



Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

Lecture with Computer Exercises: Modelling and Simulating Social Systems with MATLAB

Project Report

Strategy optimisation in Texas Hold'em

Tim Weber, Jan Speckien, Patrice Gobat, Lionel Gulich

Zurich
December 2016

Agreement for free-download

We hereby agree to make our source code for this project freely available for download from the web pages of the SOMS chair. Furthermore, we assure that all source code is written by ourselves and is not violating any copyright restrictions.

Tim Weber

Jan Speckien

Patrice Gobat

Lionel Gulich

Contents

1	Abstract	4
2	Individual contributions	4
3	Introduction and Motivations	4
4	Model and implementation	6
4.1	Overview	6
4.2	Scripts	6
4.2.1	main.m	6
4.2.2	game.m	6
4.2.3	headsup.m	7
4.2.4	adujstCardValue.m	7
4.3	Learning models	8
4.3.1	Iteration model	8
4.3.2	Count and gauge model	9
4.3.3	Threshold model	9
5	Simulation Results and Discussion	10
5.1	Results of generic model	10
5.1.1	Generic model – Results	10
5.1.2	Generic model – Interpretation	11
5.2	Results of learning algorithms	12
5.2.1	Iteration model	12
5.2.2	Count and Gauge model	12
5.2.3	Threshold model	13
5.3	Comparison of learning algorithm performances	15
6	Robustness of results and influence of sigma	16
7	Summary and Outlook	16
7.1	Summary	16
7.2	Outlook and Limitations	17
8	References	18
9	Appendix	18

1 Abstract

In the recent years Texas Holdem has been increasingly popular amongst all age groups. Therefore, the aim of this study was to research how different strategies perform against each other in a heads-up.

Hence we wrote a program that simulates a game of Texas Holdem in which all feasible strategies compete against each other. In our model the strategy is represented by a single variable called risk factor that holds the willingness of taking risks.

In the model the strategy was constant over the course of a game. For the latter part one player could adjust his risk factor based on different learning algorithms, which then are compared with each other.

The results showed, that a very passive strategy always lost. This corresponds to the outcome of a real heads up, where the passive player has a disadvantage due to the blinds. For two players in general the less extreme strategy showed to be more effective. However, this characteristic could not be verified to match reality. Concerning the learning algorithms, a significant increase in the average number of games won by the learning player could be observed. This indicates that with an adaptive strategy an important increase in the average games won can be achieved, also in real Texas Holdem.

The simple model proved to be well suited to quickly observe changes in single variables like the blind value. It can also be easily expanded to observe other variables. Nevertheless, a comparison with actual game data is needed to verify the model.

2 Individual contributions

The commonly used models for Texas Hold'em are mostly based on stochastic probability, which makes the implementation of algorithms highly complicated. Thus we decided to write our own simulation in order to simplify the problem using as few variables as possible. All further results and conclusions are based on our program.

3 Introduction and Motivations

Texas Holdem is a variation of the card game Poker and is throughout casinos, tournaments and private people amongst the most popular of its kind. The goal of the game consists in winning all the money, or from here on called chips (virtual currency), from the other players at the table. As soon as someone has no more chips left, he lost and has to leave the table. Every round consists of different stages of

card dealing and then successive betting. In each stage a player can decide if and how much chips he wants to play according to some rules and by evaluating their private cards called hand and the cards open on the table named sequentially flop, turn and river. The total amount of chips from all the players in one round is called pot. Every round ends with either one player getting the whole pot or sometimes with a split of the pot, if two or more players have an equivalent score in the end of the round. The score is the highest possible combination of the private and public cards available according to a fixed ranking.

The game can be played with a varying number of people and for simplicity be subdivided into two parts: Group stage and Heads up, whereby the main difference lies in the number of people playing. In a heads up the last two remaining players are gambling for the overall victory and therefore every game of Texas Holdem will end with a heads up.

Every player has his own strategy, where some players like to play more conservative by not betting often and waiting for better cards and others play more aggressive by betting more often and trying to bluff his opponent. An important aspect of the game is the blinds, since they force the players to bet a fixed amount before the first cards are dealt. The blinds move forward by one player after every round. They guarantee that no player only can wait for the perfect cards and that there are always chips in the pot that can be won.

We all agree that Texas Holdem is about knowing your opponent and play smarter or better than him. Therefore it is not possible to have a fixed strategy. One has to be able to adjust his strategy depending on the way the opponent is playing and vice-versa.

Hence we all often wondered what would be good strategies to best adjust to ones opponent and in particular how effective such strategies are. To make a valuable assumption about the effectiveness of such a strategy, one would need to play an enormous amount of games to formulate a statistically speaking valid statement. If a whole game or at least some parts of it can be simulated repeatedly under different circumstances and with varying strategies, then one could gain deeper insights into Texas Holdem and possibly develop an applicable strategy or at least compare different approaches.

4 Model and implementation

4.1 Overview

This section serves to understand how we implemented the simulation. Expressions are introduced that we use throughout the report. One can find the whole program on GitHub¹ in the folder code. Figure 4.1 gives an overview of all scripts used and how they are linked.

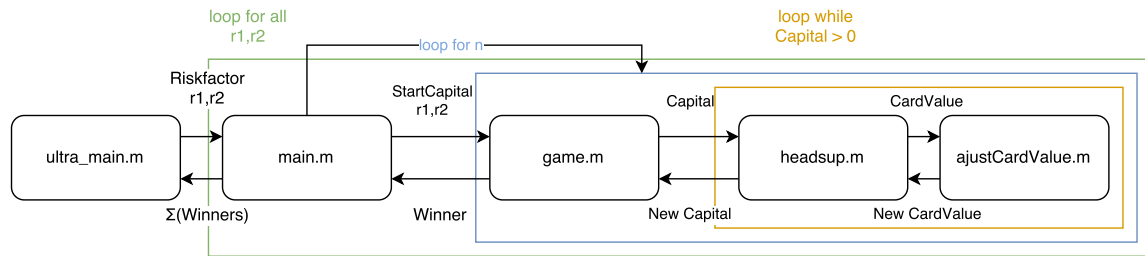


Figure 1: Overview of of the scripts used

4.2 Scripts

4.2.1 main.m

The script ultra-main plays all the different risk factors against each other and saves the amount of wins by player1 in a matrix.

The risk factor determines the character of a player. It is stored with a number between 0 and 1, whereas a smaller number represents a riskier player that plays a hand even if his score is not that high and vice versa.

The function main.m serves as an interface where one can decide on the settings through variables which then get passed on to the function game.

4.2.2 game.m

The function game.m is used to call the functions headsup/2. Like in real poker game.m continues to deal new hands until one player has run out of money. It is alternately calling two functions that differ only in which player has to pay the blind.

¹<https://github.com/atlas000/project-poker-msssm>

As described above, a blind is essential to the game of Texas Hold'em. A player has to pay some money into the pot without having seen his cards yet. That way if a player is very passive, he loses money anyway and has to become more active at some point.

In our model only a single fixed blind has been implemented, since in a heads up only the difference in blinds is relevant.

4.2.3 headsup.m

The headsup.m function simulates one hand which consists of up to four rounds. With every new card unveiled each player faces the decision whether to place a bet or not. Their action is determined by comparing their risk factors to their score. In total this scenario occurs up to four times per hand. If both players want to play until the end, a showdown determines the winner. Meaning that they have to show their cards and the one with the better score wins the hand and receives the pot. In the first round, the "preflop", one player has to pay the blind and each player gets a set of playing card.

The main aspect making poker a complicated game to simulate is its use of cards. Hence our main goal was to find a model, which did not require the implementation of a deck of cards. In order to achieve a similar effect every player is given a random number, called "card value". The number reaches from 0 to 1 and represents the quality of the hand, whereas 0 would be the worst hand and 1 the best.

In the next three rounds the "flop", "turn" and "river" new cards get unveiled on the table. So the score of the players' hands has to be adjusted.

4.2.4 adujstCardValue.m

At the heart of this simulation lies the adjustCardValue function, which recalculates the card value of a player.

Important criteria are that the cards correlate with each other, so if a player has good cards it is likely that his hand is still good when a new round is played. Also it should be possible for a value to reach the whole spectrum of card quality from 0 to 1. This is important because if one has a bad starting hand it should still be possible, with some luck, to get the best hand possible after the flop. Hence every value should be possible to receive but not with the same probability.

4.3 Learning models

An important part of real Texas Holdem is to obtain an advantage by estimating the strategy of your opponent and to react accordingly. In this model two different approaches for learning algorithms have been chosen to optimise a player's strategy, under the assumption that the opponent keeps his risk factor constant.

- One approach is to minimise the losses and thus maximise the profit of a player. This model analyses a player's own performance.
- The other approach was to first analyse the play of the opponent and then choose the optimal strategy to best counter the other player. This of course is only possible if all outcomes for all strategies are known.

4.3.1 Iteration model

This learning algorithm analyses the loss of the last hand, because if the last hand has been won there is no need to adapt the strategy. So after every hand one of the following scenarios takes place:

Validation	Interpretation	Action
$\text{Capital}(t+1) \geq \text{Capital}(t)$	The player had equal or better cards, or the hand has not been played	risk factor r will not be changed
$\text{Capital}(t+1) == \text{Capital}(t) - 4$	The player has lost after the showdown and is thus too aggressive	$r=r+0.012$
Else (Capital dropped by 1,2 or 3)	The player is too passive, since it did not come to a showdown	$r=r+0.003$

4.3.2 Count and gauge model

The main goal behind analysing the opponent's way of playing is to correctly estimate his risk factor, in order to adapt one's risk factor afterwards. This algorithm is accomplished in two consecutive steps:

Firstly it keeps track of the number of times the opponent decides to play the initial round when he is not forced to do so by the blind. The ratio between this value and the total amount of rounds played can then be used as a measurement to estimate the opponent's risk factor:

$$\text{risk factor} \approx 1 - \frac{\# \text{ rounds}}{\# \text{ rounds played}}$$

The second part consists of finding the optimal counterpart to the estimated risk factor. For that reason we evaluated results obtained by previous simulations. The data obtained indicates that for every risk factor chosen, there is at least one optimal counterpart. After extracting these values the algorithm adjusts the own risk factor.

4.3.3 Threshold model

The target of this learning model is to determine the risk factor from the opponent, which then enables to find the optimal risk factor for the learning player. Every time a showdown is conducted the learning player gets to know the score from his opponent. Because the opponent has continued the hand with this score until the showdown, it is known, that the risk factor from the opponent has to be lower than said score. With every showdown executed the possible values for risk factor of the opponent can be narrowed down, leaving a threshold which has to be undershot by the opponent's risk factor.

The optimal risk factor for the learning player can then be determined with the data obtained by previous simulations. This data includes the win probability of both players for all possible combinations of risk factors between both players. Determining the optimal risk factor then only is a matter of finding the maximum for a given opponent strategy.

5 Simulation Results and Discussion

5.1 Results of generic model

5.1.1 Generic model – Results

Figure 2 shows how many games player1 has won out of a total of 1000, when both players started out with an equal amount of chips (50) and a fixed blind equal to the betValue of 1. The colour of every cell ($x = r, y = r2$) shows the number of games won by player1 with his risk factor $r1$ against player2 with a risk factor $r2$.

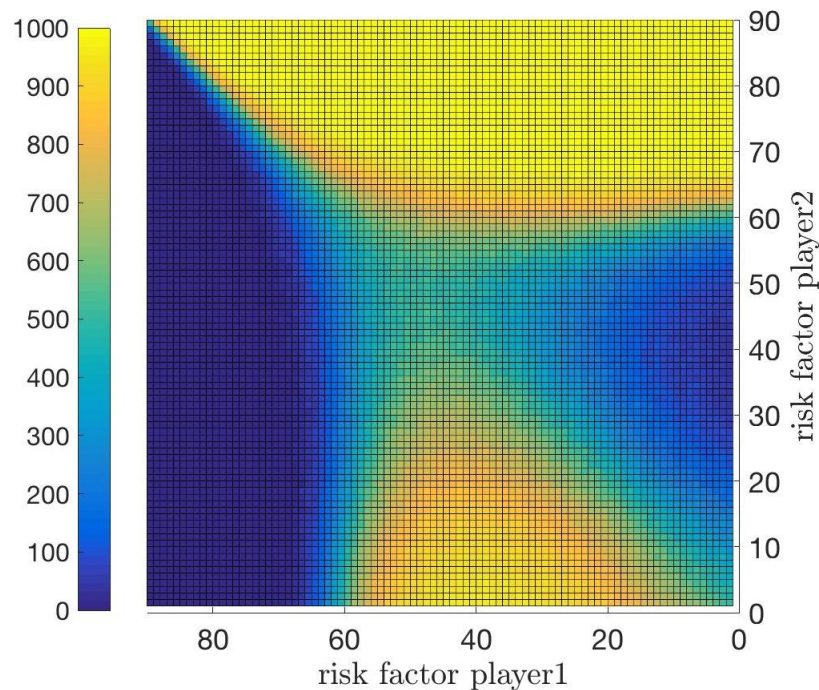


Figure 2: Games won by player1 out of 1000

On the diagonal ($r1 = r2$) a green line can be seen, which shows that player1 and player2 have each won approximately half of all games played. Around this line the games won by a player are approximately anti-symmetrically distributed, which means that both players did equally well when using the same risk factor.

A very passive strategy ($r > 0.63$) could not win any games, except against a more passive player. In this case, the less passive strategy occurred to be superior.

An aggressive strategy ($r < 0.5$) performed better. It could win against the very passive strategy ($r \text{ opponent} > 0.63$). However, it clearly lost against all less aggressive strategies under 0.63 ($r < r \text{ opponent} < 0.63$).

Thus, there are two different patterns visible in Figure 2. One pattern ($r < 0.5$), where the more passive player wins. And a second pattern ($r > 0.63$), where the more active player wins.

Between these two patterns lies a transition area, where both risk factors are between 0.5 and 0.63. In this area the games were very balanced and not even an optimal counter-strategy did win more than 60% of all games.

In the transition area the change from where it is favourable to play passive on to playing active occurs.

5.1.2 Generic model – Interpretation

The bad performance of the very passive strategy can be explained with two independent mechanisms. One being the blind that punishes the more passive player. The other being the card score dropping below the risk factor before the showdown, leading to a loss of all chips bet during that hand.

The aggressive player plays more games with weak hands, which means that he will lose more often in the case of a showdown.

Inside the area, where the game is balanced, all mentioned effects compensate each other.

The model weights the effect of a passive player dropping out after already having invested chips to heavily. In a real game of poker a player stays in the game, due to a comparison between his own odds and the pot. This effect is called pot odds.

Otherwise the simulation shows that the aggressive strategy (loose play) is favourable against an extreme passive strategy (tight play). This reflects the reality where blinds give the passive player a disadvantage in a heads up.

5.2 Results of learning algorithms

5.2.1 Iteration model

The algorithm needs around 800 hands in order to lead to a good estimation for the optimal parameter r . This leads to good results for start capitals above 70 and to very good results at start capitals above 300 even with a very unfavourable starting strategy like $r(0)=0.1$.

Advantages

The model has a great potential for different kinds of optimisation:

- Parameters used in the model can still be optimised in order to perform quicker
- More cases for different scenarios could be made (different action for loss of one and loss of two)
- The cards that the player could have taken into account in order to avoid a change in policy if he only had bad luck

Disadvantages

The algorithm is based on a simplified scenario, where the only reason for a loss is based on the strategy and never on luck. This leads in some cases to a change of policy even when it is not desired. This can be observed in the high standard deviation, even after having found a well estimated r . This behaviour leads to the algorithm loosing against extremely strong strategies due to its continuous slight deviations from the optimum risk factor.

5.2.2 Count and Gauge model

The algorithm leads to a significant increase in the number of games won. Unfortunately in the region between 0.5 and 0.7 it was rather contra productive. In this region an adjustment of the otherwise solid algorithm is needed. We relate this event to the transition from the "passive area to the active area. The impact of luck can easily shift the estimation of ones opponent's risk factor into the wrong area which ultimately ends by choosing the wrong counter strategy.

Advantages

- With an increasing number of rounds the estimation continuously becomes better
- On a long run he will always play the most suitable strategy and optimise his return
- It reaches a good estimation of the opponent's risk factor already after a few rounds
- The algorithm could also use the rounds where the other player has to pay the blind by counting the number of consecutive rounds he plays and applies the same strategy just by including the models standard deviation into the calculation. This would lead to a faster estimation of the opponent's risk factor

Disadvantages

The algorithm requires a precise knowledge of all possible outcomes and therefore is only applicable in the same circumstances the data was obtained

5.2.3 Threshold model

Figure 4 shows that the resulting wins are independent from the learning player's risk factor because results are nearly constant along a horizontal line. It is visible that the algorithm does not have the same effectiveness for all opponent-strategies. For an opponent risk factor $r \in [0, 0.4]$ and also around 0.8 the algorithm is really successful, as can be seen with the highly yellow areas in the plot, which represent a lot of wins for the learning player. From the green band at an opponents risk factor of around 0.45 it can be seen that wins are equally balanced around this point. In the blue area for an opponents risk factor between 0.5 and 0.75 wins for the learning player are infrequent. This mean that there the algorithm fails to provide an advantage over the opponent.

By subtracting the resultant matrix from generic model with blinds from that of the learning algorithm it can be visualized where the threshold learning algorithm improves the game results and where it fails to do so. This is shown in figure 3. The highly yellow areas show a drastic improvement in performance, whereas the green

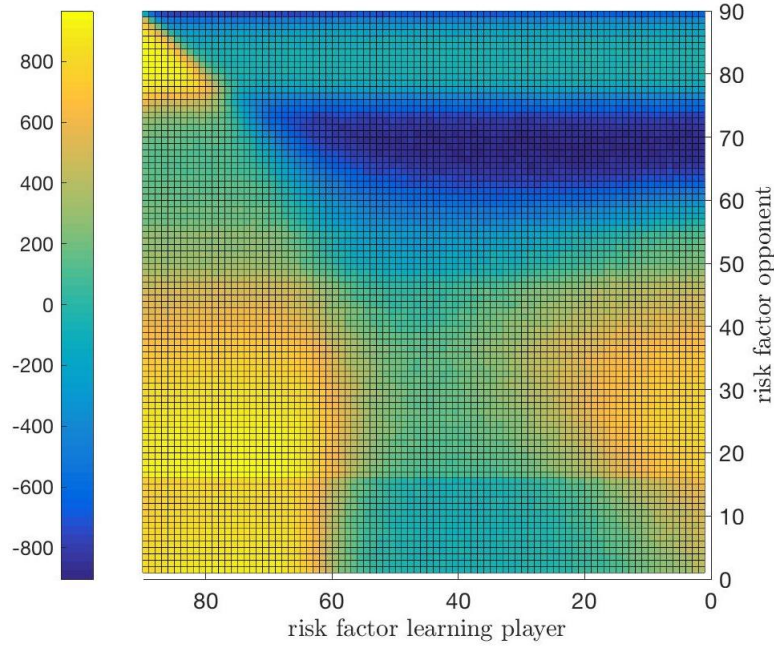


Figure 3: Topview of difference from results of threshold learning model and generic model with blinds

areas represent scopes where the performance could only be slightly changed or even remained the same. The deep blue areas show scopes in which the learning players performance has declined drastically.

The behaviour of the learning model in the green and yellow areas can be explained as follows: since the performance from the learning player in the green areas was already good, with the generic model no big improvement could be realised. The difference in wins therefore levelled off in these sections. The yellow areas underlay low performance regions in the generic model with blinds, thus a big improvement can be realised and the difference in wins rocketed.

What remains to explain are the not expected blue areas. In the blue area at a risk factor of 0.9 the opponent plays extremely passive. Hence a showdown seldom occurs and the learning algorithm cannot take effect. However this can only explain the failure of the learning aspect of the algorithm but not the loss of the learning player. The other blue area's behaviour can be explained by the addition of two effects. Firstly the gradient of the generic result matrix is extremely steep, thus already a slight misestimate in the opponent's risk factor can lead to a big change

	# Games won		# Hands played		Δ games won	
	μ	σ	μ	σ	μ	σ
Iteration model	647.55	335.5068	110020	748600	147.7569	198.3177
Count and gauge model	690.6369	343.1811	709.3543	342.5163	190.8446	477.6301
Threshold model	605.7777	301.3260	1565	1854	105.9853	501.7773

Table 1: Learning models overview figures

in the optimal risk factor for the learning player. Secondly because in the turn the opponent can also play a hand up to 0.15 lower than his risk factor the estimation for his risk factor can also be shifted downwards by up to 0.15 and therefore lead to a miscalculation in the optimal risk factor for the learning player.

5.3 Comparison of learning algorithm performances

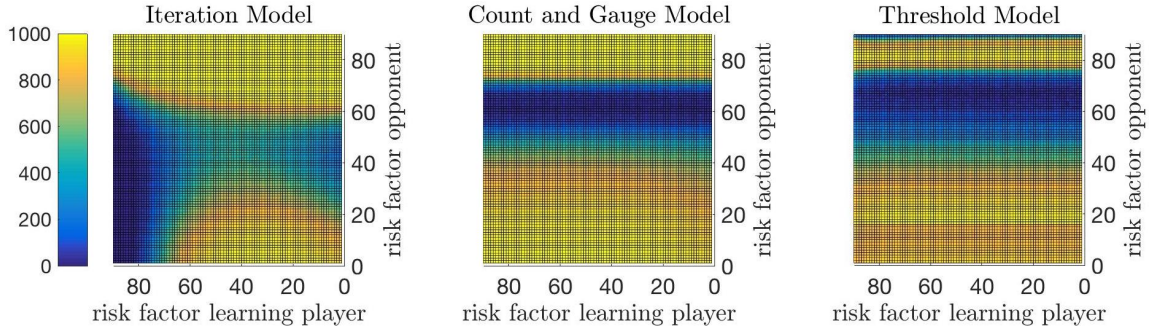


Figure 4: Result-matrices for all three learning models. From left to right: threshold, count and gauge, iteration

In table 5.3 some characteristic values for the learning algorithms are displayed. Therefore the mean and the standard deviation of the resulting-wins-matrix, total-hands-played-matrix and the difference-in-wins-matrix with and without the learning algorithm have been calculated. By comparing the mean of games won it can be seen, that the count and gauge model is the most successful algorithm in terms of most wins created. Yet due to the high standard deviation in games won, also the most inconsistent. From the standard deviation of difference in games won it can be distinguished that the iteration model is the most balanced overall, leading to small improvements for all different opponent strategies. In comparison thereto the threshold- and count and gauge model are not successful against every opponent,

but if they are, they create a high yield.

Lastly the time efficiency has to be considered: Clearly the count and gauge algorithm is the one leading to a victory the quickest, as its low mean in hands played shows. On the other hand the iteration is the most time consuming. While being the most effective for every opponent its enormous mean amount of hands played leaves it to be essentially useless in practice.

6 Robustness of results and influence of sigma

To draw some conclusions from our model we had to ensure that our model is robust enough. We performed thousand games a thousand times and found a standard deviation of 7.7 for the wins.

To the function `adjustCardValue.m` a sigma for the flop is needed that leads to a big variety of outcomes because three new cards are dealt in that round. Still the difference in the card value before and after new cards cannot be too wide, otherwise all correlation is lost. In the next two rounds, the turn and the river, where only one card is dealt the current card value has to be somewhat stable.

We chose the combination of 0.4/0.2/0.2 to be the one that represents a game of Texas Hold'em the best.

We ran some simulations with different sigma and saw that for each combination the resulting graph has the same shape but is shifted to the left or the right of the x and y scale.

Sigma have been tested on their reach from 0 to 1 and lowest/highest possible value for each risk factor. The results can be found on GitHub².

7 Summary and Outlook

7.1 Summary

The most important understanding we gained from our simulation was definitely that there is a transition between when it is favourable to play aggressively and when to play rather passively. This can be translated into the applicable technique that if your opponent plays less than every third hand, one will win more, if one plays more aggressively. On the other side of this barrier, if the opponent plays more than every

²<https://github.com/atlas000/project-poker-msssm>

second hand, playing a bit more conservatively will increase winnings in the long term. This rule can easily be implemented by just counting the number of times the opponent pays the blind when he is not forced to, divided by the amount of times it was not the case.

Our learning algorithm overall increased the performance, but all encountered big problems in the transition region. This reflects the complexity of poker even under simplified conditions and that there is no magical formula even when the opponent plays extremely balanced.

Nevertheless we learned from our approach that by keeping track of the opponent's decisions, there can be a lot of insights gained about his strategy. After a few rounds played then the easy rule mentioned above can be applied to gain an advantage. In the case that the estimated risk factor of the opponent lies within the transition area, the result has to be treated with caution, as the stubborn execution can lead to the opposite of the desired results as proved by the count and gauge algorithm. Although the findings are plausible a further comparison with real game data is still needed to verify the results.

7.2 Outlook and Limitations

This project is developable in many kind of ways. We thought of some different parameters one could implement:

One of them was a blind that would increase over time. This is how it is played in real Texas Hold'em with the motive to punish passive play. Because of that it would have been interesting to observe what difference it would have made in our simulation, also since our games lasted up to over a thousand hands and this is far away from reality.

One could also have tweaked the capital/betValue ratio in general. Possible outcomes could have been that with a certain ratio the aggressive player gets punished because his losing hands include bigger pots or the conservative player loses more often because games are more fast paced.

Another extremely interesting feature could have been giving the players the option to decide with how much they wanted to raise the pot. Currently, one chip is the only option but how would the outcome of a game be if they could bet a larger and varying amount. If that was the case, card evaluation needed to be different. A player then not only has to take his own cards in to the decision to play or drop out

but also how much money his opponent has bet in this round. Furthermore, would this change have opened the opportunity for a player to bluff his counterpart, this means pretending to have better cards than what he actually has. Now the more active player might get rewarded for his play style.

The possibility of checking could also have been included. Checking is doing nothing when your opponent has not bet any chips. This variation would need a distinct playing order and that would put the one, who has to show his cards first in a show-down, into a disadvantage. All these factors could play a role in developing a winning algorithm.

Other variations that certainly would lead to different discoveries but would make the program much more complex are coding "real" cards with their actual probabilities to show up in a hand but also including more than two players which might change the power dynamics between active and passive players.

8 References

All work was achieved by ourselves.

9 Appendix

The appendix and further files can be found on GitHub³.

³<https://github.com/atlas000/project-poker-msssm>