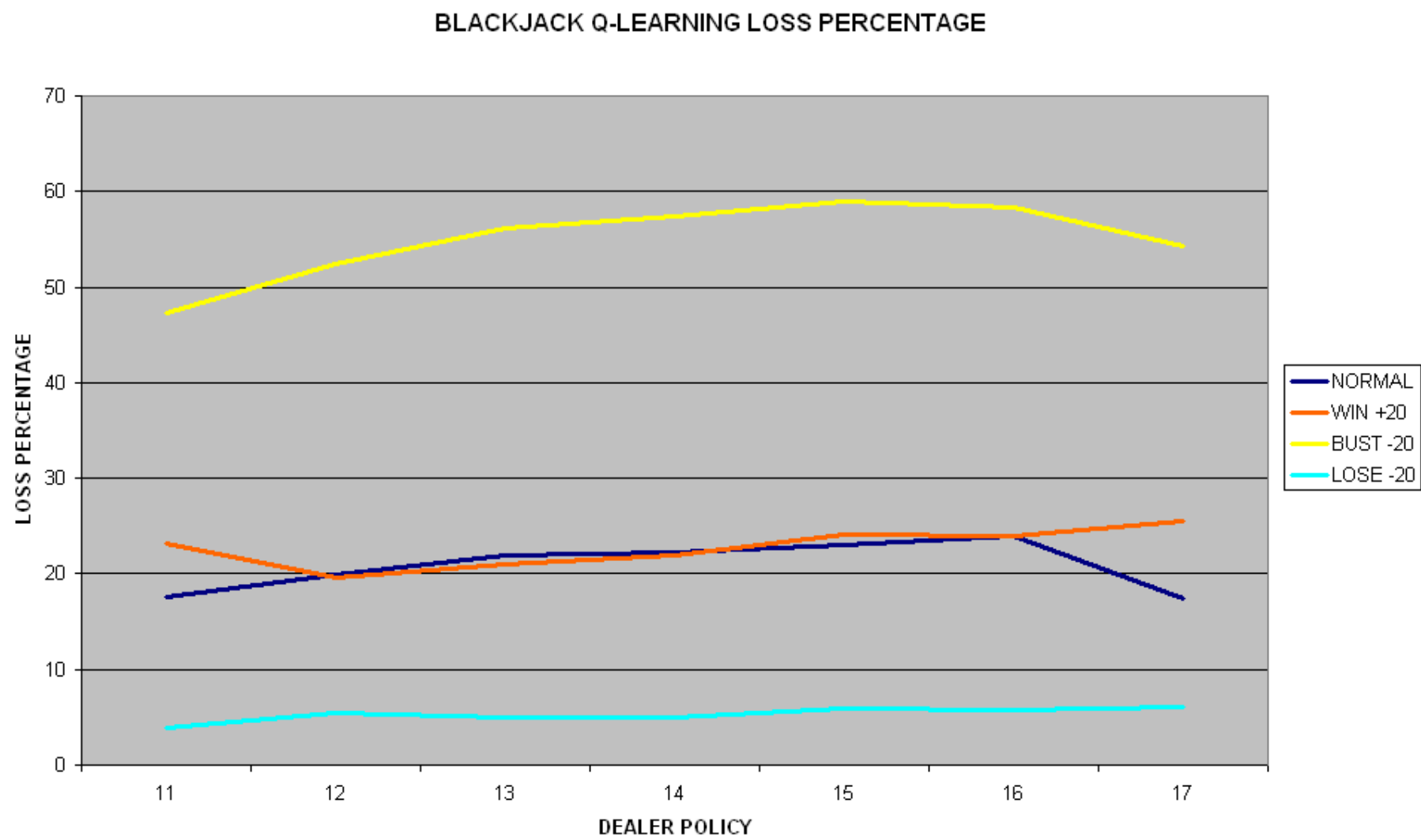


Figure 5.3: Loss Percentages

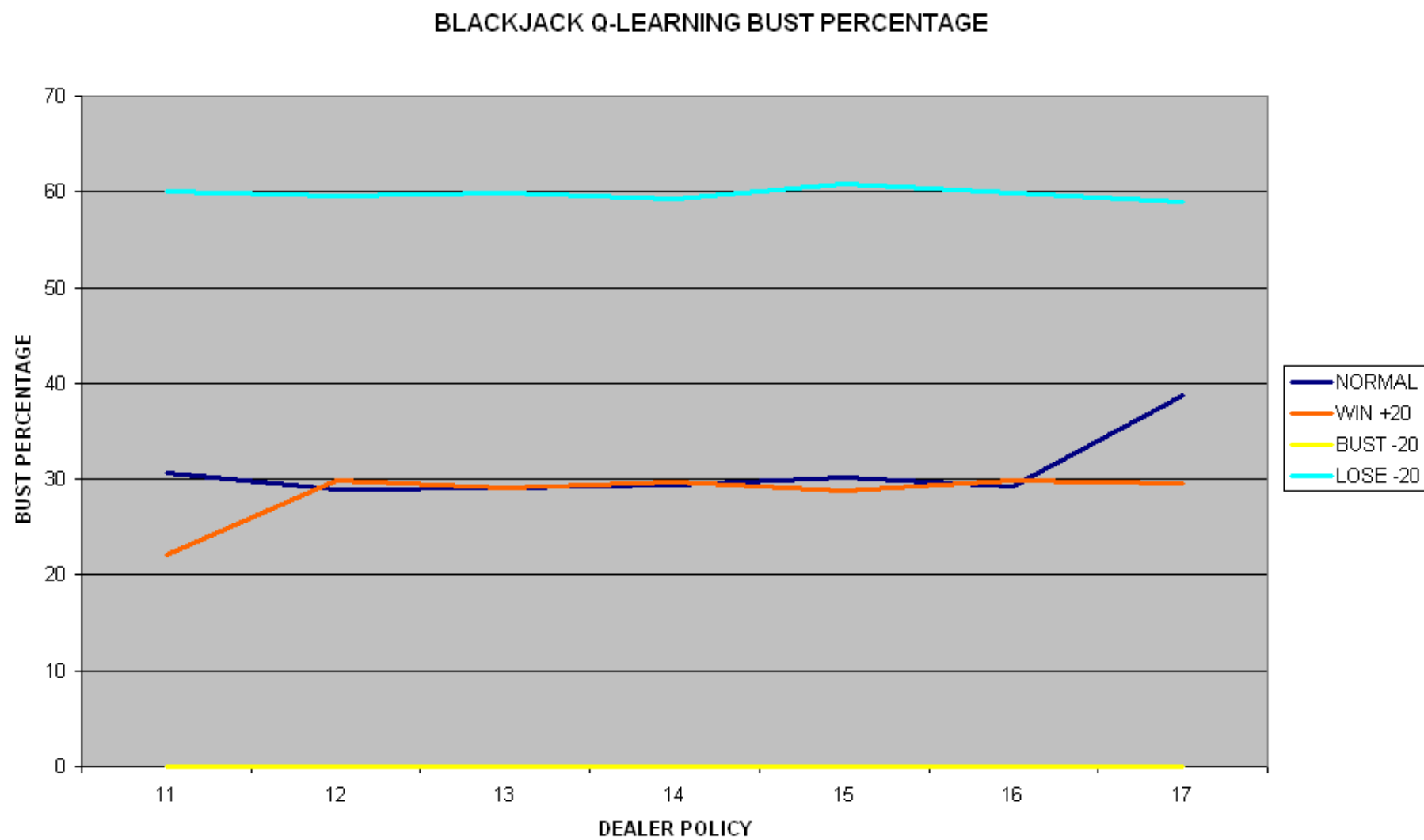


Figure 5.4: Bust Percentages

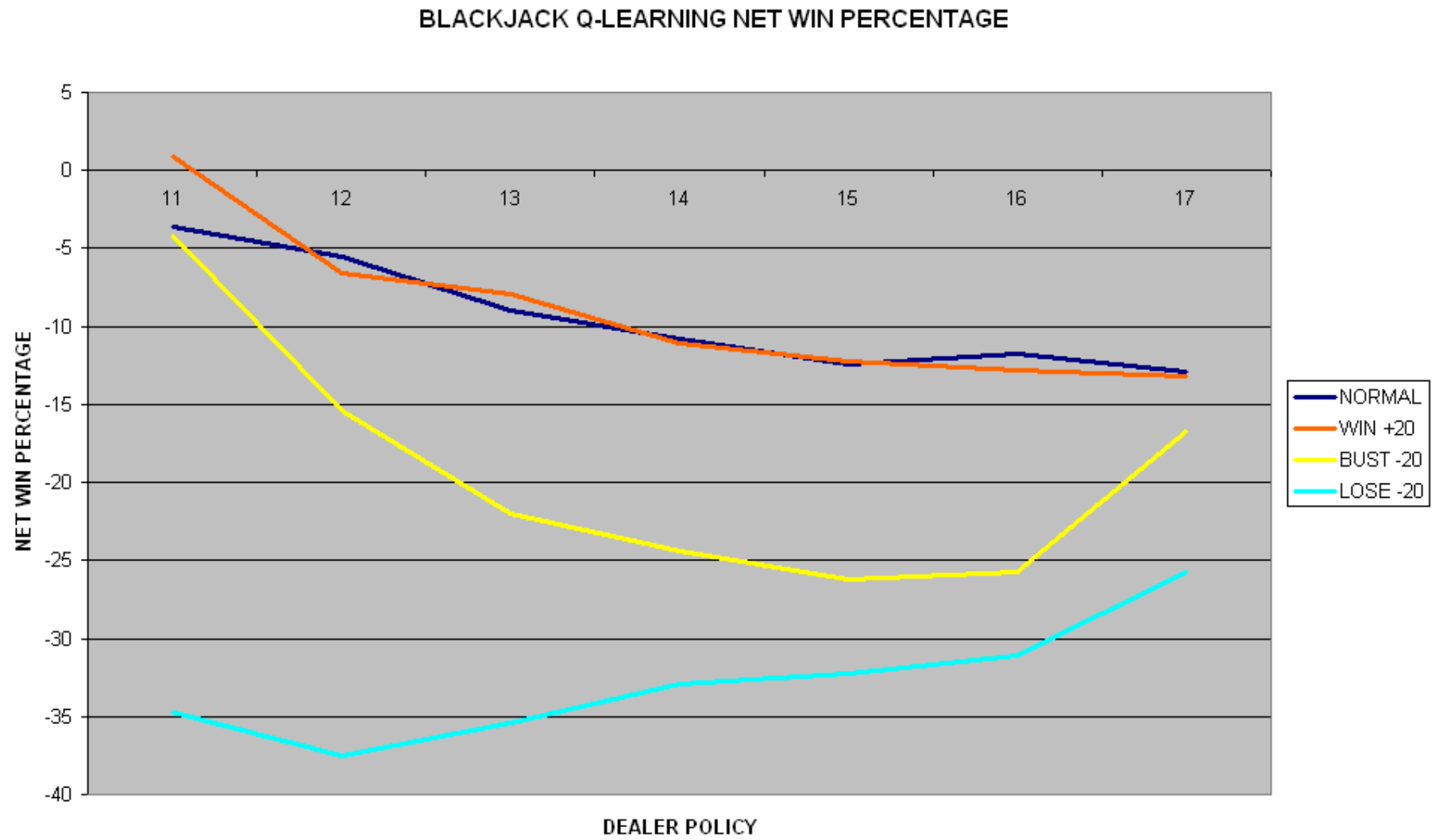


Figure 5.5: Net Win Percentages

5.3 Result Analysis

As can be seen from Figure 5.1, when the agent is given a reward of -20 for busting, it learns the most effective policy to counter it, that is hitting at 11 or below. Here we can see the agent prefers to lose (-10), rather than go bust. The best possible reward is +10 for it winning. However, it chooses not to pursue this, as it is simply tries to avoid going bust and receiving the -20 reward. However, when it is given a reward of -20 for losing, it plays a very aggressive strategy, specifically 20, hitting on everything but a 21. Here the agent tries to avoid losing, even at the risk of busting. The best way to do that is play a very aggressive strategy, which it does. By increasing the win reward to +20, the agent plays a slightly more defensive strategy. However, on several points, both agents pick the same strategy. Clearly, the negative rewards have a greater bearing on the learning than the positive rewards. Hence any tuning of the agent would concentrate more on the negative rewards than the positive.

One point that must be noted is that the agent consistently learns a threshold score at which to stop hitting. This is not specified on the algorithm. The algorithm merely provides the agent the ability to value the actions and hence build a policy based on its experience. This can be said to be one of the major successes of this project, in terms of the Reinforcement Learning Agent.

Furthermore, from Figure 5.2 it can clearly be seen that as the dealer plays a more aggressive strategy, the player will lose more often. This is as expected, because as the dealer plays a more aggressive policy, it will score higher, hence making it harder for the player to win. A point worth noting is that having a higher bust and lose penalty results in a lower winning percentage, where as increasing the win reward has a minimal effect.

As we proceed to Figure 5.3, we can see that this graph shows trends opposite to Figure 5.1. This is expected as the better policy the agent plays, it should lose less. The lines are not as smooth, due to the noise in the data, but the basic trend can be seen. Also it must be noted that the dealer policy has a minimal effect on the loss percentage of the agent.

Figure 5.4 shows the same trend we saw in Figure 5.1. The more aggressive policy the agent plays, the more it will bust. As it chooses to hit on higher scores, there is a larger probability of it going bust. Hence, the same trend repeats itself, as expected. However, it must be noted that the agent busts more than the dealer. The reason for this is not known, but it has been observed that the player busts about 5% more than the dealer playign the same policy.

From the above graphs and analysis it is clear that the best results are obtained with the winning set at +10 or +20. Changing the penalties may improve the performance in certain areas, but the overall performance is degraded. Hence the best agent under these circumstances would have a win reward of +10 or +20.

5.4 Possible Future Improvements

As can be seen from Figure 5.5 the agent tends to lose more than it wins when playing a policy of 12 or greater, no matter what the learning parameters. As this is the primary performance benchmark, any non-trivial improvement to the system would have to increase this figure. There are two possible ways this can be done, card counting and probability calculations. Both of these are discussed below.

5.4.1 Card Counting

Card counting is a tactic used by some players to try and “beat the system”. This tactic involves keeping a running total of what is know as the “count” of the deck. Each low card “counts” as +1 and each high card “counts” as -1. By adding the “count” of all the cards observed so far, the player can tell how “heavy”

or “light” a deck is, indicating a large number of high or low cards respectively. Based on this, the player can adjust their strategy.

The best way to combat this is to use a Shoe, which has been implemented (see Appendix D.1.5). Once the shoe is changed out, the player no longer has a count of the deck. This would mean that they must now begin to count again. While they are counting, they would presumably play a fixed strategy. By keeping the cut-off point as low as possible, card counting can be countered.

5.4.2 Probabilities

Another possible way to improve the system is by using probabilities. This assumes knowledge of how many decks are in play. Normally this player information is kept hidden from the player. However, if the number of the decks is known, the number of each card is known. Based on cards already played out, the agent can calculate the probability of getting the required card(s). Using these values, the agent can adjust its policy.

The simplest counter to this is to simply not reveal the number of decks. Although a well programmed agent can calculate several probabilities at the same time. It could simply maintain multiple variables and update each with based on Cards observed. The number of decks could be inferred from the frequency of Cards appearing, e.g. 3 occurrences of the King of Spades would indicate at least 3 decks. Again, if we use a Shoe instead of a Pack, we can also quite effectively counter this.

5.4.3 Busting

It has been noted that the agent seems to bust too much. There is no apparent reason for this. Based on analysis of the code in Appendix C.1.5, the system should perform equally to the dealer. However this is not the case. The software needs to be thoroughly investigated to determine the cause of this error. Once found the error needs to be rectified. This would improve the net win percentage of the system, which as stated above is the measure of improvement in this system.

5.4.4 SKCards

There is scope for improvement in the SKCards framework. The Blackjack only implements the core rules. The extended rules can be implemented for the player and make it a complete Blackjack game. Also, the GUIs are very functional but do not look very attractive. These could be made more attractive. There is also the possibility of adding more games to the framework.

Chapter 6

Summary and Conclusions

6.1 Summary of the Project

6.1.1 Goals

There were two main aims of this project. The main goal was to develop a Reinforcement Learning agent for Blackjack. The second goal was to refactor the SKCards system and incorporate the Reinforcement Learning agent into it.

6.1.2 Outcomes

As we can see, the Reinforcement Learning agent was successfully developed. Although it was unable to beat the dealer, this is due to the nature of the game and not any shortcomings in the learning agent. However, it must also be noted that the large number of busts lead it perform lower than expected.

The refactoring was successful. The new system was created, and conforms the the model shown in Figure 3.2. However due the large amount of time spent on trying to implement the AI, there was little time for the final system. Although the system was built, it has only the most basic functionality and supports only the core rules as stated in Section 2.1.2. Although it was planned to have all the rules included in the system, this was not possible.

6.2 Conclusion

Based on the results gained from the system, it can be said that the Reinforcement Learning Agent learns an effective policy. Although the algorithm does not specify that the agent should learn a threshold score for hitting, the agent does learn this. This can be seen a large success of the learning algorithm. From figure 5.1 we can see that modifying the values of the bust and lose reward has a much greater effect than the win parameter. It can be concluded that the agent does learn to play Blackjack and that Reinforcement Learning is suitable for Blackjack.

Also we notice that as the dealer plays more aggressively, even though the player is playing an equal or better policy, they tend to lose more. It can be inferred from here that through pure strategy alone, the casino policy is very difficult to beat.

Also we can see there is some noise affecting the system, as can be seen in Figure 5.4. The graph is fairly consistent with what is expected as the number of busts relates to the policy the agent plays. Both figure 5.4 and 5.1 show the same basic trend. However in Figure 5.4 the last point for the agent with -20 reward for loss is significantly higher than the others. This is the noise and the random element of the game clouding the results.

Although the Law of Large Numbers mitigates this effect to a certain extent, it is not sufficient. The Pack is shuffled using the `Math.random()` function in Java, which is not very effective at generating very large numbers of data. This factor is not very significant and does not affect most of the tests. However, this, combined with noise and the random element in Blackjack, can explain erroneous results such as this.

Chapter 7

Project Self-Evaluation

7.1 Introduction

This project aimed to build a Reinforcement Learning agent for Blackjack and build it into an existing framework. The framework had to be refactored to make it more usable and conform with software engineering norms. It can be said that this was very successful to a large extent.

7.2 Project outcomes

The first outcome, which was to develop a Reinforcement Learning agent, was successful. The agent does learn a policy of hitting until a threshold score. This score varies on the rewards it is given. However it can be seen from Figure 5.1 that the agent does learn a suitable policy based on what it is told to do. There is however, the problem of the agent busting more than the dealer.

The second outcome was to refactor the old system and include the Reinforcement Learning agent. This was successfully achieved. However, it was only built to support the core rules of Blackjack. This was due to a planning oversight on my part, which meant that I spent too much time on the AI development.

7.3 Project planning

The initial planning of my project was quite challenging. There were several tasks to be completed and these all needed to be scheduled. The 3 main tasks were:

- Background reading and research
- Prototyping
- System development

Although I had allocated time for all these tasks, the prototyping stage spilled over into the system development stage. This gave me less time to work on the system. Barring this, the project went according to plan. There was a maximum deviation from plan of 1 or 2 days. I now realise that I should have allowed a larger buffer time in between tasks for spill over and unforeseen errors

7.4 Self-evaluation

During the course of this project, I learnt the limitations of my programming skills. Although I consider myself highly competent in terms of programming, developing graphical user interfaces is an area where

there is great room for improvement. My interfaces tend to lean towards functionality at the expense of appearances. This makes the interface less aesthetically pleasing.

Also, I have realised that my planning tends to be too optimistic and does not account for errors and other spill over. This is due to my underestimation of buffer time between tasks, which would be taken up by error correction.

Finally I have learnt that I am quite competent at analysis of numerical data. Based on the prototypes I have built and the results thereof, (see Appendix A.2 & B.2) it has become apparent to me that my mathematical skills are better than I thought before the project. However, if I work with the same data for too long, I tend to miss the bigger picture.

Appendix A

Card High/Low prototype

A.1 Source Code

```
1  //Adding HighLow to SKCards
2  package SKCards;
3
4  /*
5   * Importing the ArrayList class
6   * and the IO components
7   */
8  import java.util.ArrayList;
9  import java.io.*;
10
11
12 /**
13  * This is the 1st prototype for the Reinforcement Learning agent
14  * for the SKCards system, specifically SKBlackJack. This agent
15  * attempts to guess if the next card is higher or lower than the
16  * current card. It achieves this by using a simple sum of all
17  * the past rewards.
18  *
19  * @author Saqib A Kakvi
20  * @version 2.0A
21  * @since 2.0A
22  * @see SKCards.Card Card
23  * @SK.bug system does not consider the case when the Cards are equal.
24  */
25 public class HighLow{
26
27     public final static int HIGH = 1;
28     public final static int LOW = 0;
29
30     double[] highv, highp;
31     double[] lowv, lowp;
32     double tau;
33     Pack pack;
34     int cscore, hscore;
35
36 }
```

```

37  /**
38   * The main method. Creates a new HighLow.
39   */
40  public static void main(String[] args){
41      new HighLow();
42  }
43
44
45  /**
46   * Default constructor. Creates a new learning agent, with no previous
47   * knowledge. It begins learning for 10,000 episodes. After it has
48   * it has learnt, it plays against the user. As part of the testing,
49   * the values of the learning are written to seperate files.
50   */
51  public HighLow(){
52
53      pack = new Pack(2, true);
54      highv = new double[14];
55      highp = new double[14];
56      lowv = new double[14];
57      lowp = new double[14];
58      tau = 0.1;
59
60      for(int i = 0; i < highv.length; i++){
61          highv[i] = 0.0;
62          highp[i] = 0.0;
63          lowv[i] = 0.0;
64          lowp[i] = 0.0;
65      }
66
67      System.out.println("*****LEARNING*****");
68      long start = System.currentTimeMillis();
69      try{
70
71          BufferedWriter[] wr = new BufferedWriter[14];
72
73          for(int i = 0; i < wr.length; i++){
74
75              String name = "values_" + i + ".csv";
76
77              wr[i] = new BufferedWriter(new FileWriter(name));
78              wr[i].write("low" + Integer.toString(i));
79              wr[i].write("\t");
80              wr[i].write("high" + Integer.toString(i));
81              wr[i].newLine();
82          }
83          for(int i = 0; i < 10000; i++){
84
85              episode();
86
87              for(int j = 0; j < wr.length; j++){
88
89                  wr[j].write(Double.toString(lowv[j]));
90                  wr[j].write("\t");

```

```

91         wr[j].write(Double.toString(highv[j]));
92         wr[j].newLine();
93     }
94 }
95
96     for(int i = 0; i < wr.length; i++){
97         wr[i].flush();
98         wr[i].close();
99     }
100 }
101 catch(Exception ex){
102 }
103 long stop = System.currentTimeMillis();
104 long time = stop - start;
105 System.out.println("*****DONE_LEARNING*****");
106 System.out.println("*****TOOK_" + time + "_ms_TO_LEARN*****");
107 try{
108     BufferedWriter bw = new BufferedWriter(new FileWriter("values.txt"));
109
110     bw.write("num" + "\t" + "high" + "\t" + "low");
111     bw.newLine();
112
113     for(int i = 0; i < highv.length; i++){
114
115         bw.write(i + "\t" + highv[i] + "\t" + lowv[i]);
116         bw.newLine();
117     }
118
119     bw.flush();
120     bw.close();
121 }
122 catch(Exception ex){
123     ex.printStackTrace();
124 }
125 System.out.println();
126 System.out.println("NOW_LETS_PLAY!!!");
127 System.out.println("You_need_to_guess_if_the_next_card_is_higher_or" +
128     "lower_than_the_current_card");
129 System.out.println("Enter_1_for_higher,_0_for_lower,_" +
130     "equals_is_not_an_option!");
131
132 cscore = 0;
133 hscore = 0;
134
135 try{
136
137     BufferedReader br = new BufferedReader(new InputStreamReader(System.in));
138
139     for(int i = 1 ; i < 10; i++){
140         Card card1 = pack.nextCard();
141         Card card2 = pack.nextCard();
142
143         int c1 = card1.getPoints();
144         int c2 = card2.getPoints();

```

```

145
146         int caxn;
147
148         if (highv[c1] > lowv[c1]){
149             caxn = 1;
150         }
151         else{
152             caxn = 0;
153         }
154
155         System.out.println("Number_" + i + ":_ " + card1.toString());
156         System.out.print("Enter_your_guess:_");
157         int haxn = Integer.parseInt(br.readLine());
158         System.out.println("I_guessed_" + caxn);
159         cplay(caxn,c1,c2);
160         hplay(haxn,c1,c2);
161         System.out.println("The_actual_card_is:_ " + card2.toString());
162     }
163
164     System.out.println("YOU_SCORED:_ " + hscore);
165     System.out.println("I_SCORED:_ " + cscore);
166 }
167 catch(Exception ex){
168 }
169 }
170
171
172 /**
173  * This method goes through a single episode of the learning. This
174  * involves the agent picking an action based on the value of the
175  * card. The selected action is then taken.
176  */
177 public void episode(){
178
179
180     Card card1 = pack.nextCard();
181     Card card2 = pack.nextCard();
182
183     int c1 = card1.getPoints();
184     int c2 = card2.getPoints();
185
186     //Jo < A < 2 < 3 < 4 < 5 < 6 < 7 < 8 < 9 < T < J < Q < K
187     //0 < 1 < 2 < 3 < 4 < 5 < 6 < 7 < 8 < 9 < 10 < 11 < 12 < 13
188
189     double hv = Math.exp(highv[c1]/tau);
190     double lv = Math.exp(lowv[c2]/tau);
191     double sum = hv + lv;
192
193     highp[c1] = hv/sum;
194     lowp[c1] = lv/sum;
195
196     double r = Math.random();
197
198     if (r > lowp[c1]){

```

```

199         action(HIGH,c1,c2);
200     }
201     else{
202         action(LOW,c1,c2);
203     }
204 }
205
206
207 /**
208  * Carries out a single action, high or low, based on the parameters
209  * recieved from the episode method. The values of the two cards are
210  * also passed for verification. The score is adjusted accordingly,
211  * based on the values of the cards and the axn taken.
212  *
213  * @param axn the axn to be taken
214  * @param c1 the value of Card1
215  * @param c2 the value of Card2
216  */
217 public void action(int axn, int c1, int c2){
218
219     if(axn == HIGH){
220         if(c1 < c2){
221             highv[c1] += 1;
222         }
223         else{
224             highv[c1] -= 1;
225         }
226     }
227     else if(axn == LOW){
228         if(c1 > c2){
229             lowv[c1] += 1;
230         }
231         else{
232             lowv[c1] -= 1;
233         }
234     }
235 }
236
237
238 /**
239  * This method is used for a single guess made by the computer in the
240  * head to head game with the human. If the computer guessed correctly
241  * it gets a point.
242  *
243  * @param axn the axn to be taken
244  * @param c1 the value of Card1
245  * @param c2 the value of Card2
246  */
247 public void cplay(int axn, int c1, int c2){
248
249     if(axn == 1){
250         if(c1 < c2){
251             cscore++;
252         }

```

```

253     }
254     else if (axn == 0){
255         if (c1 > c2){
256             cscore++;
257         }
258     }
259 }
260
261
262 /**
263  * This method is used for a single guess made by the human in the
264  * head to head game with the computer. If the human guessed correctly
265  * it gets a point.
266  *
267  * @param axn the axn to be taken
268  * @param c1 the value of Card1
269  * @param c2 the value of Card2
270  */
271 public void hplay(int axn, int c1, int c2){
272
273     if (axn == 1){
274         if (c1 < c2){
275             hscore++;
276         }
277     }
278     else if (axn == 0){
279         if (c1 > c2){
280             hscore++;
281         }
282     }
283 }
284 }

```

A.2 Results

This section contains the results of one run of the High/Low prototype. Several runs were conducted and the results were found to be fairly consistent. The reason for variations in the results was the fact that no proper learning algorithm was used, hence there was no supporting proof of convergence. The following graphs illustrate the values of guessing high or low for a given card, with respect to the number of occurrences.

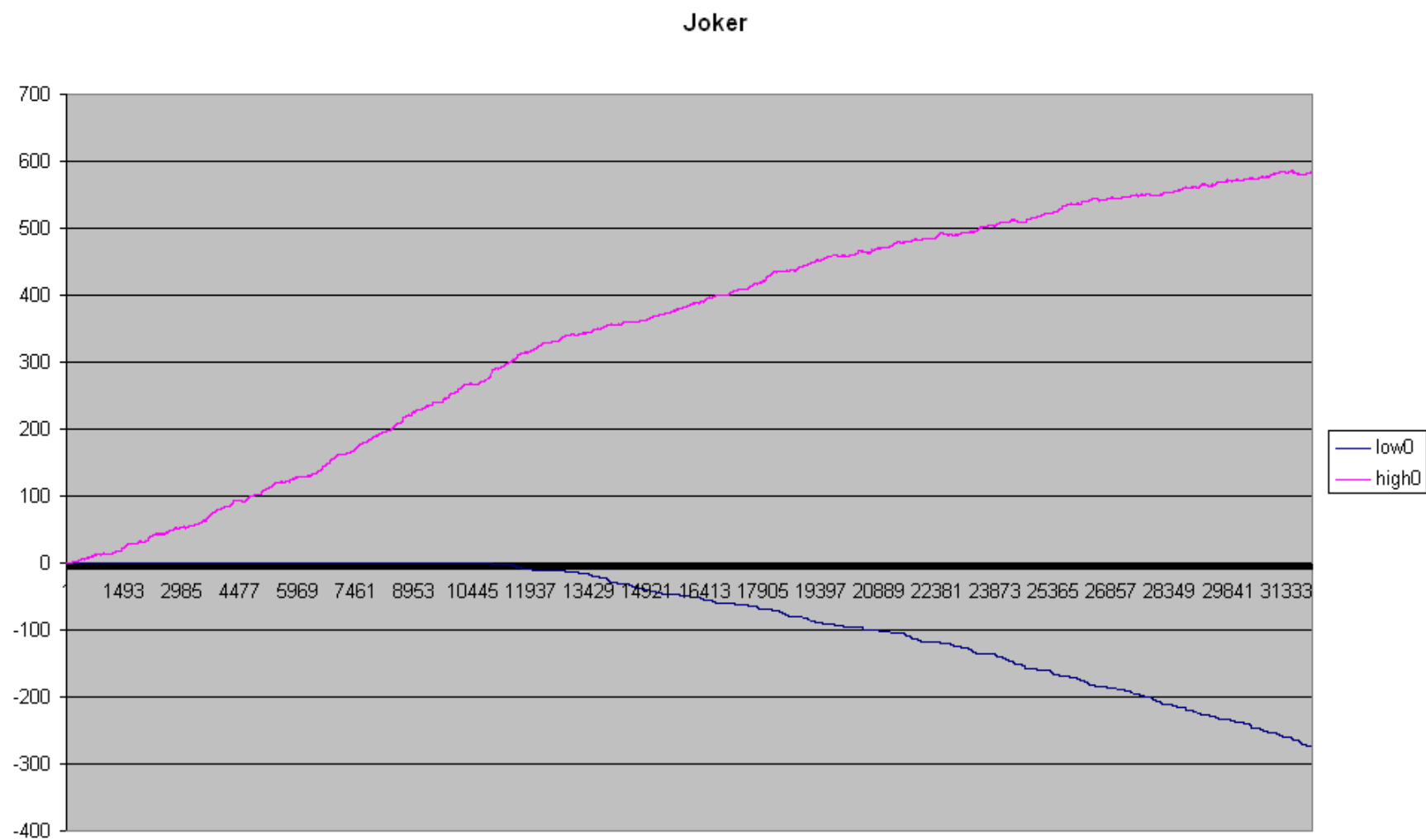


Figure A.1: Values for Joker

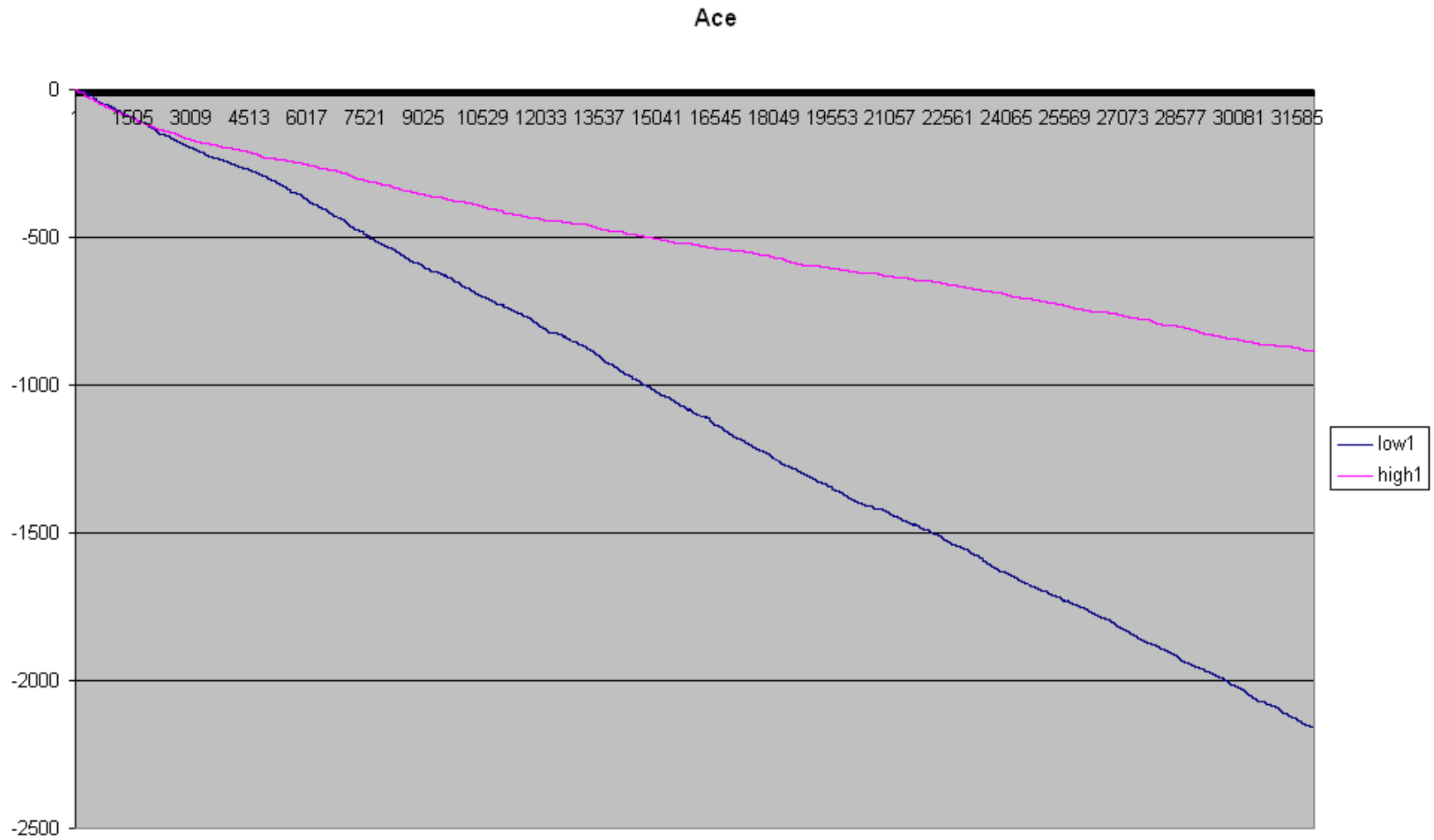


Figure A.2: Values for Ace

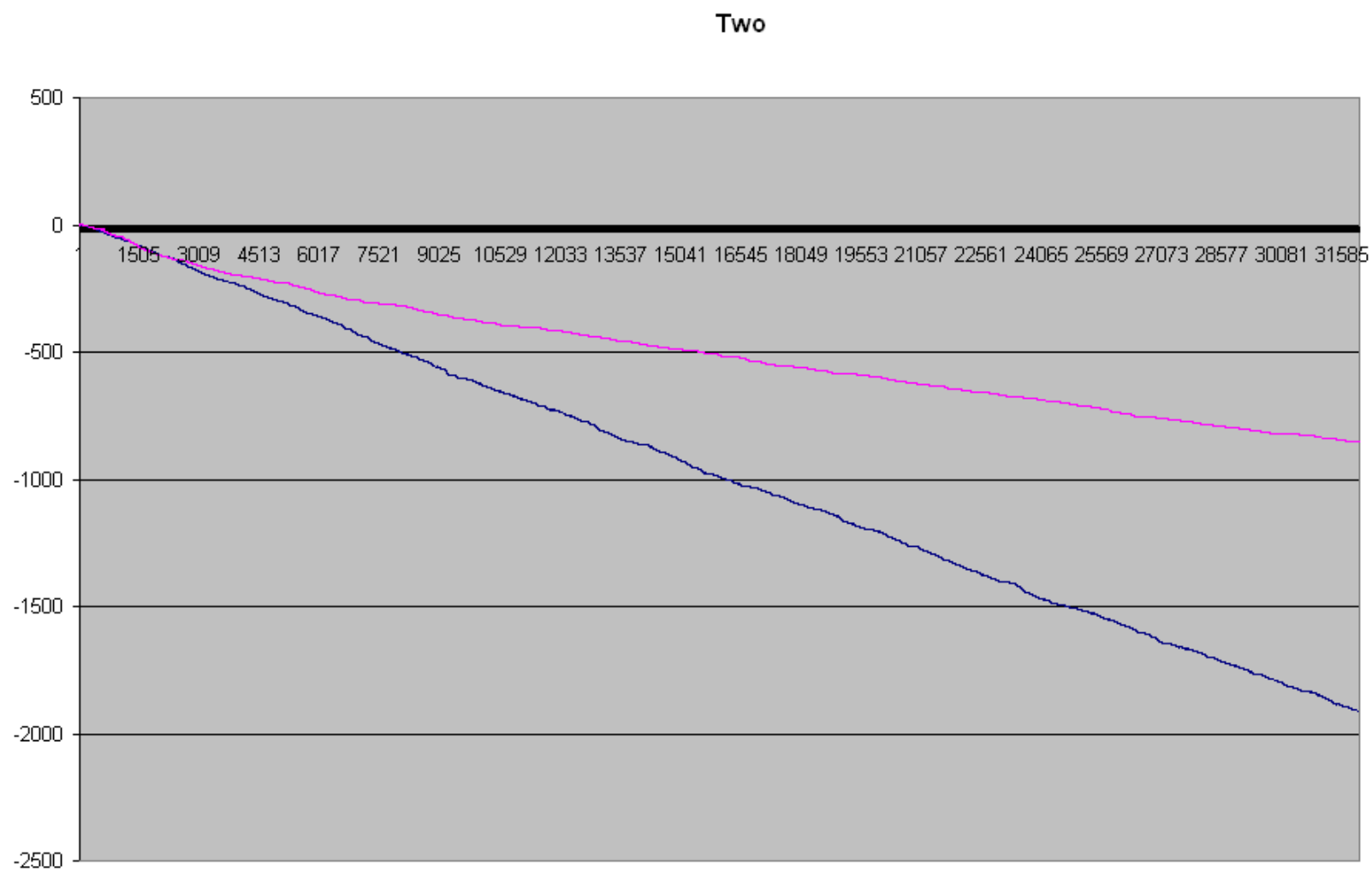


Figure A.3: Values for 2

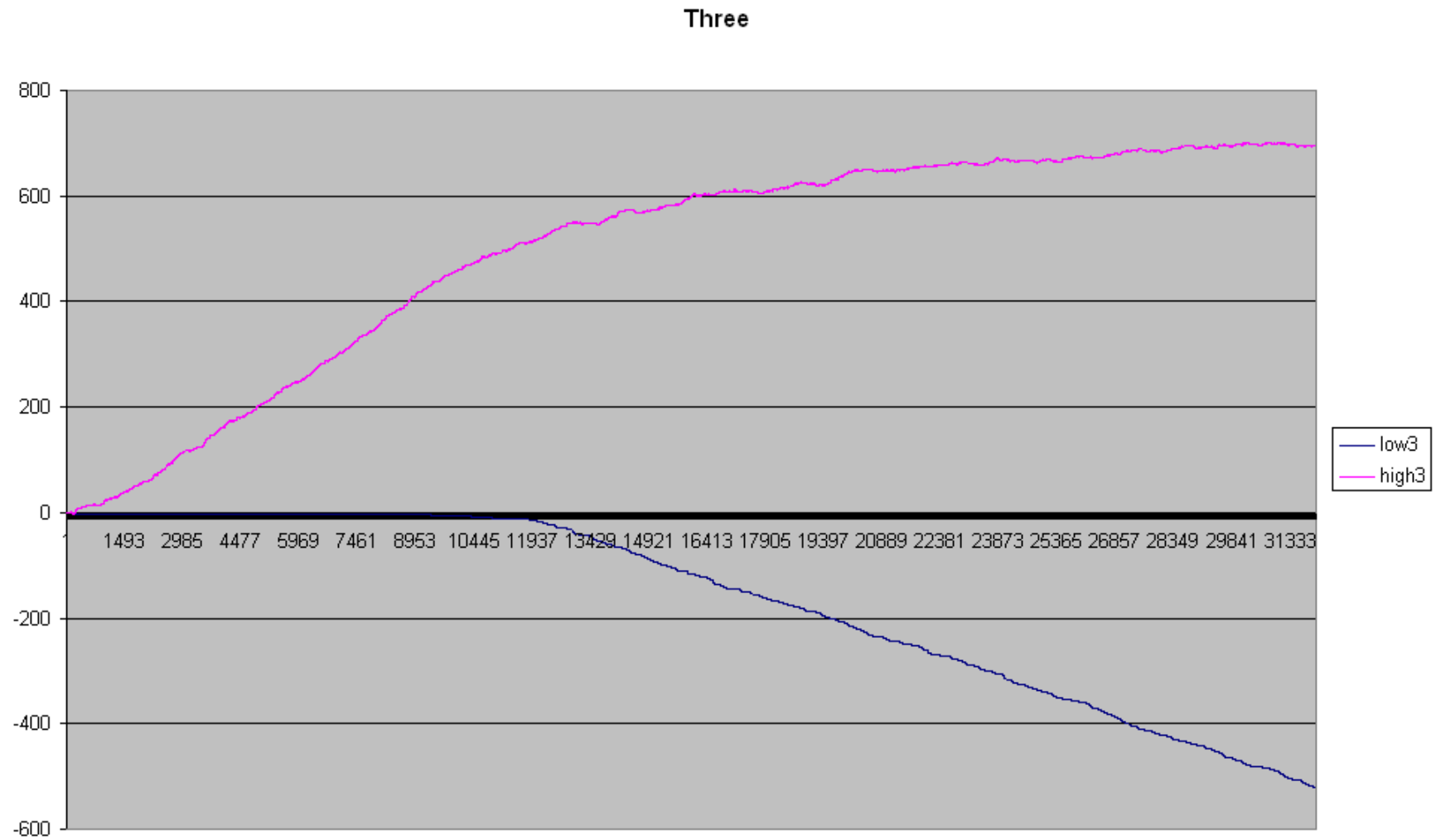


Figure A.4: Values for 3

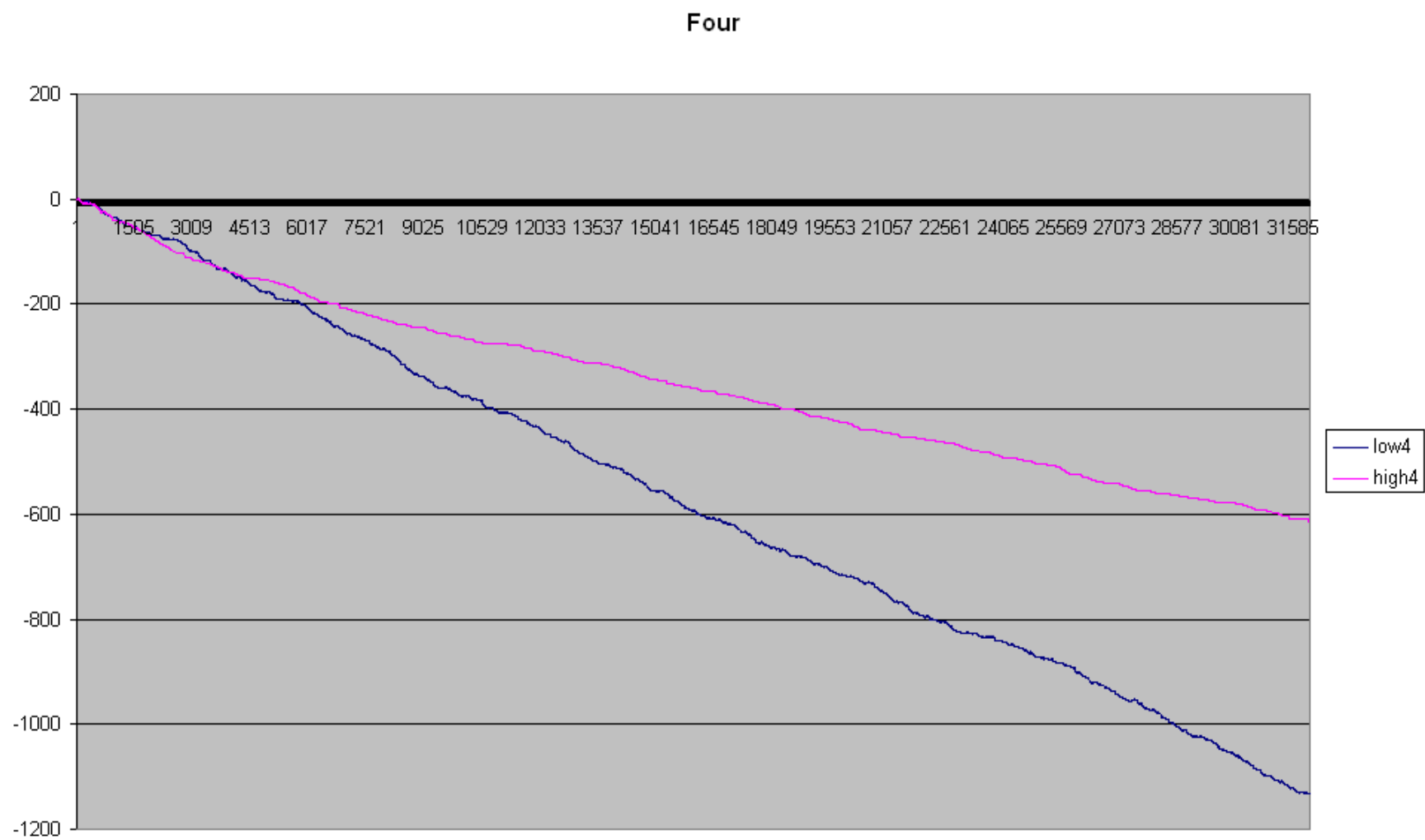


Figure A.5: Values for 4

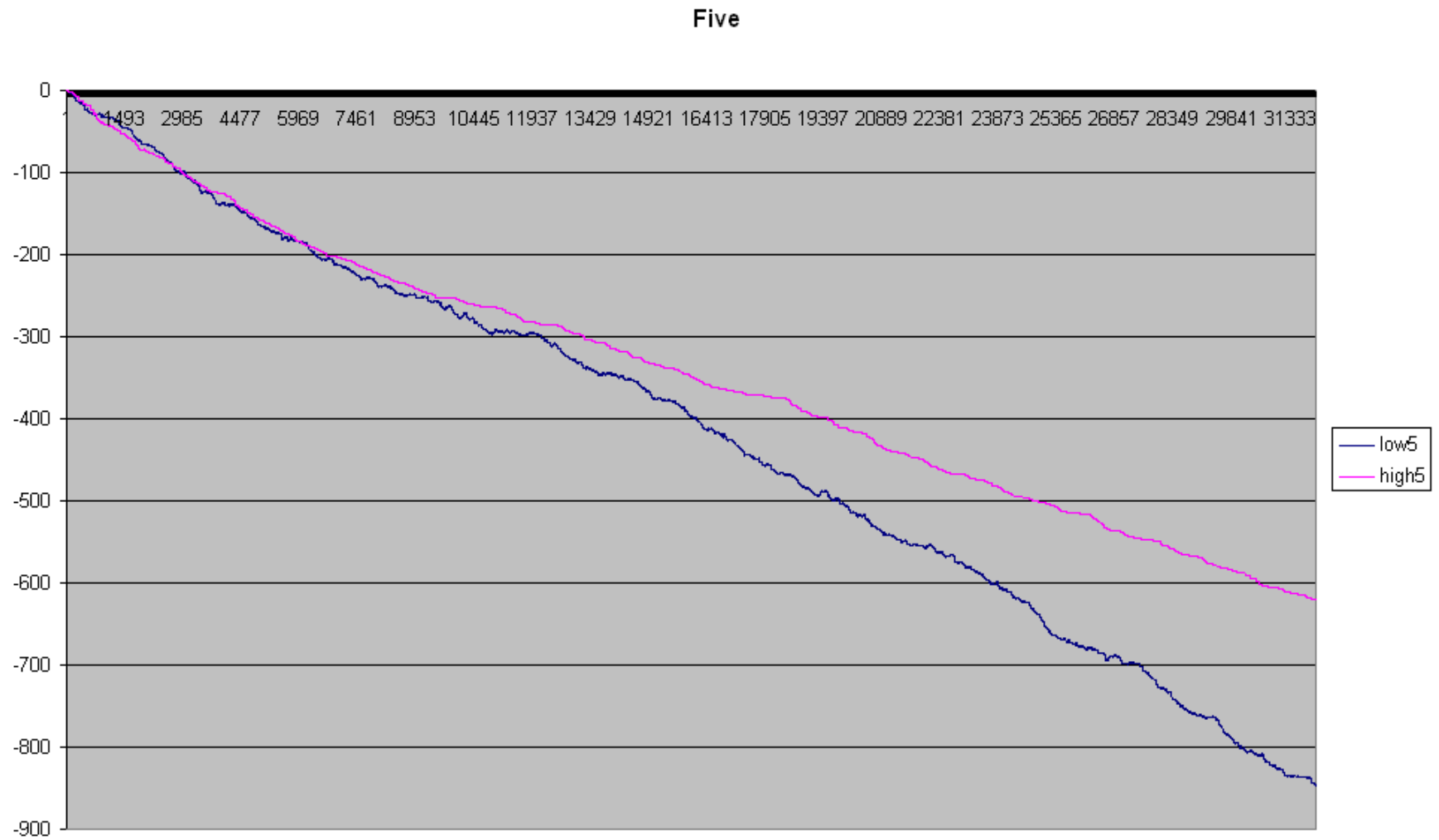


Figure A.6: Values for 5

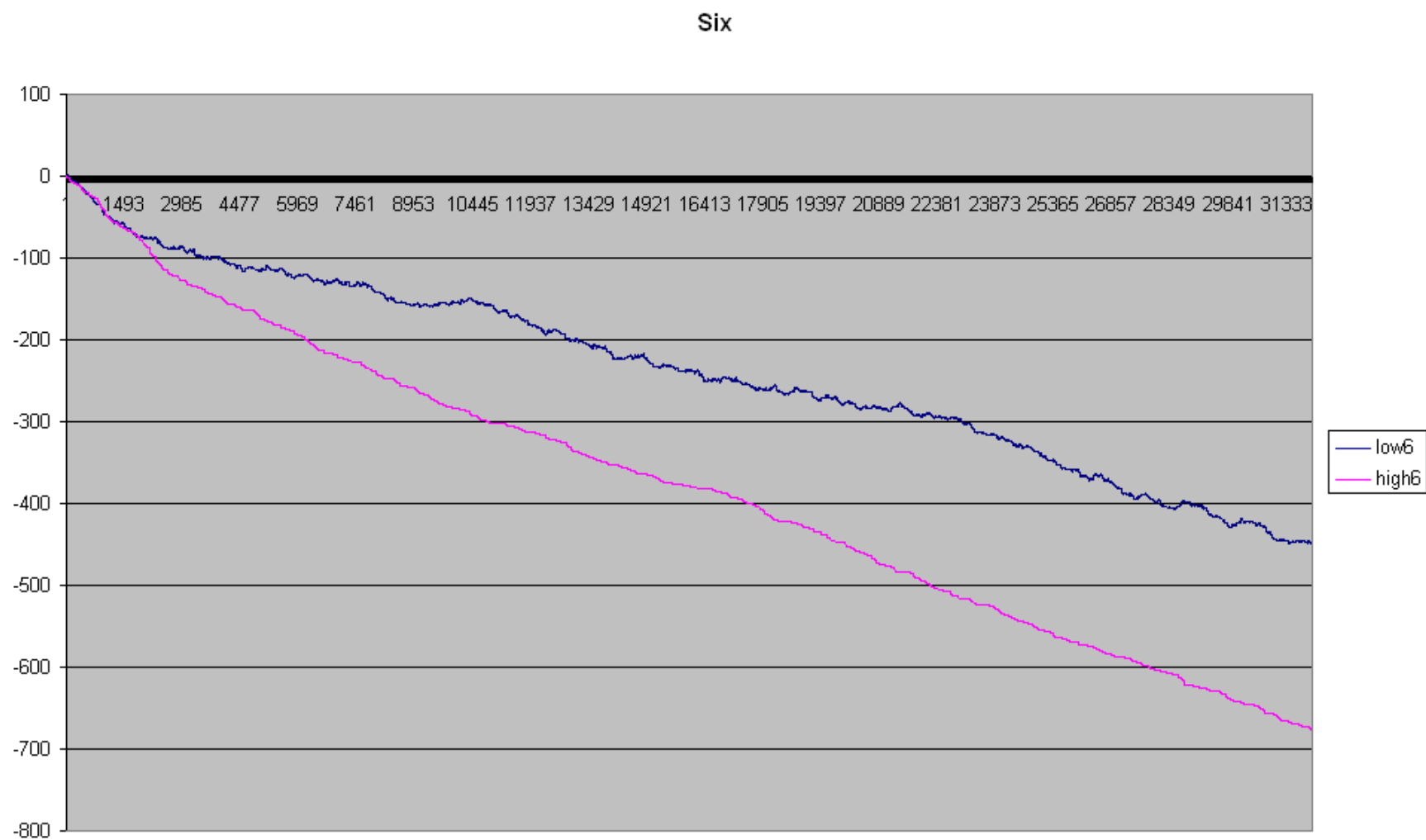


Figure A.7: Values for 6

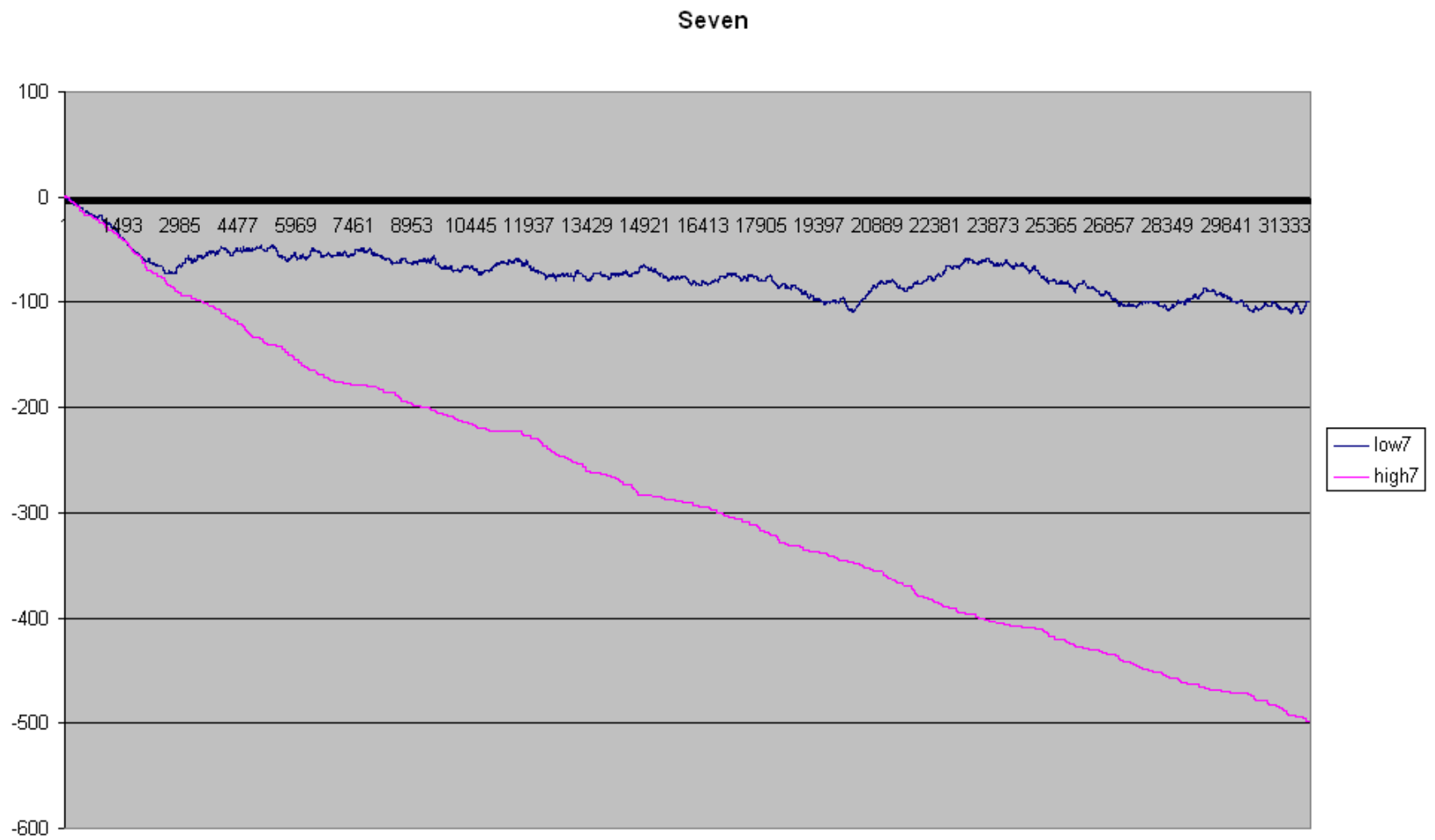


Figure A.8: Values for 7

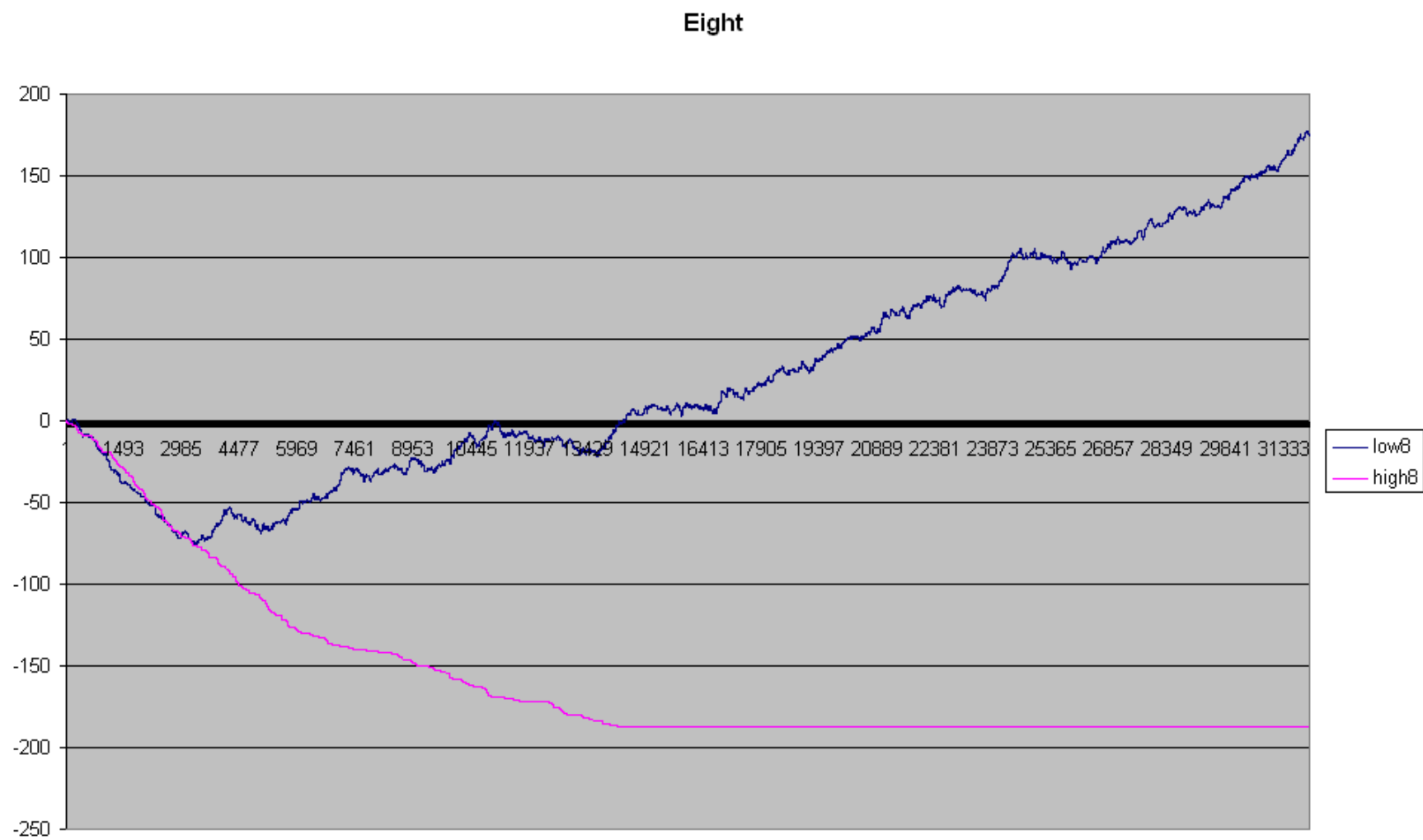


Figure A.9: Values for 8

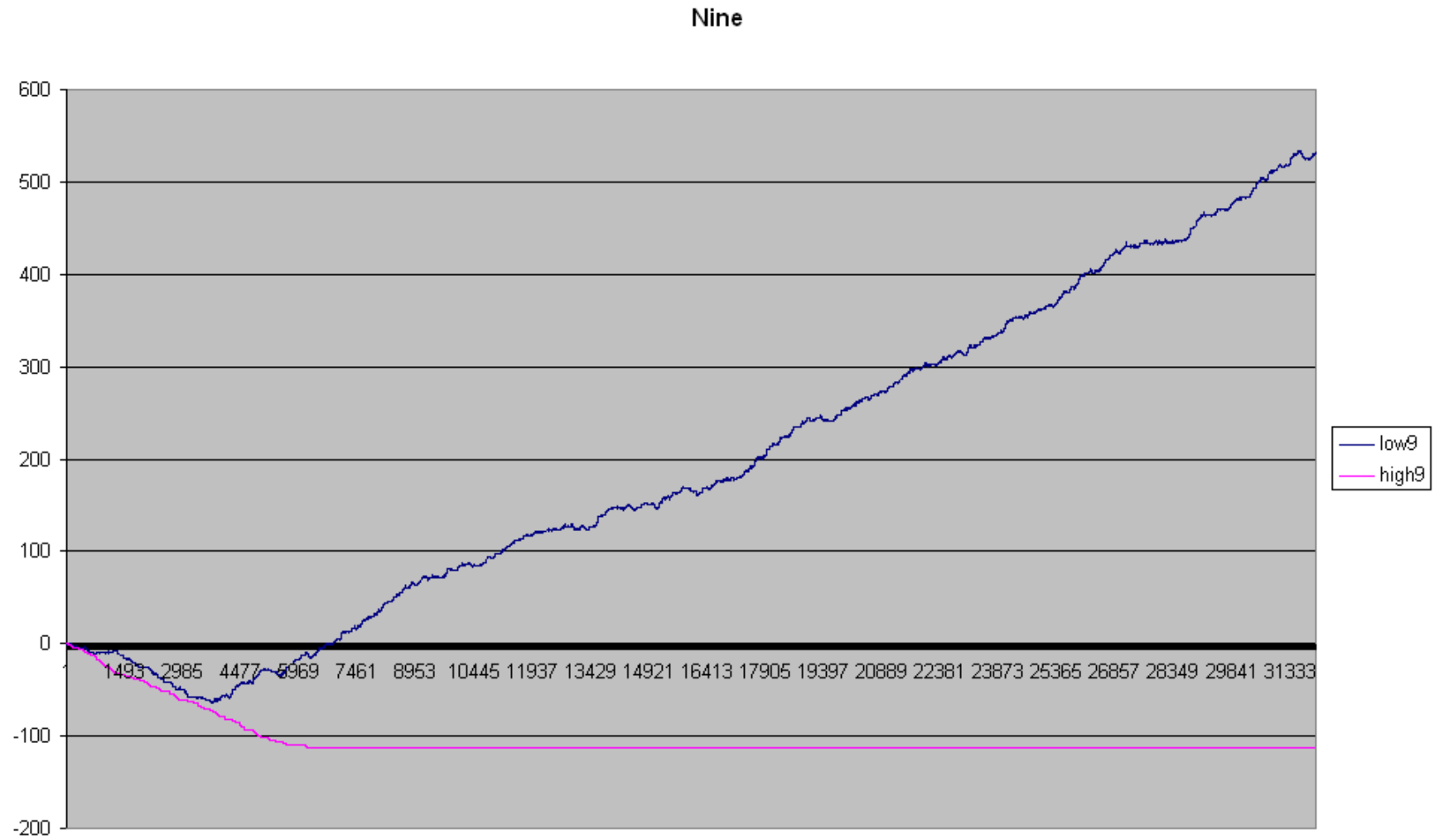


Figure A.10: Values for 9

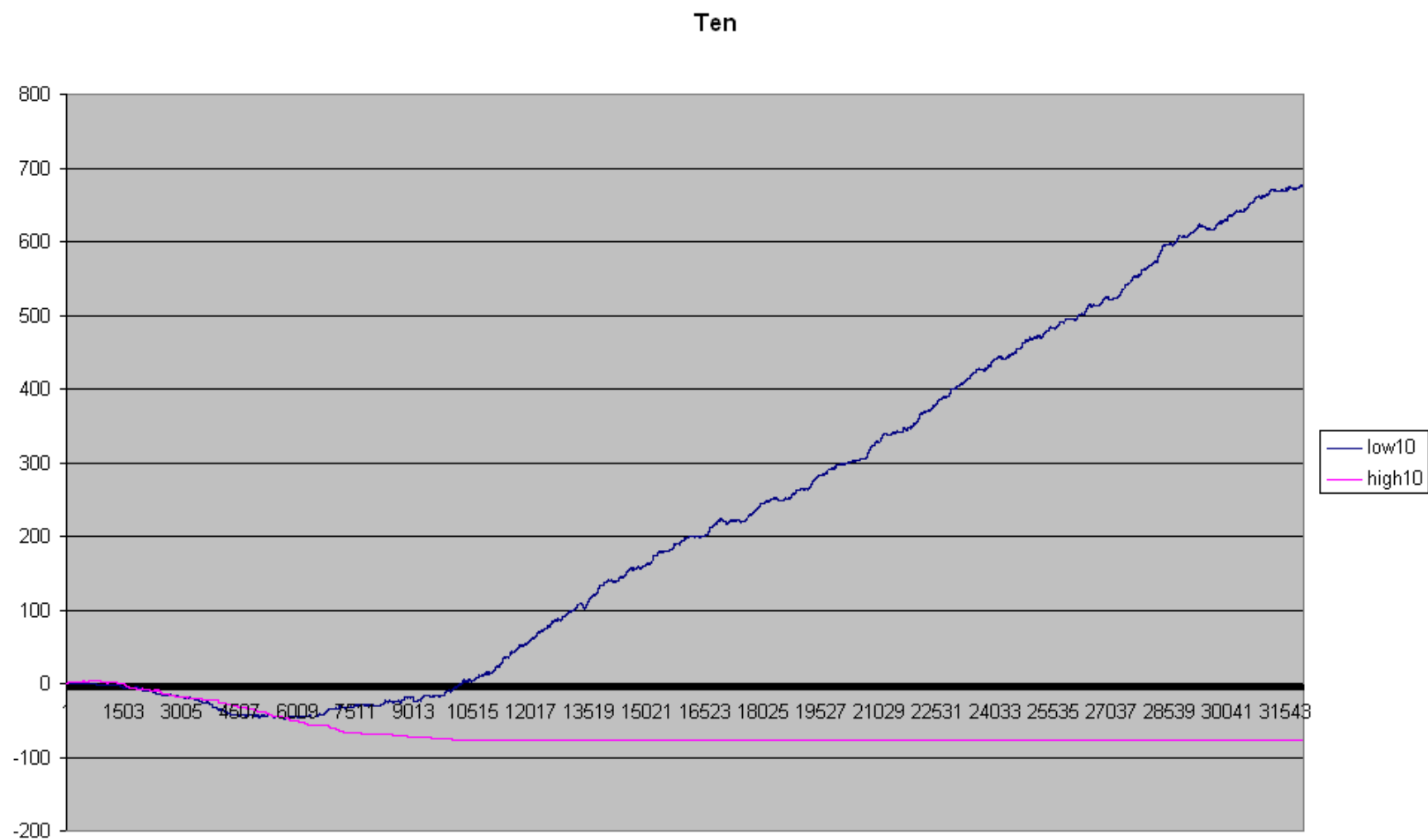


Figure A.11: Values for 10

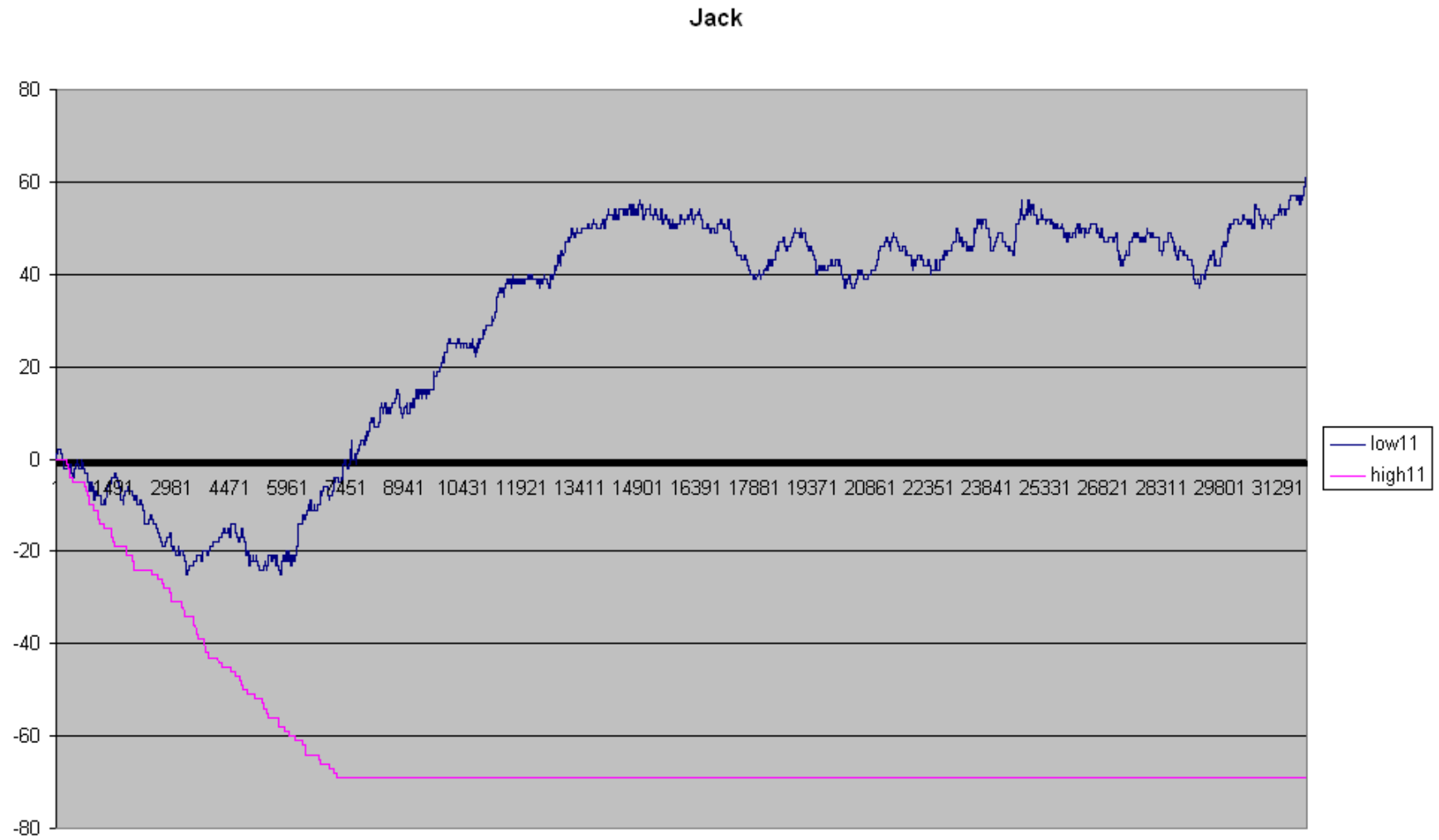


Figure A.12: Values for Jack

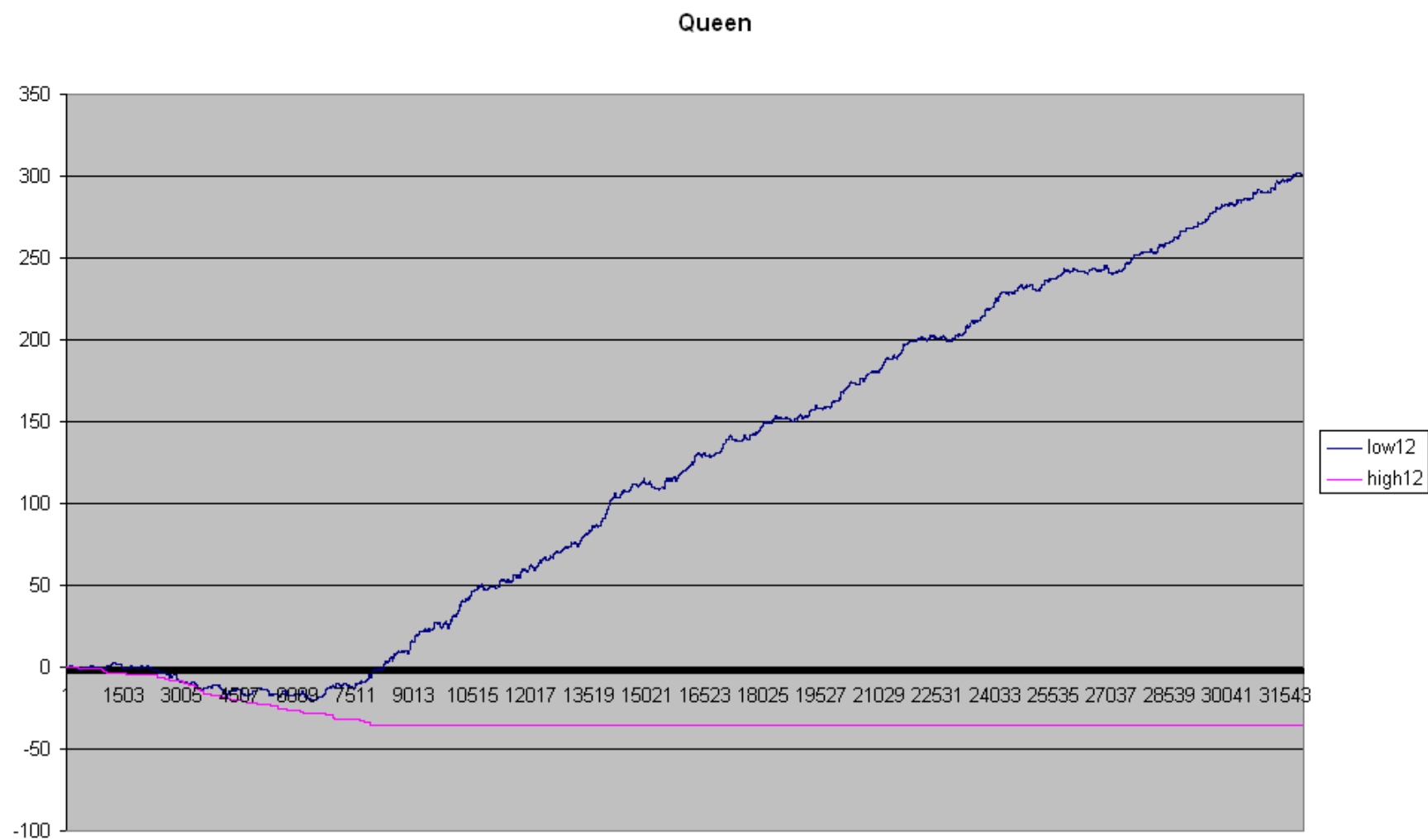


Figure A.13: Values for Queen

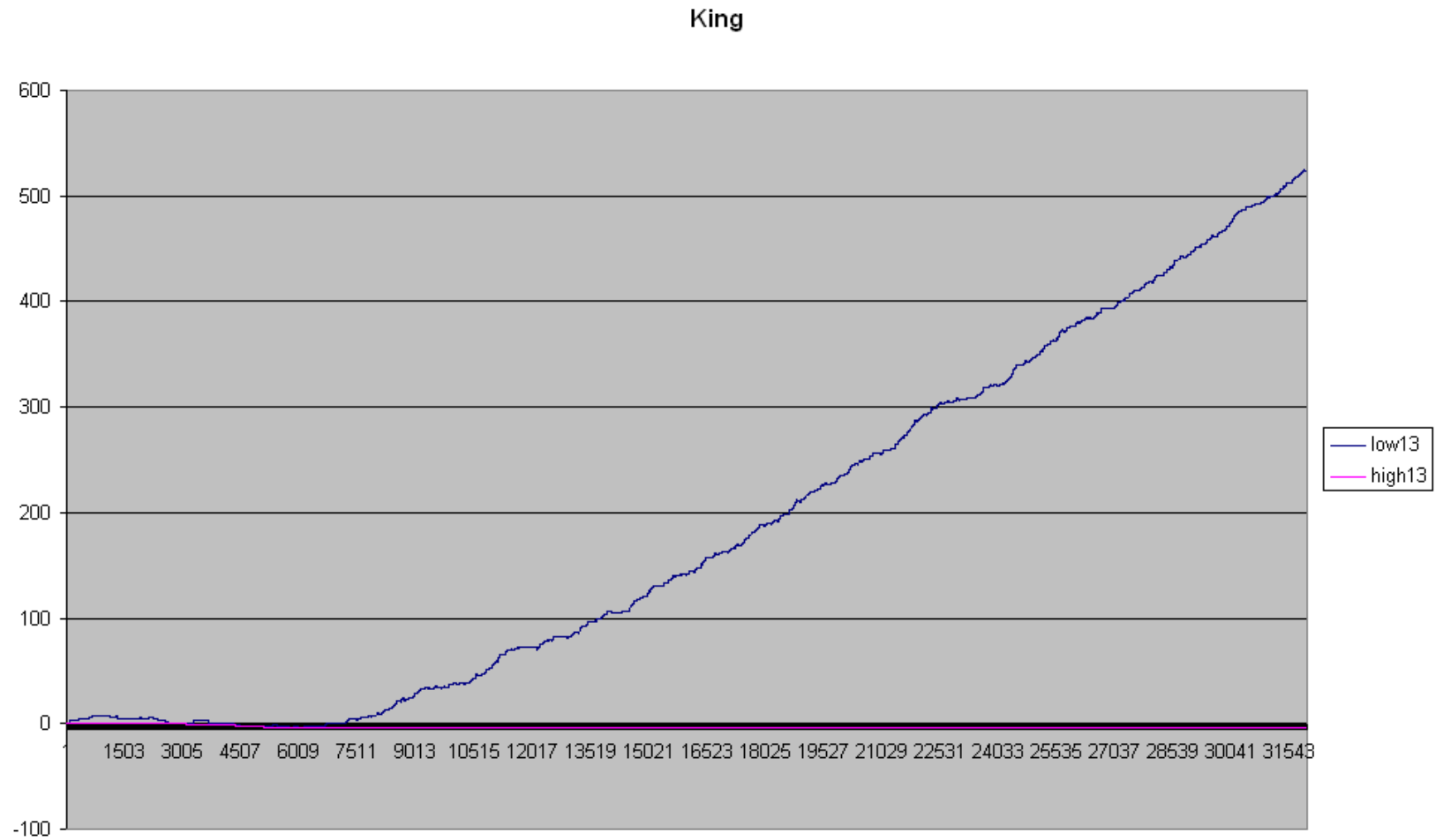


Figure A.14: Values for King

VALUE	FINAL POLICY
Joker	HIGH
2	HIGH
3	HIGH
4	HIGH
5	HIGH
6	HIGH
7	LOW
8	LOW
9	LOW
10	LOW
Jack	LOW
Queen	LOW
King	LOW

Table A.1: Final Policy

As we can see, The agent learns policies at the extremities (Joker and King) very easily. However with the middle values, it has some problems. Using a proper learning algorithm would help smooth out the learning in the middle values. It can be seen from Table A.1 below that the agent has learnt a fairly good policy.

Appendix B

GridWorld prototype

B.1 Source Code

```
1  //Adding GridWorld to SKCards
2  package SKCards;
3
4  /*
5   * importing the ArrayList class
6   * and the swing/awt components
7   * and the border components
8   * and the BufferedImage class
9   * and the IO components
10  * and the ImageIO class
11  */
12  import java.util.ArrayList;
13  import java.awt.*;
14  import javax.swing.*;
15  import javax.swing.border.*;
16  import java.awt.image.BufferedImage;
17  import java.io.*;
18  import javax.imageio.ImageIO;
19
20
21  /**
22   * This is the second prototype for the SKCards system. This
23   * class attempts to solve the GridWorld problem, using the
24   * Reinforcement Learning Monte Carlo methods.
25   *
26   * @author Saqib A Kakvi
27   * @version 2.0A
28   * @since 2.0A
29   * @SK.bug none
30   */
31  public class GridWorld extends JFrame{
32
33      private JLabel[][] grid;
34      private Point loc;
35      private double tau, gamma, alpha;
36      private double[][] val, prob;
```

```

37     private boolean[][] pen;
38     private ArrayList<Point> path;
39     private JFrame values;
40     private JLabel[][] vlab;
41     private String filename;
42     private int h,w, p;
43
44     /**
45      * main method. Creates a new instance of GridWorld, using the default
46      * constructor.
47      *
48      * @param args not used
49      */
50     public static void main(String [] args){new GridWorld().run();}
51
52
53     /**
54      * Default constructor. Creates a GridWorld instance b passing the
55      * parameters 10 and 10 to the main constructor.
56      *
57      * @see GridWorld
58      */
59     public GridWorld(){ this (10,10);}
60
61
62     /**
63      * Main Constructor. Creates a grid of size <i>n</i>x<i>m</i>.
64      * All the grid squares are created and penalties assigned to the
65      * appropriate squares. The system is then run for 2500 episodes.
66      *
67      * @param n the length of the grid
68      * @param m the width of the drid
69      */
70     public GridWorld(int n, int m){
71
72         super("GridWorld");
73         Dimension d = (Toolkit.getDefaultToolkit()).getScreenSize();
74         h = (int)(d.getHeight()) / 4 * 3;
75         w = (int)(d.getWidth()) / 2;
76
77         setSize(w,h);
78         setDefaultCloseOperation(DISPOSE_ON_CLOSE);
79
80         grid = new JLabel[n][m];
81         tau = 0.75;
82         gamma = 0.75;
83         alpha = 0.5;
84         val = new double[n][m];
85         pen = new boolean[n][m];
86
87         for(int i = 0 ; i < n; i++){
88             for(int j = 0; j < m; j++){
89                 grid[i][j] = new JLabel();
90                 grid[i][j].setBorder(new LineBorder(Color.BLUE));

```



```

91         grid[i][j].setHorizontalAlignment(SwingConstants.CENTER);
92         int r = new Double(Math.random() * 10).intValue();
93         if ((i+j == ((n+m)/2)-1) && (i != n/2)){
94             pen[i][j] = true;
95             grid[i][j].setForeground(Color.RED);
96         }
97         else{
98             pen[i][j] = false;
99         }
100     }
101 }
102
103 Container cp = getContentPane();
104 cp.setLayout(new GridLayout(n,m));
105
106 for(int i = 0 ; i < n; i++){
107     for(int j = 0; j < m; j++){
108         cp.add(grid[i][j]);
109     }
110 }
111
112 grid[0][0].setForeground(Color.GREEN);
113 grid[n-1][m-1].setForeground(Color.GREEN);
114
115 for(int i = 0; i < val.length; i++){
116     for(int j = 0; j < val[i].length; j++){
117         if(pen[i][j]){
118             val[i][j] = -100;
119         }
120         else{
121             val[i][j] = 0;
122         }
123     }
124 }
125
126 val[0][0] = 100;
127 val[n-1][m-1] = 100;
128
129 values = new JFrame("VALUES");
130 values.setSize(w,h);
131 values.setLocation(w,0);
132 values.setDefaultCloseOperation(DISPOSE_ON_CLOSE);
133
134 Container c = values.getContentPane();
135
136 c.setLayout(new GridLayout(n, m));
137
138 vlab = new JLabel[n][m];
139
140 for(int i = 0; i < vlab.length; i++){
141     for(int j = 0; j < vlab[i].length; j++){
142
143         vlab[i][j] = new JLabel(Double.toString(val[i][j]));
144         vlab[i][j].setHorizontalAlignment(SwingConstants.CENTER);

```

```

145         vlab[i][j].setOpaque(true);
146         vlab[i][j].setForeground(Color.WHITE);
147         vlab[i][j].setBackground(Color.BLACK);
148         c.add(vlab[i][j]);
149     }
150 }
151
152     reset();
153     this.setVisible(true);
154     values.setVisible(true);
155     this.repaint();
156     values.repaint();
157 }
158
159 /**
160  * Runs the GridWorld program.
161  */
162 public void run(){
163
164     for(int i = 1; i <= 2500; i++){
165         filename = new String("TEST" + Integer.toString(i) + ".jpg");
166         setTitle("Episode_" + i);
167         episode();
168         reset();
169         screenShot();
170     }
171     System.out.println(p + "_PENALTIES");
172 }
173
174
175 /**
176  * Sets the GridWorld and the values frame's visibility using
177  * the value passed.
178  *
179  * @param vis the visibility of the windows.
180  */
181 public void setVisibility(boolean vis){
182
183     this.setVisible(vis);
184     values.setVisible(vis);
185 }
186
187
188 /**
189  * This method runs a single episode through the system. It
190  * continues to move a single step until it reaches a terminal
191  * square. When this happens, a screenshot of the windows is taken.
192  *
193  * @see step
194  */
195 private void episode(){
196
197     int x = (int)((Math.random() * 100) % grid.length);
198     int y = (int)((Math.random() * 100) % grid[0].length);

```

```

199
200     loc = new Point(x,y);
201
202     while(checkloc() != 0){
203         x = (int)((Math.random() * 100) % grid.length);
204         y = (int)((Math.random() * 100) % grid[0].length);
205
206         loc = new Point(x,y);
207     }
208
209     path = new ArrayList<Point>();
210
211     grid[x][y].setText("X");
212     while(checkloc() == 0){
213         step();
214     }
215
216     reward();
217 }
218
219
220 /**
221  * Testing method. Used to take a screenshot of the system
222  * after each episode. This allows the progress and results
223  * of the system to be documented.
224  */
225 private void screenShot(){
226     //CODE TAKEN FROM
227     //www.java-tips.org/java-se-tips/java.awt/how-to-capture-screenshot.html
228     try {
229         Robot robot = new Robot();
230         //Capture the screen shot of the area of the screen
231         //defined by the rectangle
232         BufferedImage bi=robot.createScreenCapture(new Rectangle(w*2,h));
233         ImageIO.write(bi, "jpg", new File(filename));
234     }
235     catch (Exception ex) {
236         ex.printStackTrace();
237     }
238 }
239
240
241 /**
242  *
243  */
244 private void reset(){
245
246     for(int i = 0; i < grid.length; i++){
247         for(int j = 0; j < grid[i].length; j++){
248             if(pen[i][j]){
249                 grid[i][j].setText("PEN");
250             }
251             else{
252                 grid[i][j].setText("");

```

```

253     }
254     }
255 }
256
257 grid[0][0].setText("GOAL");
258 grid[grid.length-1][grid[0].length-1].setText("GOAL");
259
260 for(int i = 0; i < vlab.length; i++){
261     for(int j = 0; j < vlab[i].length; j++){
262         vlab[i][j].setText(Double.toString(val[i][j]));
263         if(val[i][j] > 0){
264             int v = (int)val[i][j] * 255 / 100;
265             vlab[i][j].setBackground(new Color(0, v, 0));
266         }
267         else{
268             int v = (int)val[i][j] * 255 / 100 * -1;
269             vlab[i][j].setBackground(new Color(v, 0, 0));
270         }
271     }
272 }
273
274 repaint();
275 }
276
277
278 /**
279  *
280  */
281 private int checkloc(){
282
283     if(getx() == 0 && gety() == 0){
284         return 1;
285     }
286     else if(getx() == grid.length-1 && gety() == grid[0].length-1){
287         return 1;
288     }
289     else if(pen[getx()][gety()]){
290         return -1;
291     }
292     else{
293         return 0;
294     }
295 }
296
297
298 /**
299  * This method moves the agent a single step. The direction
300  * is decided on using the softmax probability distributions.
301  *
302  * @see action(int)
303  */
304 private void step(){
305
306     double lim = 0.0;

```

```

307     double sum = 0.0;
308     int x = getx();
309     int y = gety();
310     double[] prob = new double[4];
311
312     double up,down,left,right;
313
314     boolean uv = true, lv = true, dv = true, rv = true;
315
316     if(getx() == 0){
317         uv = false;
318     }
319     else if(getx() == grid.length-1){
320         dv = false;
321     }
322
323     if(gety() == grid[0].length-1){
324         rv = false;
325     }
326     else if(gety() == 0){
327         lv = false;
328     }
329
330     if(uv){
331         up = Math.exp(val[x-1][y]/tau);
332         sum += up;
333     }
334     if(rv){
335         right = Math.exp(val[x][y+1]/tau);
336         sum += right;
337     }
338     if(dv){
339         down = Math.exp(val[x+1][y]/tau);
340         sum += down;
341     }
342     if(lv){
343         left = Math.exp(val[x][y-1]/tau);
344         sum += left;
345     }
346
347     if(uv){
348         prob[0] = Math.exp(val[x-1][y]/tau) / sum;
349     }
350     if(rv){
351         prob[1] = Math.exp(val[x][y+1]/tau) / sum;
352     }
353     if(dv){
354         prob[2] = Math.exp(val[x+1][y]/tau) / sum;
355     }
356     if(lv){
357         prob[3] = Math.exp(val[x][y-1]/tau) / sum;
358     }
359
360     double r = Math.random();

```

```

361
362     for(int i = 0; i < prob.length; i++){
363
364         lim += prob[i];
365         if(r < lim){
366             action(i);
367             break;
368         }
369     }
370
371     repaint();
372 }
373
374
375 /**
376  * Takes an action, that is moving 1 step in any of the
377  * valid directions.
378  *
379  * @param axn the direction to move.
380  */
381 private void action(int axn){
382
383     path.add(new Point(loc));
384
385     switch(axn){
386         case 0: //up
387             loc.translate(-1,0);
388             break;
389         case 1: //right
390             loc.translate(0,1);
391             break;
392         case 2: //down
393             loc.translate(1,0);
394             break;
395         case 3: //left
396             loc.translate(0,-1);
397             break;
398     }
399
400     grid[getx()][gety()].setText("X");
401 }
402
403
404 /**
405  *
406  */
407 private void reward(){
408
409     if(checkloc() == -1) p++;
410     double reward = val[getx()][gety()];
411
412     for(int i = path.size()-1; i >=0; i--){
413
414         reward *= gamma;

```

```

415
416         Point p = path.get(i);
417         int x = (int)p.getX();
418         int y = (int)p.getY();
419         val[x][y] += alpha * (reward - val[x][y]);
420
421         reward = val[x][y];
422     }
423 }
424
425
426 /**
427  * Returns the x co-ordinate of the agent's current
428  * position, as an int so it can be used as an
429  * index in the grid array.
430  *
431  * @return the x co-ordinate
432  * @see gety
433  */
434 private int getX(){ return (int)loc.getX();}
435
436
437 /**
438  * Returns the y co-ordinate of the agent's current
439  * position, as an int so it can be used as an
440  * index in the grid array.
441  *
442  * @return the y co-ordinate
443  * @see getX
444  */
445 private int getY(){ return (int)loc.getY();}
446 }

```

B.2 Results

The GridWorld problem can be said to be mathematically trivial, as the values can be calculated very easily. Given the final reward r and the discount rate γ , we can calculate the value v of a square n steps away from the goal using the formula:

$$v = r \times \gamma^n$$

Based on this formula, we can calculate the values of all the grid squares based on their distance from the goal square as given in the table below:

Included below are some screen shots taken during the execution of the program. These screen shots selected were before episodes 1, 25, 50, 100, 150, 200, 250, 500, 750, 1000, 1500 & 2500.

STEPS TO GOAL	VALUE
0	100
1	75
2	56.25
3	42.1875
4	31.6406
5	23.7305
6	17.7979
7	13.3484
8	10.0113
9	7.5085
10	5.6314
11	4.2235
12	3.1676

Table B.1: Predicted values for the GridWorld problem



Figure B.1: Episode 1



Figure B.2: Episode 25



Figure B.3: Episode 50



Figure B.4: Episode 100



Figure B.5: Episode 150



Figure B.6: Episode 200



Figure B.7: Episode 250



Figure B.8: Episode 500



Figure B.9: Episode 750



Figure B.11: Episode 1500



Figure B.12: Episode 2000



Figure B.13: Episode 2500

Appendix C

Blackjack AI

C.1 Source Code

C.1.1 Face

```
1 //Adding Face to SKCards
2 package SKCards;
3
4
5 /**
6  * Defines an enum of face values. The values are 2–10(the words),
7  * Jack, Queen, King, Ace and Joker.
8  *
9  * @author SK
10 * @version 2.0A
11 * @since 1.0A
12 * @see SKCards.Suit Suit values
13 * @SK.bug none
14 */
15 public enum Face{TWO, THREE, FOUR, FIVE, SIX, SEVEN, EIGHT, NINE, TEN,
16     JACK, QUEEN, KING, ACE, JOKER};
```

C.1.2 Suit

```
1 //Adding Suit to SKCards
2 package SKCards;
3
4
5 /**
6  * Defines an enum of the playing card suits. The values are Diamonds, Clubs,
7  * Hearts, Spades and Null for jokers.
8  *
9  * @author SK
10 * @version 2.0A
11 * @since 1.0A
12 * @see SKCards.Face Face values
13 * @SK.bug none
14 */
15 public enum Suit{DIAMONDS, CLUBS, HEARTS, SPADES, NULL}
```

C.1.3 Card

```

1  //Adding Card to SKCards
2  package SKCards;
3
4
5  /*
6   * Importing the ArrayList class
7   * and the swing/awt components
8   */
9  import java.util.ArrayList;
10 import java.awt.*;
11 import javax.swing.*;
12
13
14 /**
15  * Card represents a playing card and has two sates – suit and face value.
16  * Suit is defined in an enum and can be Diamonds, Clubs, Hearts, Spades or
17  * Null for Jokers. The face value is also defined in an enum and is the
18  * string representation of the card's face value.
19  *
20  * @author SK
21  * @version 2.0A
22  * @since 1.0A
23  * @SK.bug none
24  */
25 public class Card{
26
27     private Suit suit;
28     private Face faceValue;
29
30     /**
31      * Constructs a new card with the values passed as arguments.
32      * @param v the face value of the card.
33      * @param s the suit of the card.
34      * @see SKCards.Face Face Values
35      * @see SKCards.Suit Suit Vaules
36      */
37     public Card(Face v, Suit s){
38
39         faceValue = v;
40         suit = s;
41     }
42
43
44     /**
45      * Returns the face value of the active Card.
46      * @return the code of the Face value.
47      */
48     public String getValue(){
49
50         if(this == null){
51             return "␣";
52         }

```



```

53         else if (faceValue == Face.TWO){
54             return "2";
55         }
56         else if (faceValue == Face.THREE){
57             return "3";
58         }
59         else if (faceValue == Face.FOUR){
60             return "4";
61         }
62         else if (faceValue == Face.FIVE){
63             return "5";
64         }
65         else if (faceValue == Face.SIX){
66             return "6";
67         }
68         else if (faceValue == Face.SEVEN){
69             return "7";
70         }
71         else if (faceValue == Face.EIGHT){
72             return "8";
73         }
74         else if (faceValue == Face.NINE){
75             return "9";
76         }
77         else if (faceValue == Face.TEN){
78             return "T";
79         }
80         else if (faceValue == Face.JACK){
81             return "J";
82         }
83         else if (faceValue == Face.QUEEN){
84             return "Q";
85         }
86         else if (faceValue == Face.KING){
87             return "K";
88         }
89         else if (faceValue == Face.ACE){
90             return "A";
91         }
92         else if (faceValue == Face.JOKER){
93             return "Joker";
94         }
95         else{
96             return "";
97         }
98     }
99
100
101     /**
102     * Returns the face value of the active Card.
103     * @return the code of the Suit value.
104     */
105     public String getSuit(){
106

```

```

107         if(this == null){
108             return "";
109         }
110         else if(suit == Suit.NULL){
111             return "";
112         }
113         else{
114             return (new Character(suit.toString().charAt(0))).toString();
115         }
116     }
117
118
119     /**
120      * Gets the String value of the card int the form FS(FaceSuit), except the
121      * joker. eg the 2 of Hearts is 2H, the 10 of Diamonds is TD, the King of
122      * Spades is KS, the Joker is Joker.
123      * @return the String value of the card
124      */
125     public String toString(){
126
127         if(this == null){
128             return "␣";
129         }
130         else if(faceValue == Face.JOKER){
131             return "Joker";
132         }
133         else{
134             return (getValue() + getSuit());
135         }
136     }
137
138
139     /**
140      * Gets the image associated with the given card, using the cardSet provided.
141      * @param cs the cardSet
142      * @return the card's image
143      */
144     public ImageIcon toImage(String cs){
145
146         return new ImageIcon("images/cards/" + cs + "/" + toString() + ".jpg");
147     }
148
149
150     /**
151      * Returns the image of the card face down, using the back provided
152      * @param b the back
153      * @return the card's back image
154      */
155     public ImageIcon toBackImage(String b){
156
157         return new ImageIcon("images/backs/" + b);
158     }
159
160

```

```
161     /**
162      * Prints out the description of the active card.
163      */
164     public void print(){
165
166         System.out.println(this.toString());
167     }
168
169
170     /**
171      * Returns the point value of the card for general uses. This methods counts
172      * Ace as 1, number cards at face value, Jack as 11, Queen as 12, King as 13
173      * and Joker as 0.
174      * @return the point value of the card.
175      */
176     public int getPoints(){
177
178         if(faceValue == Face.ACE){
179             return 1;
180         }
181         else if(faceValue == Face.TWO){
182             return 2;
183         }
184         else if(faceValue == Face.THREE){
185             return 3;
186         }
187         else if(faceValue == Face.FOUR){
188             return 4;
189         }
190         else if(faceValue == Face.FIVE){
191             return 5;
192         }
193         else if(faceValue == Face.SIX){
194             return 6;
195         }
196         else if(faceValue == Face.SEVEN){
197             return 7;
198         }
199         else if(faceValue == Face.EIGHT){
200             return 8;
201         }
202         else if(faceValue == Face.NINE){
203             return 9;
204         }
205         else if(faceValue == Face.TEN){
206             return 10;
207         }
208         else if(faceValue == Face.JACK){
209             return 11;
210         }
211         else if(faceValue == Face.QUEEN){
212             return 12;
213         }
214         else if(faceValue == Face.KING){
```

```
215         return 13;
216     }
217     else{
218         return 0;
219     }
220 }
221
222 /**
223  * Gets the point value of a face value card.
224  * @param faceValue the face value of the cards.
225  */
226 public static int getPoints(Face faceValue){
227
228     if(faceValue == Face.ACE){
229         return 1;
230     }
231     else if(faceValue == Face.TWO){
232         return 2;
233     }
234     else if(faceValue == Face.THREE){
235         return 3;
236     }
237     else if(faceValue == Face.FOUR){
238         return 4;
239     }
240     else if(faceValue == Face.FIVE){
241         return 5;
242     }
243     else if(faceValue == Face.SIX){
244         return 6;
245     }
246     else if(faceValue == Face.SEVEN){
247         return 7;
248     }
249     else if(faceValue == Face.EIGHT){
250         return 8;
251     }
252     else if(faceValue == Face.NINE){
253         return 9;
254     }
255     else if(faceValue == Face.TEN){
256         return 10;
257     }
258     else if(faceValue == Face.JACK){
259         return 11;
260     }
261     else if(faceValue == Face.QUEEN){
262         return 12;
263     }
264     else if(faceValue == Face.KING){
265         return 13;
266     }
267     else{
268         return 0;
```

```

269     }
270 }
271
272 /**
273  * Returns the points value of the Card for Blackjack. This
274  * method is only used by the BlackJack prototype. The SKCards
275  * system uses the BJRules.getBJPoints(Card) method, which is
276  * almost identical.
277  *
278  * @return the Blackjack points value of the Card.
279  */
280 public int getBJPoints(){
281
282     int points = getPoints();
283
284     if(points == 10 || points == 11 || points == 12 || points == 13){
285         return 10;
286     }
287     else{
288         return points;
289     }
290 }
291 }

```

C.1.4 Pack

```

1  //Adding Pack to SKCards
2  package SKCards;
3
4
5  /*
6   * Importing the ArrayList class
7   */
8  import java.util.ArrayList;
9
10
11 /**
12  * Represents a pack of cards. A pack is composed of a number of decks, with
13  * or without jokers, both to be specified by the user.
14  *
15  * @author SK
16  * @version 2.0A
17  * @since 1.0A
18  * @see SKCards.Card Card
19  * @SK.bug none
20  */
21 public class Pack{
22
23     private final static Suit[] suits = Suit.values();
24     private final static Face[] vals = Face.values();
25
26     private ArrayList<Card> pack;
27     protected int decks, currentCard;
28     private boolean jokers;

```

```

29
30
31  /**
32   * The default constructor. Creates a pack of Cards, with
33   * 1 deck and no jokers
34   */
35  public Pack(){
36      this(1, false);
37  }
38
39  /**
40   * Creates a new pack of cards. This pack is defined by
41   * the parameters passed to it.
42   *
43   * @param decks the number of decks in the pack.
44   * @param jokers whether the decks have jokers or not.
45   * @see SKCards.Card#makeDeck
46   */
47  public Pack(int d, boolean j){
48
49      decks = d; jokers = j;
50
51      pack = new ArrayList<Card>();
52
53      ArrayList<Card> deck = new ArrayList<Card>();
54
55      for(int i = 1; i <= decks; i++){
56
57          deck = makeDeck(jokers);
58          deck = shuffle(deck);
59
60          for(Card c: deck){
61
62              pack.add(c);
63          }
64      }
65
66      pack = Pack.shuffle(pack);
67      currentCard = 0;
68  }
69
70
71  /**
72   * Creates a new ordered deck of cards, with or without jokers depending
73   * on the arguments passed.
74   * @param jokers a boolean which determines if the deck has jokes or not.
75   * @return deck an ArrayList of the cards.
76   */
77  public static ArrayList<Card> makeDeck(boolean jokers){
78
79      ArrayList<Card> deck = new ArrayList<Card>();
80
81      for(int i = 0; i < suits.length - 1; i++){
82          for(int j = 0; j < vals.length - 1; j++){

```

```

83
84         deck.add(new Card( vals[j], suits[i]));
85     }
86 }
87
88
89     if(jokers){
90         deck.add(new Card(Face.JOKER, Suit.NULL));
91         deck.add(new Card(Face.JOKER, Suit.NULL));
92     }
93
94     return deck;
95 }
96
97
98 /**
99  * Prints out all the cards in the deck.
100 */
101 public void print(){
102
103     for(Card c: pack){
104
105         c.print();
106     }
107 }
108
109
110 /**
111  * Takes a pack of cards and shuffles them.
112  * @param pack the pack of cards to be shuffle.
113  * @return the shuffledPack.
114  */
115 public static ArrayList<Card> shuffle(ArrayList<Card> pack){
116
117     ArrayList<Boolean> shuffled = new ArrayList<Boolean>();
118     ArrayList<Card> shuffledPack = new ArrayList<Card>();
119
120     for(int i = 0; i < pack.size(); i++){
121
122         shuffled.add(false);
123     }
124
125     while(shuffled.contains(false)){
126
127         Double numb = Math.random();
128         numb *= 125;
129         int numbr = numb.intValue();
130         int num = numbr % pack.size();
131
132         if(!shuffled.get(num)){
133
134             shuffledPack.add(pack.get(num));
135             shuffled.set(num, true);
136         }

```

```

137     }
138
139     return shuffledPack;
140 }
141
142 /**
143  * Regenerates the Pack, by creating a new Pack. This is identical to the
144  * previous only in the number of decks and the presence of jokers, but
145  * the order of the Cards is different.
146  */
147 protected void newPack(){
148
149     pack = new ArrayList<Card>();
150
151     ArrayList<Card> deck = new ArrayList<Card>();
152
153     for(int i = 1; i <= decks; i++){
154
155         deck = makeDeck(jokers);
156         deck = shuffle(deck);
157
158         for(Card c: deck){
159
160             pack.add(c);
161         }
162     }
163
164     pack = Pack.shuffle(pack);
165     currentCard = 0;
166 }
167
168 /**
169  * Gives the next Card in the Pack. If the Pack has been used
170  * up, a new Pack is generated.
171  *
172  * @return the next Card
173  */
174 public Card nextCard(){
175
176     if(currentCard >= size()){
177         newPack();
178     }
179
180     currentCard++;
181     return pack.get(currentCard - 1);
182 }
183
184
185 /**
186  * Returns the size of the Pack
187  *
188  * @return size of the Pack
189  */
190 public int size(){

```



```

191
192         return pack.size();
193     }
194
195     /**
196      * Returns the number of decks in the Pack
197      *
198      * @return the number of decks in the Pack
199      */
200     public int decks(){
201
202         return decks;
203     }
204
205
206     /**
207      * Returns a boolean to indicate if there are
208      * Jokers in the Pack or not
209      *
210      * @return true if there are Jokers
211      * @return false if there are no Jokers
212      */
213     public boolean jokers(){
214
215         return jokers;
216     }
217 }

```

C.1.5 Blackjack

```

1  //Adding BlackJack to SKCards
2  package SKCards;
3
4  /*
5   * Importing the swing/awt components and the event models
6   * and the I/O classes
7   */
8  import java.awt.*;
9  import java.awt.event.*;
10 import javax.swing.*;
11 import javax.swing.event.*;
12 import java.io.*;
13
14
15 /**
16  * Blackjack game with Temporal Difference AI.
17  */
18 public class BlackJack extends JFrame{
19
20     /**
21      * This class represents a player in the BlackJack game. It can be either a human
22      * player or a Reinforcement Learning player, using Temporal Difference.
23      * @author Saqib A Kakvi (33106087)
24      * @version 1.0

```

```

25      * @rubish yes it is
26      */
27      private class BPlayer{
28
29          //ACTIONS
30          private static final int HIT = 1;
31          private static final int STAND = 0;
32          private static final int DEAL = -1;
33
34          //CONSANTS
35          private final double TAU = 0.5;
36          private final double G = 0.75;
37          private final double A = 0.6;
38
39          //the player's hand
40          private Card[] hand;
41          //the player's funds
42          private int funds;
43          //the players scores counting low/high when applicable
44          private int score, hscore;
45          //the values attached of the state-action pairs
46          private double[][] val;
47          //*****TESTING***** the number of times the agent takes a certain action
48          private int[][] count;
49          //the players last action and score
50          private int lastaxn, lastscore;
51          //flags to indicate if the player is playing an ace as high, is a greedy agent
52          //or is playing a fixed heuristic policy
53          boolean high, greedy, fixed;
54          //*****TESTING***** writing all the action data to file
55          BufferedWriter p;
56
57
58          /**
59           * Constructor.
60           */
61          public BPlayer(){
62              hand = new Card[5];
63              funds = 1000000;
64              score = 0;
65              hscore = 10;
66              val = new double[2][21];
67              count = new int[2][21];
68              lastaxn = DEAL;
69              lastscore = 0;
70              greedy = false;
71              fixed = false;
72              try{
73                  p = new BufferedWriter(new FileWriter("playerdump.csv"));
74                  p.write("SCORE_");
75                  p.write("HITPROB_");
76                  p.write("STANDPROB_");
77                  p.write("RNDML_");
78                  p.write("AXN_");

```

```

79         p.newLine();
80     }
81     catch (Exception ex){
82         ex.printStackTrace();
83     }
84 }
85
86
87 /**
88  * This method uses the softmax method of exploration to pick an action in the
89  * current state. After the action is picked, the previous state is updated based
90  * on the Temporal-Difference Q-Learning algorithm.
91  */
92 private void explore(){
93
94     greedy = false;
95     fixed = false;
96
97     if (high && hscore <= 21){
98         score = hscore;
99     }
100    else if (high && hscore > 21 && score == hscore){
101        high = false;
102        score -= 10;
103    }
104    double hit = Math.exp(val[HIT][score-1]/TAU);
105    double stand = Math.exp(val[STAND][score-1]/TAU);
106
107    double sum = hit + stand;
108
109    double hprob = hit / sum;
110    double sprob = stand / sum;
111
112    double maxval = Math.max(val[STAND][score-1], val[HIT][score-1]);
113
114    if (lastaxn != DEAL){
115        double rd = A * ((G * maxval) - val[lastaxn][lastscore-1]);
116        val[lastaxn][lastscore-1] += rd;
117    }
118
119    double r = Math.random();
120
121    lastscore = score;
122
123    try{
124        p.write(Integer.toString(score)); p.write(" ");
125        p.write(Double.toString(hprob)); p.write(" ");
126        p.write(Double.toString(sprob)); p.write(" ");
127        p.write(Double.toString(r)); p.write(" ");
128        if (r > hprob){
129            stand();
130            lastaxn = STAND;
131            p.write("STAND");
132        }

```

```

133         else{
134             hit();
135             lastaxn = HIT;
136             p.write("HIT");
137         }
138         p.newLine();
139     }
140     catch(Exception ex){
141         ex.printStackTrace();
142     }
143     count[lastaxn][lastscore-1]++;
144 }
145
146
147 /**
148  * This method provides a greedy agent. This agent is defensive in its behavior
149  * as it stands when the value of hitting is equal to the value of standing.
150  */
151 private void greedy(){
152
153     greedy = true;
154     fixed = false;
155
156     if(high && hscore <= 21){
157         score = hscore;
158     }
159     else if(high && hscore > 21 && score == hscore){
160         high = false;
161         score -= 10;
162     }
163
164     lastscore = score;
165
166     if(val[HIT][score-1] > val[STAND][score-1]){
167         hit();
168         lastaxn = HIT;
169     }
170     else{
171         stand();
172         lastaxn = STAND;
173     }
174     count[lastaxn][lastscore-1]++;
175 }
176
177
178 /**
179  * This method provides a fixed heuristic agent for the Reinforcement Learning
180  * agent to learn against.
181  */
182 private void fixed(){
183
184     greedy = false;
185     fixed = true;
186

```

```

187         if (high && hscore <= 21){
188             score = hscore;
189         }
190         else if (high && hscore > 21 && score == hscore){
191             high = false;
192             score -= 10;
193         }
194
195         if (score < 17){
196             hit();
197         }
198         else{
199             stand();
200         }
201     }
202
203
204     /**
205      * This method is used to allocate final rewards to the last state visited.
206      * agents do not receive rewards, as they are not learning. This reward func
207      * based on the Temporal Difference Q-Learning algorithm
208      */
209     private void reward(double r){
210
211         if (greedy || fixed){
212             return;
213         }
214
215         double maxval = Math.max(val[STAND][lastscore - 1], val[HIT][lastscore - 1]);
216
217         if (lastaxn != DEAL){
218             double rd = A * (r + (G * maxval) - val[lastaxn][lastscore - 1]);
219             val[lastaxn][lastscore - 1] += rd;
220         }
221     }
222 }
223
224
225 //the playing agents
226 private BPlayer player, dealer;
227 //the pack of card
228 private Pack pack;
229 //the players' hands
230 private JLabel[] pcards, dcards;
231 //display of the users funds
232 private JLabel fundsLabel;
233 //command buttons
234 private JButton hit, stand, betButton, clear;
235 //panels used to display user components
236 private JPanel buttonPanel, playerPanel, dealerPanel;
237 //the content pane
238 private Container cp;
239 //textfield for user to enter amount to bet
240 private JTextField betField;

```

```

241 //amount the user has bet
242 private int bet;
243 //used to indicate current position in the players hand
244 private int currentCard;
245 //flag to indicate if the player is playing or not
246 private boolean play;
247 ****TEST**** stores number of times each result is recorded.
248 private int[] result;
249 ****TEST**** stores number of times each result is recorded.
250 private int[] deals;
251 ****TEST**** records the scores of each player when the draw.
252 private BufferedWriter push;
253 ****TEST**** stores number of times the player busts at each score.
254 private int[] busts;
255
256 /**
257  * The main method. Creates a new BlackJack interface.
258  * @param args not used
259  */
260 public static void main(String[] args){new BlackJack();}
261
262
263 /**
264  * Sole constructor. Creates the interface and the players. Begins the learning
265  * for the agents(s).
266  */
267 public BlackJack(){
268
269     super("BlackJack");
270     setSize(600,600);
271     Dimension d = (Toolkit.getDefaultToolkit()).getScreenSize();
272     int h = (int)(d.getHeight());
273     int w = (int)(d.getWidth());
274     int x = (w - 600) / 2;
275     int y = (h - 600) / 2;
276     setLocation(x,y);
277     setDefaultCloseOperation(EXIT_ON_CLOSE);
278
279     pack = new Pack();
280     player = new BPlayer();
281     dealer = new BPlayer();
282     play = true;
283     result = new int[6];
284     deals = new int[21];
285     busts = new int[21];
286     currentCard = 0;
287     player.score = 0;
288
289     try{
290         push = new BufferedWriter(new FileWriter("push.txt"));
291     }
292     catch(Exception ex){
293         ex.printStackTrace();
294     }

```

```

295
296     cp = getContentPane();
297     cp.setLayout(new BorderLayout());
298
299     hit = new JButton("HIT");
300     hit.addMouseListener(new MouseAdapter(){
301
302         public void mouseClicked(MouseEvent e){
303             hit();
304         }
305     });
306     hit.setEnabled(false);
307
308     stand = new JButton("STAND");
309     stand.addMouseListener(new MouseAdapter(){
310
311         public void mouseClicked(MouseEvent e){
312             stand();
313         }
314     });
315     stand.setEnabled(false);
316
317     betField = new JTextField("10");
318     betButton = new JButton("BET");
319     betButton.addMouseListener(new MouseAdapter(){
320
321         public void mouseClicked(MouseEvent e){
322             bet();
323         }
324     });
325
326     clear = new JButton("CLEAR");
327     clear.addMouseListener(new MouseAdapter(){
328
329         public void mouseClicked(MouseEvent e){
330             clear();
331         }
332     });
333
334     buttonPanel = new JPanel();
335     buttonPanel.setLayout(new GridLayout(5,1));
336
337     buttonPanel.add(hit);
338     buttonPanel.add(stand);
339     buttonPanel.add(betField);
340     buttonPanel.add(betButton);
341     buttonPanel.add(clear);
342
343     pcards = new JLabel[5];
344     dcards = new JLabel[5];
345     fundsLabel = new JLabel(Integer.toString(player.funds));
346
347     playerPanel = new JPanel();
348     playerPanel.setLayout(new GridLayout(1,6));

```

```

349     dealerPanel = new JPanel();
350     dealerPanel.setLayout(new GridLayout(1,6));
351
352     for(int i = 0; i < 5; i++){
353         pcards[i] = new JLabel();
354         dcards[i] = new JLabel();
355         playerPanel.add(pcards[i]);
356         dealerPanel.add(dcards[i]);
357     }
358     playerPanel.add(fundsLabel);
359     dealerPanel.add(new JLabel());
360
361     cp.add(dealerPanel, BorderLayout.NORTH);
362     cp.add(buttonPanel, BorderLayout.EAST);
363     cp.add(playerPanel, BorderLayout.SOUTH);
364
365     setVisible(true);
366
367     int c = 0;
368
369     for(int i = 0; i < 1000; i++){
370         bet();
371         while(play){
372             player.explore();
373         }
374         while(!play){
375             dealer.fixed();
376         }
377         System.out.println("played_" + c);
378         c++;
379     }
380
381     write(1);
382
383     for(int i = 0; i < 9000; i++){
384         bet();
385         while(play){
386             player.explore();
387         }
388         while(!play){
389             dealer.fixed();
390         }
391         System.out.println("played_" + c);
392         c++;
393     }
394
395     write(2);
396
397     for(int i = 0; i < 40000; i++){
398         bet();
399         while(play){
400             player.explore();
401         }
402         while(!play){

```



```

403         dealer.fixed();
404     }
405     System.out.println("played_" + c);
406     c++;
407 }
408
409 write(3);
410
411 result = new int[6];
412 busts = new int[21];
413 player.count = new int[2][21];
414
415 for(int i = 0; i < 5000; i++){
416     bet();
417     while(play){
418         player.greedy();
419     }
420     while(!play){
421         dealer.fixed();
422     }
423     System.out.println("played_" + c);
424     c++;
425 }
426
427 write(4);
428 writecount();
429 try{
430     push.flush();
431     push.close();
432     player.p.flush();
433     player.p.close();
434 }
435 catch(Exception ex){
436     ex.printStackTrace();
437 }
438 }
439
440
441 /**
442  * *****TESTING METHOD***** <br/>
443  * Writes all the data collected to .csv files , which can be imported into an
444  * spreadsheet package for further data analysis.
445  * @param x the number tag for the file.
446  */
447 private void write(int x){
448
449     try{
450         BufferedWriter bw = new BufferedWriter(new FileWriter("results" + x + ".c
451
452         bw.write("Player_Bust");
453         bw.write("\t");
454         bw.write("Dealer_Bust");
455         bw.write("\t");
456         bw.write("Both_Bust");

```

```

457         bw.write("\t");
458         bw.write("Player_wins");
459         bw.write("\t");
460         bw.write("Player_Loses");
461         bw.write("\t");
462         bw.write("Push");
463         bw.newLine();
464
465         for(int i = 0; i < result.length; i++){
466             bw.write(Integer.toString(result[i]));
467             bw.write("\t");
468         }
469
470         bw.flush();
471         bw.close();
472
473         bw = new BufferedWriter(new FileWriter("pvals" + x + ".csv"));
474
475         bw.write("score");
476         bw.write("\t");
477         bw.write("stand");
478         bw.write("\t");
479         bw.write("hit");
480         bw.write("\t");
481         bw.newLine();
482
483         for(int i = 0; i < player.val[0].length; i++){
484             bw.write(Integer.toString(i));
485             bw.write("\t");
486             for(int j = 0 ; j < player.val.length; j++){
487                 bw.write(Double.toString(player.val[j][i]));
488                 bw.write("\t");
489             }
490             bw.newLine();
491         }
492
493         bw.flush();
494         bw.close();
495
496         bw = new BufferedWriter(new FileWriter("dvals" + x + ".csv"));
497
498         bw.write("score");
499         bw.write("\t");
500         bw.write("stand");
501         bw.write("\t");
502         bw.write("hit");
503         bw.write("\t");
504         bw.newLine();
505
506         for(int i = 0; i < player.count[0].length; i++){
507             bw.write(Integer.toString(i));
508             bw.write("\t");
509             for(int j = 0 ; j < player.count.length; j++){
510                 bw.write(Double.toString(player.count[j][i]));

```

```

511         bw.write("\t");
512     }
513     bw.newLine();
514 }
515
516 bw.flush();
517 bw.close();
518
519 bw = new BufferedWriter(new FileWriter("busts" + x + ".csv"));
520 for(int i = 0; i < busts.length; i++){
521     bw.write(Integer.toString(i));
522     bw.write("\t");
523     bw.write(Integer.toString(busts[i]));
524     bw.newLine();
525 }
526
527 bw.flush();
528 bw.close();
529 }
530 catch(Exception ex){
531     ex.printStackTrace();
532 }
533 }
534
535 public void writecount(){
536
537     try{
538         BufferedWriter bw = new BufferedWriter(new FileWriter("count.txt"));
539
540         for(int i = 0; i < 2; i++){
541             for(int j = 0; j < 21; j++){
542                 bw.write(new Integer(player.count[j][i]).toString());
543                 bw.write("_");
544             }
545             bw.newLine();
546         }
547         bw.flush();
548         bw.close();
549     }
550     catch(Exception ex){
551
552     }
553 }
554
555 /**
556  * This method is used to update the GUI to the most current state. Every call of
557  * methods redraws all elements on screen and then repaints the JFrame.
558  * @see clear
559  */
560 private void refresh(){
561
562     if(play){
563         for(int i = 0; i < 5; i++){
564             if(player.hand[i] != null){

```

```

565         pcards[i].setIcon(player.hand[i].toImage("Normal"));
566     }
567 }
568 dcards[0].setIcon(dealer.hand[0].toImage("Normal"));
569 dcards[1].setIcon(dealer.hand[1].toBackImage("Trickster"));
570
571 }
572 else{
573     for(int i = 0; i < 5; i++){
574         if(dealer.hand[i] != null){
575             dcards[i].setIcon(dealer.hand[i].toImage("Normal"));
576         }
577     }
578 }
579
580 fundsLabel.setText(Integer.toString(player.funds));
581 ((JComponent)(cp)).revalidate();
582 repaint();
583 }
584
585
586 /**
587  * Bets money for the player. The amount to be bet is parsed from betField. betF
588  * is set to 10 by default. This is so the AI does not have to be made to bet, as
589  * can simply call the method.
590  */
591 private void bet(){
592
593     bet = Integer.parseInt(betField.getText());
594
595     if(bet <= player.funds){
596         player.funds -= bet;
597         deal();
598         hit.setEnabled(true);
599         stand.setEnabled(true);
600         betButton.setEnabled(false);
601     }
602     refresh();
603 }
604
605
606 /**
607  * Deals a new hand. The player/dealer hands are reinitialised and the pointer an
608  * all reset. Both scores are then updated.
609  * @see hit
610  * @see stand
611  */
612 private void deal(){
613
614     dealer.hand[0] = pack.nextCard();
615     dealer.score += dealer.hand[0].getBJPoints();
616     if(dealer.hand[0].getValue().equals("A")){
617         dealer.high = true;
618         dealer.hscore = dealer.score + 10;

```

```

619     }
620     dealer.hand[1] = pack.nextCard();
621     dealer.score += dealer.hand[0].getBJPoints();
622     if(dealer.hand[1].getValue().equals("A")){
623         dealer.high = true;
624         dealer.hscore = dealer.score + 10;
625     }
626     dealer.lastaxn = BPlayer.DEAL;
627     if(dealer.high && dealer.hscore <= 21){
628         dealer.score = dealer.hscore;
629     }
630
631     player.hand[0] = pack.nextCard();
632     player.score += player.hand[0].getBJPoints();
633     if(player.hand[0].getValue().equals("A")){
634         player.high = true;
635         player.hscore = player.score + 10;
636     }
637     player.hand[1] = pack.nextCard();
638     player.score += player.hand[1].getBJPoints();
639     if(player.hand[1].getValue().equals("A")){
640         player.high = true;
641         player.hscore = player.score + 10;
642     }
643     if(player.high && player.hscore <= 21){
644         player.score = player.hscore;
645     }
646     deals[player.score - 1]++;
647     player.lastaxn = BPlayer.DEAL;
648     currentCard = 2;
649 }
650
651
652 /**
653  * This method adds a Card to the active BPlayer's hand. The pointer is moved to
654  * next Card. The player is forced to stand if they exceed a score of 21 or have
655  * cards.
656  * @see deal
657  * @see stand
658  */
659 private void hit(){
660
661     if(play){
662         player.hand[currentCard] = pack.nextCard();
663         player.score += player.hand[currentCard].getBJPoints();
664         if(player.hand[currentCard].getValue().equals("A")){
665             player.high = true;
666             player.hscore += player.hand[currentCard].getBJPoints();
667         }
668         refresh();
669         currentCard++;
670         if(player.score > 21 || currentCard == 5){
671             stand();
672         }

```

```

673     }
674     else{
675         dealer.hand[currentCard] = pack.nextCard();
676         dealer.score += dealer.hand[currentCard].getBJPoints();
677         if(dealer.hand[currentCard].getValue().equals("A")){
678             dealer.high = true;
679             dealer.hscore += dealer.hand[currentCard].getBJPoints();
680         }
681         refresh();
682         currentCard++;
683         if(dealer.score > 21 || currentCard == 5){
684             stand();
685         }
686     }
687 }
688
689
690 /**
691  * Stands for the agent that called the method. If it is the player, then the pla
692  * is passed on to the dealer. If it is the dealer, then the result is calculated
693  * and the hand ends
694  * @see deal
695  * @see hit
696  */
697 private void stand(){
698     if(play){
699         play = false;
700         currentCard = 2;
701         refresh();
702     }
703     else{
704         if(player.score > 21 && dealer.score <= 21){
705             //player bust
706             player.reward(-20);
707             result[0]++;
708             busts[player.lastscore-1]++;
709         }
710         else if(player.score <= 21 && dealer.score > 21){
711             //dealer bust
712             player.funds += bet * 2;
713             player.reward(10);
714             result[1]++;
715         }
716         else if(player.score > 21 && dealer.score > 21){
717             //dealer and player bust
718             player.reward(-10);
719             result[2]++;
720             busts[player.lastscore-1]++;
721         }
722         else if(player.score > dealer.score){
723             //player wins
724             player.funds += bet * 2;;
725             player.reward(10);
726             result[3]++;

```

```

727     }
728     else if(player.score < dealer.score){
729         //player loses
730         player.reward(-10);
731         result[4]++;
732     }
733     else{
734         //push
735         player.funds += bet;
736         player.reward(5);
737         result[5]++;
738         try{
739             push.write(Integer.toString(player.score));
740             push.write(" ");
741             push.write(Integer.toString(dealer.score));
742             push.newLine();
743         }
744         catch(Exception ex){
745             ex.printStackTrace();
746         }
747     }
748     bet = 0;
749     clear();
750     play = true;
751 }
752 }
753
754
755 /**
756  * Clears all GUI elements and resets all values to thier intitial values.
757  * @see refresh
758  */
759 private void clear(){
760
761     playerPanel.removeAll();
762     dealerPanel.removeAll();
763
764     for(int i = 0; i < 5; i++){
765         player.hand[i] = null;
766         dealer.hand[i] = null;
767         pcards[i] = new JLabel();
768         dcards[i] = new JLabel();
769         playerPanel.add(pcards[i]);
770         dealerPanel.add(dcards[i]);
771     }
772
773     playerPanel.add(fundsLabel);
774     bet = 0;
775     currentCard = 0;
776     player.score = 0;
777     dealer.score = 0;
778     dealer.high = false;
779     player.high = false;
780     fundsLabel.setText((new Integer(player.funds)).toString());

```

```
781         betButton.setEnabled(true);
782         hit.setEnabled(false);
783         stand.setEnabled(false);
784         repaint();
785     }
786 }
```


Appendix D

Final System

D.1 Source Code

D.1.1 Face

See Appendix C.1.1

D.1.2 Suit

See Appendix C.1.2

D.1.3 Card

See Appendix C.1.3

D.1.4 Pack

See Appendix C.1.2

D.1.5 Shoe

```
1 //Adding Shoe to the SKCards package
2 package SKCards;
3
4
5 /**
6  * This object represents a Dealer's shoe. This is a number of decks
7  * of cards in a closed container. There is a special cut card, which
8  * is inserted randomly into the Shoe. When this is dealt, the shoe is
9  * discarded and a new one is used. Shoes are used to counter Card counting.
10 *
11 * @author Saqib A Kakvi
12 * @version 2.0A
13 * @since 2.0A
14 * @see SKCards.Pack Pack
15 * @SK.bug none
16 */
17 public class Shoe extends Pack{
18
19     //the position of the cutoff
```

```

20     int cutoff;
21     //the size of the Shoe
22     int size;
23
24     /**
25      * Constructor. Creates a pack with the specified parameters by calling
26      * super(d,j). The cutoff Card's position is calculated and is stored.
27      * No extra Card is added to the Pack, so the size remains the same.
28      *
29      * @param d the number of decks
30      * @param j if there are jokers in the deck
31      */
32     public Shoe(int d, boolean j){
33
34         super(d,j);
35         size = super.size();
36         calculateCutoff();
37     }
38
39
40     /**
41      * Calculates the cutoff Card position. The Card will be in the range of
42      * 0-85% of the Shoe size. Hence a maximum of of 85% of the Pack is
43      * played.
44      */
45     private void calculateCutoff(){
46         cutoff = new Double(Math.random() * 100 * 0.85).intValue() * size;
47     }
48
49
50     /**
51      * Checks if the current Card is the cutoff Card. If it is then a new Pack
52      * is created and the cutoff Card's position is recalculated. Then the
53      * method returns the Card from super.nextCard().
54      *
55      * @override SKCards.Pack nextCard
56      * @return the next Card in the Shoe
57      */
58     public Card nextCard(){
59
60         if(currentCard == cutoff){
61             newPack();
62             calculateCutoff();
63         }
64         return super.nextCard();
65     }
66 }

```

D.1.6 BJPlayer

```

1  //Adding Player to SKCards
2  package SKCards;
3
4

```

```
5  /*
6   * Importing the ArrayList class
7   * and the Arrays class
8   * and the swing/awt components
9   */
10 import java.util.ArrayList;
11 import java.util.Arrays;
12 import java.awt.*;
13 import javax.swing.*;
14
15
16 /**
17  * This is a player in the BlackJack game
18  *
19  * @author SK
20  * @version 2.0A
21  * @since 1.0C
22  */
23 public class BJPlayer extends Player{
24
25     protected int score , hscore;
26
27     protected Card[] hand;
28
29     protected BJRules rules;
30
31     boolean high;
32
33     /**
34      * Creates a new player object with a given details.
35      *
36      * @param id the player name
37      * @param cash the amount of money the player starts with
38      * @param setPack the pack of cards to be used in the game
39      * @param bck the card back the player uses
40      * @param cs the card set used by the player
41      */
42     public BJPlayer(String id , int cash , String bck , String cs ,
43         BJRules r){
44
45         rules = r;
46         funds = cash;
47         name = id;
48         back = bck;
49         cardSet = cs;
50         high = false;
51         score = 0;
52         hscore = 10;
53
54         hand = new Card[5];
55     }
56
57     /**
58      * Basic constructor. Simply takes in the rules. Used for
```

```

59      * BJRLPlayer
60      *
61      * @param r the rules
62      */
63      public BJPlayer(BJRules r){
64
65          rules = r;
66      }
67
68      /**
69       * Creates a new hand of cards for the current player and one for the computer
70       * and then deals two cards to each hand
71       */
72      public void newHand(){
73          hand = new Card[5];
74          score = 0;
75      }
76
77      /**
78       * Hits the player's hand
79       *
80       * @param currentCard the current card
81       * @param nextCard the card to be added
82       */
83      public int hit(int currentCard, Card nextCard){
84
85          hand[currentCard] = nextCard;
86
87          score += rules.getBJPoints(hand[currentCard]);
88
89          if(hand[currentCard].getValue().equals("A")){
90              high = true;
91              hscore += rules.getBJPoints(hand[currentCard]);
92          }
93
94          if(high && hscore <= 21){
95              score = hscore;
96          }
97          else if(high && hscore > 21){
98              if(hscore == score){
99                  score -= 10;
100             }
101             high = false;
102         }
103
104         return score;
105     }
106
107
108     /**
109      * Returns the players current score
110      *
111      * @return score
112      */

```

```

113     public int getScore(){
114         return score;
115     }
116
117     /**
118      * Ends the current hand
119      */
120     public void endHand(){
121
122         int comp = 0;
123         int player = 0;
124
125
126         System.out.println("player——>" + player);
127         System.out.println("computer——>" + comp);
128
129         if(comp > 0 || player > 0){
130             if(comp > player){
131                 JOptionPane.showMessageDialog(null , "HOUSE_WINS!!!");
132                 rules.display();
133             }
134             else if(player > comp){
135                 JOptionPane.showMessageDialog(null , "PLAYER_WINS!!!");
136                 funds += bet;
137                 rules.display();
138             }
139             else if(comp == player){
140                 JOptionPane.showMessageDialog(null , "PUSH!!!");
141                 funds += bet;
142                 rules.display();
143             }
144             else{
145                 JOptionPane.showMessageDialog(null , "NO_RESULT_FOUND");
146             }
147         }
148         else{
149             JOptionPane.showMessageDialog(null , "NO_RESULT_FOUND");
150             rules.display();
151         }
152     }
153
154
155     /**
156      * Clears all the hands and clears the interface
157      */
158     public void clear(){
159
160         for(int i = 0; i < hand.length; i++){
161             hand[i] = null;
162         }
163
164         rules.clear();
165     }
166

```

```

167
168     /**
169      * Gets the player's current hands
170      * @return the player's hands
171      */
172     public Card[] getHand(){
173
174         return hand;
175     }
176 }

```

D.1.7 BJUI

```

1  //adding BJUI to SKCards
2  package SKCards;
3
4  /*
5   *
6   */
7  import java.awt.*;
8  import java.awt.event.*;
9  import javax.swing.*;
10 import javax.swing.event.*;
11
12 /**
13  * This is the User Interface for the Blackjack Game.
14  * It allows the human player to access the functions
15  * ans shows the current state of the game.
16  *
17  * @author SK
18  * @version 2.0A
19  * @since 2.0A
20  * @SK.bug none
21  */
22 public class BJUI extends JFrame implements UI{
23
24     //the players' hands
25     private JLabel[] pcards, dcards;
26     //display of the users funds
27     private JLabel fundsLabel;
28     //command buttons
29     private JButton hit, stand, bet, clear, exit;
30     //panels used to display user components
31     private JPanel buttonPanel, playerPanel, dealerPanel;
32     //the content pane
33     private Container cp;
34     //textfield for user to enter amount to bet
35     private JTextField betField;
36     //the game rules
37     BJRules rules;
38
39     /**
40      * Creates a new UI. Takes in the rules for the current game.
41      * This allows the two objects to communicate.

```

```

42      *
43      * @param r the Rules for the game
44      */
45      public BJUI(BJRules r){
46
47          super("BlackJack");
48          setSize(600,600);
49          Dimension d = (Toolkit.getDefaultToolkit()).getScreenSize();
50          int h = (int)(d.getHeight());
51          int w = (int)(d.getWidth());
52          int x = (w - 600) / 2;
53          int y = (h - 600) / 2;
54          setDefaultCloseOperation(DO_NOTHING_ON_CLOSE);
55          setLocation(x,y);
56
57          rules = r;
58
59          cp = getContentPane();
60          cp.setLayout(new BorderLayout());
61
62          hit = new JButton("HIT");
63          hit.addMouseListener(new MouseAdapter(){
64
65              public void mouseClicked(MouseEvent e){
66                  rules.hit();
67              }
68          });
69          hit.setEnabled(false);
70
71          stand = new JButton("STAND");
72          stand.addMouseListener(new MouseAdapter(){
73
74              public void mouseClicked(MouseEvent e){
75                  rules.stand();
76              }
77          });
78          stand.setEnabled(false);
79
80          betField = new JTextField("10");
81          bet = new JButton("BET");
82          bet.addMouseListener(new MouseAdapter(){
83
84              public void mouseClicked(MouseEvent e){
85                  rules.bet(Integer.parseInt(betField.getText()));
86              }
87          });
88          bet.setEnabled(true);
89
90          clear = new JButton("CLEAR");
91          clear.addMouseListener(new MouseAdapter(){
92
93              public void mouseClicked(MouseEvent e){
94                  clear();
95              }

```

```

96         });
97
98         exit = new JButton("EXIT");
99         exit.addMouseListener(new MouseAdapter(){
100
101             public void mouseClicked(MouseEvent e){
102                 BJUI.this.dispose();
103             }
104         });
105
106         buttonPanel = new JPanel();
107         buttonPanel.setLayout(new GridLayout(6,1));
108
109         buttonPanel.add(hit);
110         buttonPanel.add(stand);
111         buttonPanel.add(betField);
112         buttonPanel.add(bet);
113         buttonPanel.add(clear);
114         buttonPanel.add(exit);
115
116         pcards = new JLabel[5];
117         dcards = new JLabel[5];
118         fundsLabel = new JLabel();
119
120         playerPanel = new JPanel();
121         playerPanel.setLayout(new GridLayout(1,6));
122         dealerPanel = new JPanel();
123         dealerPanel.setLayout(new GridLayout(1,6));
124
125         for(int i = 0; i < 5; i++){
126             pcards[i] = new JLabel();
127             dcards[i] = new JLabel();
128             playerPanel.add(pcards[i]);
129             dealerPanel.add(dcards[i]);
130         }
131         playerPanel.add(fundsLabel);
132         dealerPanel.add(new JLabel());
133
134         cp.add(dealerPanel, BorderLayout.NORTH);
135         cp.add(buttonPanel, BorderLayout.EAST);
136         cp.add(playerPanel, BorderLayout.SOUTH);
137
138         setVisible(true);
139     }
140
141
142     /**
143      * Refreshes the GUI to show the current situation.
144      * @override display
145      */
146     public void display(){
147
148         boolean play = rules.getPlay();
149

```



```

150         Card[] phand = rules.getPlayerHand(rules.PLAYER);
151
152         Card[] chand = rules.getPlayerHand(rules.DEALER);
153
154         String cs = rules.getPlayerCardSet(rules.PLAYER);
155         String bk = rules.getPlayerBack(rules.PLAYER);
156
157         for(int i = 0; i < phand.length; i++){
158             if(phand[i] != null){
159                 pcards[i].setIcon(phand[i].toImage(cs));
160             }
161         }
162         dcards[0].setIcon(chand[0].toImage(cs));
163
164         if(play){
165             dcards[1].setIcon(chand[1].toBackImage(bk));
166         }
167         else{
168             for(int i = 1; i < chand.length; i++){
169                 if(chand[i] != null){
170                     dcards[i].setIcon(chand[i].toImage(cs));
171                 }
172             }
173         }
174
175         fundsLabel.setText(rules.getPlayerFunds(rules.PLAYER).toString());
176         repaint();
177     }
178
179
180     /**
181     * Displays a new hand on the GUI.
182     * @override newDisplay
183     */
184     public void newDisplay(){
185
186         Card[] phand = rules.getPlayerHand(rules.PLAYER);
187
188         Card[] chand = rules.getPlayerHand(rules.DEALER);
189
190         String cs = rules.getPlayerCardSet(rules.PLAYER);
191         String bk = rules.getPlayerBack(rules.PLAYER);
192
193         for(int i = 0; i < phand.length; i++){
194             if(phand[i] != null){
195                 pcards[i].setIcon(phand[i].toImage(cs));
196             }
197         }
198
199         dcards[0].setIcon(chand[0].toImage(cs));
200         dcards[1].setIcon(chand[1].toBackImage(bk));
201
202         fundsLabel.setText(rules.getPlayerFunds(rules.PLAYER).toString());
203         repaint();

```

```

204     }
205
206
207     /**
208      * Clears the GUI
209      * @override clear
210      */
211     public void clear(){
212
213         System.out.println();
214
215         playerPanel.removeAll();
216         dealerPanel.removeAll();
217
218         for(int i = 0; i < 5; i++){
219             pcards[i] = new JLabel();
220             dcards[i] = new JLabel();
221             playerPanel.add(pcards[i]);
222             dealerPanel.add(dcards[i]);
223         }
224
225         playerPanel.add(fundsLabel);
226         betEnabled(true);
227         hitEnabled(false);
228         standEnabled(false);
229         repaint();
230     }
231
232     /**
233      * Gets the amount bet from the betField.
234      *
235      * @return bet
236      */
237     public int getBet(){
238         return(Integer.parseInt(betField.getText()));
239     }
240
241     /**
242      * Enables or disables the hit button
243      *
244      * @boolean hit button state
245      */
246     public void hitEnabled(boolean en){
247         hit.setEnabled(en);
248     }
249
250     /**
251      * Enables or disables the stand button
252      *
253      * @boolean stand button state
254      */
255     public void standEnabled(boolean en){
256         stand.setEnabled(en);
257     }

```

```

258
259     /**
260      * Enables or disables the bet button
261      *
262      * @boolean bet button state
263      */
264     public void betEnabled(boolean en){
265         bet.setEnabled(en);
266     }
267 }

```

D.1.8 BJRules

```

1  //adding BJRules to SKCards
2  package SKCards;
3
4  /*
5   * importing the swing classes.
6   */
7  import javax.swing.*;
8
9  /**
10   * The rules of the Blackjack game.
11   *
12   * @author SK
13   * @version 2.0A
14   * @since 2.0A
15   * @SK.bug none
16   */
17  public class BJRules extends Rules{
18
19      public static final int PLAYER = 0;
20      public static final int DEALER = 1;
21
22      private boolean play;
23      int currentCard, pscore, dscore, bet;
24
25      public BJRules(String id, int cash, Pack setPack, String bck, String cs,
26                      SKMain m){
27
28          main = m;
29
30          ui = new BJUI(this);
31
32          pack = setPack;
33
34          players = new Player[2];
35
36          players[PLAYER] = new BJPlayer(id, cash, bck, cs, this);
37          players[DEALER] = new BJRLPlayer(this);
38
39          play = true;
40      }
41

```

```

42     public void newGame(){
43
44         clear();
45         ((BJUI)ui).hitEnabled(true);
46         ((BJUI)ui).standEnabled(true);
47         ((BJUI)ui).betEnabled(false);
48         players[PLAYER].newHand();
49         players[DEALER].newHand();
50         ((BJPlayer)(players[PLAYER])).hit(0, pack.nextCard());
51         pscore = ((BJPlayer)(players[PLAYER])).hit(1, pack.nextCard());
52         ((BJPlayer)(players[DEALER])).hit(0, pack.nextCard());
53         dscore = ((BJPlayer)(players[DEALER])).hit(1, pack.nextCard());
54         currentCard = 2;
55         display();
56     }
57
58     public void endGame(){
59
60         Card[] phand = ((BJPlayer)players[PLAYER]).getHand();
61         Card[] chand = ((BJPlayer)players[DEALER]).getHand();
62
63         for(int i = 0; i < phand.length; i++){
64             if(phand[i] != null){
65                 System.out.print(phand[i].toString());
66             }
67         }
68         System.out.println("_PLAYER_" + pscore);
69
70         for(int i = 0; i < chand.length; i++){
71             if(chand[i] != null){
72                 System.out.print(chand[i].toString());
73             }
74         }
75         System.out.println("_DEALER_" + dscore);
76
77         if(pscore > 21 && dscore <= 21){
78             JOptionPane.showMessageDialog(null, "Player_is_bust!");
79             ((BJRLPlayer)(players[DEALER])).reward(10);
80         }
81         else if(dscore > 21 && pscore <= 21){
82             JOptionPane.showMessageDialog(null, "Dealer_is_bust!");
83             players[PLAYER].setFunds(bet*2);
84             ((BJRLPlayer)(players[DEALER])).reward(-10);
85         }
86         else if(pscore > 21 && dscore > 21){
87             JOptionPane.showMessageDialog(null, "Both_bust!");
88             players[PLAYER].setFunds(bet);
89             ((BJRLPlayer)(players[DEALER])).reward(-10);
90         }
91         else if(pscore > dscore){
92             JOptionPane.showMessageDialog(null, "Player_wins!");
93             players[PLAYER].setFunds(bet*2);
94             ((BJRLPlayer)(players[DEALER])).reward(-10);
95         }

```

```

96         else if(pscore == dscore){
97             JOptionPane.showMessageDialog(null, "Push!");
98             players[PLAYER].setFunds(bet);
99             ((BJRLPlayer)(players[DEALER])).reward(5);
100        }
101        else{
102            JOptionPane.showMessageDialog(null, "Player loses!");
103        }
104        clear();
105        play = true;
106    }
107
108    public void hit(){
109        if(play){
110            pscore = ((BJPlayer)(players[PLAYER])).hit(currentCard, pack.nextCard());
111            currentCard++;
112            if(pscore > 21 || currentCard == 5){
113                display();
114                stand();
115            }
116        }
117        else{
118            dscore = ((BJPlayer)(players[DEALER])).hit(currentCard, pack.nextCard());
119            currentCard++;
120            if(pscore > 21 || currentCard == 5){
121                display();
122                stand();
123            }
124        }
125        display();
126    }
127
128    public void stand(){
129        System.out.println(play);
130        if(play){
131            play = false;
132            display();
133            currentCard = 2;
134            while(!play){
135                ((BJRLPlayer)players[DEALER]).learn();
136            }
137        }
138        else{
139            play = true;
140            endGame();
141        }
142        display();
143    }
144
145    public void bet(int amt){
146
147        if(amt <= players[PLAYER].getFunds()){
148            bet = amt;
149            players[PLAYER].setFunds(-1 * bet);

```

```

150     }
151     else {
152         JOptionPane.showMessageDialog(null, "INSUFFICIENT_FUNDS!!!");
153     }
154     newGame();
155     display();
156 }
157
158 public int getBJPoints(Card c){
159
160     int points = c.getPoints();
161
162     if(points == 10 || points == 11 || points == 12 || points == 13){
163         return 10;
164     }
165     else {
166         return points;
167     }
168 }
169
170 public boolean getPlay(){
171     return play;
172 }
173
174 public Card[] getPlayerHand(int p){
175
176     return ((BJPlayer) players[p]).getHand();
177 }
178 }

```

D.1.9 RLPlayer

```

1  //adding BJRLPlayer to SKCards
2  package SKCards;
3
4  /*
5   * importing the IO classes
6   */
7  import java.io.*;
8
9
10 /**
11  * The RL Player for the Blackjack game.
12  *
13  * @author SK
14  * @version 2.0A
15  * @since 2.0A
16  * @SK.bug none
17  */
18 public class BJRLPlayer extends BJPlayer{
19
20     private final double ALPHA = 0.6;
21     private final double GAMMA = 0.75;
22     private final double TAU = 0.5;

```

```

23
24     private static final int HIT = 1;
25     private static final int STAND = 0;
26     private static final int DEAL = -1;
27
28     private double [][] val;
29     private int lastaxn, lastscore;
30
31     /**
32      * Sole constructor. Takes in the rules and passes it to
33      * super(r).
34      *
35      * @param r the rules
36      */
37     public BJRLPlayer(BJRules r){
38
39         super(r);
40
41         val = new double [2][21];
42
43         try{
44             BufferedReader br = new BufferedReader(new FileReader("rlvals.txt"));
45
46             String data = br.readLine();
47             int i = 0;
48
49             while(data != null && i < val[0].length){
50
51                 System.out.println(data);
52
53                 String [] ds = data.split(",");
54
55                 val[STAND][i] = Double.parseDouble(ds[0]);
56                 val[HIT][i] = Double.parseDouble(ds[1]);
57
58                 data = br.readLine();
59                 i++;
60             }
61         }
62         catch(Exception ex){
63             ex.printStackTrace();
64         }
65     }
66
67     /**
68      * Starts a new hand.
69      */
70     public void newHand(){
71
72         lastaxn = DEAL;
73         super.newHand();
74     }
75
76     /**

```

```

77      * Plays a greedy policy, but updates the values as it plays
78      */
79      public void learn(){
80
81          if(score > 21){
82              return;
83          }
84
85          System.out.println("HIT_=" + val[HIT][score - 1]);
86          System.out.println("STAND_=" + val[STAND][score - 1]);
87
88          double maxval = Math.max(val[STAND][score - 1], val[HIT][score - 1]);
89
90          if(lastaxn != DEAL){
91              double rd = ALPHA * ((GAMMA * maxval) - val[lastaxn][lastscore - 1]);
92              val[lastaxn][lastscore - 1] += rd;
93          }
94
95          lastscore = score;
96
97          if(val[HIT][score - 1] > val[STAND][score - 1]){
98              rules.hit();
99              lastaxn = HIT;
100              System.out.println("HIT");
101          }
102          else{
103              rules.stand();
104              lastaxn = STAND;
105              System.out.println("STAND");
106          }
107      }
108
109
110      /**
111       * Allocates the final reward to the agent.
112       *
113       * @param r the reward
114       */
115      public void reward(double r){
116
117          double maxval = Math.max(val[STAND][lastscore - 1], val[HIT][lastscore - 1]);
118
119          if(lastaxn != DEAL){
120              double rd = ALPHA * (r + (GAMMA * maxval) - val[lastaxn][lastscore - 1]);
121              val[lastaxn][lastscore - 1] += rd;
122          }
123      }
124  }

```


Appendix E

Sample Blackjack Game

The following notation is used in the sample game:

[] - a face down card

— - Separation of two hands

Cards - Face values - 2 - 9 \Rightarrow 2 - 9, T \Rightarrow 10, J \Rightarrow Jack, Q \Rightarrow Queen, K \Rightarrow King

Suits - D - Diamonds, C - Clubs, H - Hearts, S - Spades

e.g. KS \Rightarrow King of Spades, 2D \Rightarrow 2 of Diamonds, TH \Rightarrow 10 of Hearts

```
Dealer AC []
Bet
Player [] []
Player bets $10, and takes an insurance
Dealer AC []
Bet $10 $10
Player AD AH
Player splits
Dealer AC []
Bet $10 $10 $10
Player AD QS | AH AS
(BLACKJACK)
Player stands on hand 1, splits hand 2
Dealer AC []
Bet $10 $10 $10 $10
Player AD QS | AH 2C | AS 4S
Player hits hand 2
Dealer AC []
Bet $10 $10 $10 $10
Player AD QS | AH 2C 4D | AS 4S
Player hits hand 2
Dealer AC []
Bet $10 $10 $10 $10
Player AD QS | AH 2C 4D KD | AS 4S
(Ace now counts ``low" (1))
Player hits hand 2
Dealer AC []
Bet $10 $10 $10 $10
Player AD QS | AH 2C 4D KD 5S | AS 4S
(BUST!)
Player doubles down on hand 3 (counting ace ``high")
```

```
Dealer AC []  
Bet $10 $10 $10 $20  
Player AD QS | AH 2C 4D KD 5S | AS 4S 6C  
Player stands on hand 3, dealer plays by turning over his second card  
Dealer AC JC (BLACKJACK!)  
Bet $10 $10 $10 $20  
Player AD QS | AH 2C 4D KD 5S | AS 4S 6C
```

The dealer has hit a Blackjack, so the player receives their \$10. The player loses the \$10 from hand 2 because they went bust. As for hands 1 and 3, the dealer and player have tied. However some casinos enforce the rule the “House wins a push”. Under this rule, the player would also lose the \$30. If the house treats pushes normally, then the player will get his \$30 back. Although there are many different variations on the rules, known as “house rule”, they are all well known and accepted by the players.

Bibliography

Parlett, D. (1994). *Teach Yourself Card Games*. Hodder Headline Plc, Euston Road, London, NW1 3BH.

Sutton, R. S. and Barto, A. G. (1998). *Reinforcement Learning: An Introduction*. MIT Press, Cambridge, MA, USA.