

A Statistical Exploitation Module for Texas Hold'em

And It's Benefits When Used With an Approximate Nash Equilibrium Strategy

Kevin Norris

Department of Computer Science
University of Auckland
Auckland, New Zealand
knor031@aucklanduni.ac.nz

Ian Watson

Department of Computer Science
University of Auckland
Auckland, New Zealand
ian@cs.auckland.ac.nz

Abstract— An approximate Nash equilibrium strategy is difficult for opponents of all skill levels to exploit, but it is not able to exploit opponents. Opponent modeling strategies on the other hand provide the ability to exploit weak players, but have the disadvantage of being exploitable to strong players. We examine the effects of combining an approximate Nash equilibrium strategy with an opponent based strategy. We present a statistical exploitation module that is capable of adding opponent based exploitation to any base strategy for playing No Limit Texas Hold'em. This module is built to recognize statistical anomalies in the opponent's play and capitalize on them through the use of expert designed statistical exploitations. Expert designed statistical exploitations ensure that the addition of the module does not increase the exploitability of the base strategy. The merging of an approximate Nash equilibrium strategy with the statistical exploitation module has shown promising results in our initial experiments against a range of static opponents with varying exploitabilities. It could lead to a champion level player once the module is improved to deal with dynamic opponents.

Keywords— artificial intelligence, game AI, opponent modeling, nash equilibrium, exploitation, poker, Texas Hold'em

I. TERMINOLOGY

The term "Opponent Modeling" has been ambiguously used in the literature to mean playing an opponent based strategy. Intuitively one would think that the opponent model would be the collection of information about the opponent that represents the way in which the opponent plays, a model of the opponents play style. From this definition opponent modeling would just be the creation of this model, but in the literature it has been used to also incorporate the way in which the model is used to play an opponent based strategy. We believe that for ease of understanding, the ambiguous usage of the term "Opponent Modeling" should be broken into three terms:

- **Opponent Model:** the set of information that represents the opponent's play style.
- **Opponent based actions:** the way in which you use the given Opponent Model to determine actions that are suited towards the opponent you are facing.
- **Opponent based strategy:** The combination of creating an opponent model for your opponent and

using this model to determine opponent based actions. There are two types of opponent based strategies that can be employed:

- **Fully opponent based strategy:** The agent uses the opponent model to play an opponent based action for every action.
- **Partially opponent based strategy:** The agent plays opponent based actions only some of the time.

II. INTRODUCTION

Many online poker players use Heads up Displays (HUD). A HUD shows the player the recorded frequency statistics for each of his/her opponents and him/her self. Many novice online players play a relatively static strategy as they learn the game, which they alter only if they see anomalies in the statistics of their opponents. These anomalies indicate their opponent is playing sub-optimally. When they see these anomalies they devise statistical exploits: ways to alter their strategy to take advantage of the statistical anomaly such that they increase their win rate against the opponent. The statistical exploitation module presented in this paper is built around the idea of recognizing statistical anomalies and allowing the agent to capitalize on them through the use of expert designed statistical exploitations in the domain of heads up no-limit Texas Hold'em poker.

The module alone is not a complete strategy, it is designed to be added to a base strategy. The resulting strategy is an enhanced version of the base strategy which is able to capitalize on statistical anomalies, exploiting opponents. We decided on using an approximate Nash equilibrium strategy as our base strategy for testing the module. A Nash equilibrium strategy is described as "a strategy for each player of the game, with the property that no single player can do better by changing to a different strategy" in [1]. Johanson [1] reasons that while playing such a strategy one can do no worse than tie the game, since the opponent cannot do better by playing a strategy other than the equilibrium. Johanson [1] recognizes that "using such a strategy allows us to defend against any opponent, or allows us to learn an opponent's tendencies safely for several hands before attempting to exploit them".

Finding a Nash equilibrium in a complex game is very difficult, so instead they are approximated, providing a suboptimal strategy, a strategy which is exploitable. However the closer the approximation, the closer the strategy is to being un-exploitable, therefore, approximate Nash equilibrium strategies are very difficult to exploit. We felt that such a defensive strategy would be the ideal base strategy. We thus tried to preserve the desirable aspect of the base strategy, being difficult to exploit; and added the ability to exploit opponents, producing a strong agent that can do well against both weak and strong opponents.

The combination of a Nash equilibrium strategy and exploitation abilities was also done in the poker agent Polaris created in [1]. We will give an overview of Polaris. The tradeoff between *difficulty to exploit* and *exploitability of opponents* is the main challenge faced when trying to add exploitation abilities to a Nash equilibrium strategy. In section VIII we will compare our technique of combining exploitation abilities with a Nash equilibrium strategy to the technique used by Polaris, and discuss why our statistical exploitation module dealt better with this challenge than the Polaris agents.

Section III provides a brief overview of the game of No Limit Texas Hold'em, which our module has been designed for. In section IV we discuss related work, particularly in the area of combining a Nash equilibrium strategy and opponent exploitation abilities. Section V discusses the role of frequency statistics in opponent modeling, how they are used by many opponent based strategies, and how our usage of them differs. Section VI provides an overview of the module's design and a look into the various parts of the module. Then, we illustrate our experimental methodology and present our results in Section VII. Section VIII provides conclusions and an explicit comparison between our technique for adding exploitation to a Nash equilibrium strategy versus the technique described in [1]. Section IX discusses avenues for future work.

III. NO LIMIT TEXAS HOLD'EM

We describe briefly the game of Texas Hold'em focusing on two-player no limit Hold'em as our module has been specialized for this domain. If a game consists of only two players, it is described as being a heads-up match. The game of heads-up no limit Texas Hold'em consists of four stages: pre-flop, flop, turn and river. During the pre-flop stage each player is dealt two hole cards, which only they can see. Two forced bets are contributed to the pot, these being the small blind (SB) and the big blind (BB) before any betting takes place. The big blind is usually double the value of the small blind. In the game of heads-up Texas Hold'em the dealer contributes the small blind and the non-dealer contributes the big blind. The dealer signifies the player who is first to act during the pre-flop stage of the game and last to act for each of the other stages of the game. The betting actions, which are common to all variations of poker, are as follows:

- **Fold:** When a player abandons his hand, no longer committing any chips to the pot and giving up any right to contest the chips that make up the pot.
- **Check/Call:** When a player commits the minimum amount of chips with which he/she is able to continue to contest the pot. A check requires zero chips to be committed, and a call requires an amount greater than zero to be committed.
- **Bet/Raise:** When a player commits a larger number of chips than the amount necessary to continue to contest the pot, this is known as a bet. If a player is in the position where he/she must call a bet to continue, but then decides to invest more than the call amount in the pot, this is known as a raise.

In a no limit game a player may bet any amount they desire up to the total value of chips they possess. Once the betting in one stage of the game is complete and as long as no players have folded, play continues to the next stage. Each further stage after the pre-flop stage involves the drawing of community cards from the shuffled deck of cards as follows:

- **Flop:** 3 community cards
- **Turn:** 1 community card
- **River:** 1 community card

In a standard heads-up no-limit poker game the chip stacks of each player would fluctuate between hands depending on who won the previous hand. To reduce the variance of this structure a variation known as Doyle's Game is played during our experiments where the starting stacks of both players are reset to a specified amount at the beginning of every hand.

IV. RELATED WORK

A. Polaris

Polaris is a collection of techniques for creating poker agents, which are described in [1]. Johanson [1] outlines new approaches to calculating approximate Nash equilibrium strategies and abstract game best response strategies called Counterfactual Regret Minimization and Frequentist Best Response respectively. A technique called Restricted Nash Response is also illustrated that creates an agent which is a compromise between the two previously mentioned techniques. The last technique presented in [1] creates a meta-agent which is made up of agents created from the previous techniques, along with a "coach" agent which decides which agent to use for every hand. We will delve deeper into the last two techniques as they are the ones which deal with combining an approximate Nash equilibrium strategy with exploitation abilities.

1) *Restricted Nash Response (RNR)*

Counterfactual Regret Minimization creates an agent that is difficult to defeat, but is unable to exploit its opposition, thus does not win very much. Frequentist Best Response creates an agent that can exploit specific opponents, but is easily defeated by opponents it is not designed to defeat. The Restricted Nash Response technique uses Regret Minimization to find a compromise between these two extremes, creating agents that exploit particular opponents or classes of opponents and still providing a bound on their exploitability. This strategy is constructed by finding a Nash equilibrium in a restricted game, where the opponent must play according to a fixed strategy with probability p . p is chosen when creating the strategy and determines the proportion of time the opponent must use the fixed strategy. p ranges between zero and one: a p of zero means that the opponent never plays the fixed strategy so a Nash equilibrium is computed; and a p of one means that the opponent only uses the fixed strategy so a best response is computed. All values for p between zero and one represent a tradeoff between exploitation and exploitability.

2) *Meta-Agent*

In competitions the opponents will be unknown. So which of the previously discussed agents should be used against each opponent? Each of the previous agents have their pros and cons. In order to obtain the benefits of each, [1] created a team of several agents of various types, thereby creating a meta-agent. The problem faced by the meta-agent in regards to which agent to choose when selecting an action is solved by expert algorithms. In this case the algorithm UCB1 was chosen. UCB1 is designed to trade off exploration and exploitation when choosing the agent. Since there are various types of agents in the team, the UCB1 algorithm uses different costs of exploration for different types. For example the cost of exploring the use of a Frequency Best Response agent is very high, whereas the cost of exploring a Restricted Nash Response agent is lower and the cost of using the Nash equilibrium agent is the lowest. Johanson [1] found that using a team of agents provided better results than using only one of its parts.

3) *Discussion*

A single RNR agent performs worse against arbitrary opponents overall in comparison to an agent using an equilibrium strategy. RNR agents are only able to exploit opponents that play strategies similar to the strategy they were trained against and are easier for opponents to exploit. This means that the majority of opponents the RNR agent plays will not be exploited, and the opponents will be able to exploit the RNR agent better than they would be able to exploit an equilibrium strategy. In the experimental results shown in [1] a team of RNR, and a team of Frequentist Best Response (FBR) agents played 4 opponents that they had a counter-strategy against and 2 opponents that were unknown. The FBR team performed worse than an equilibrium agent against the opponents, while the RNR team performed better. There is a cost associated with exploring the various agents in the team.

The FBR team performed worse because the cost of choosing the wrong FBR agent is high, since FBR agents are easily exploited by opponents that they are not designed against. The RNR team performed better because the cost of exploring suboptimal RNR agents is low, since the agents play strategies that are close to an approximate equilibrium strategy.

Creating a meta-agent using this technique that performs better than an equilibrium strategy on average over arbitrary opponents, should be relatively straightforward. A number of RNR agents along with a Nash equilibrium agent should perform better than a Nash equilibrium agent alone. The counter strategies the RNR agents employ will allow this meta-agent to exploit certain play styles. If the meta-agent comes across a play style that it does not have a counter-strategy for, it should simply play the Nash equilibrium strategy. The RNR agent's strategies are similar to the strategy employed by a Nash equilibrium agent and therefore also similar to one another's strategies. Due to this it may take the "coach" agent a long time to determine which agent is best used against an opponent. The more RNR agents in the meta-agent the longer the exploration phase would take. Although the cost of exploring an RNR agent is low it does add up, and if the exploration phase is long enough, may cause the meta-agent to perform worse than a single Nash equilibrium agent. Thus it may be difficult to create a meta-agent that is able to exploit a large number of varying strategies.

V. FREQUENCY STATISTICS

Poker is an incomplete information game as well as a stochastic game, making it very difficult for a player to know where they stand in a hand. How can a player make good decisions without knowing what the opponent has, or what they will do? If a player wishes to make decisions based on the opponent they must have knowledge of the opponent. This knowledge comes in the form of history, the history of actions that the opponent carried out in their previous hands. Statistical opponent models record and use past play by keeping statistics of the opponent's action frequencies for use in later similar situations to improve the quality of decision making and to maximize profit. Keeping statistics on the opponent's action frequencies gives the player a good indication of how the opponent will act and react in given situations. Such indication follows from the probability distributions for the possible actions an opponent can make in each situation.

Schauenberg [3] describes two problems faced when using frequency statistics for opponent modeling. The first problem is that without enough observations of a situation the action distribution statistics can be largely inaccurate. There are two possible solutions to this problem: using priors, or not using the statistics to impact the action decisions until the game states they correlate to have been observed enough times to make their probability distributions statistically significant. The second problem is that opponents often alter their strategy

throughout the game, invalidating the statistics that have been observed. If the opponent's play does not match the frequencies the player is using to make their decisions the opponent will be exploiting the player's incorrect information, causing them to lose money. This problem can be solved through decaying history, a technique through which the statistics are only affected by the latest x hands. x needs to be set to a small enough number, so that, if the opponent changes his/her strategy, the statistic will depict the alteration quickly.

The statistics used in opponent modeling research are usually limited in number and simplistic in nature, each generalizing greatly over a large set of game states. Many researchers are only using simple statistics such as: raise, call, fold percentage per street, VPIP (voluntary put money in pot) how often one calls or raises pre-flop, PFR(pre-flop raise) how often one raises pre-flop and aggression factor: the ratio of the number of times the player is aggressive vs. the number of times he is passive postflop. Opponent model research that has utilized only such simple statistics includes [3], [4], [5], [6], and [7]. Adding in more statistics, while possibly making opponent models more expressive and better at capturing the nuances of an opponent play style, will also increase the computation and complexity of the program. Such complexity is needed to extract relevant data from the statistics to make betting decisions and leads to the tradeoff between program complexity and required computation verses opponent model expressivity.

HUD's currently prevalent in online poker, provide a full observation model. They provide the player with any thinkable statistic because they save every hand. Online players have been using HUD's for some time and have discovered many statistics, not generally used in opponent modeling research, that are very helpful for determining an opponent's play style. We believe opponent models can be greatly improved by adding some of the statistics that have become popular in online play through the HUD's. For example: **3bet**, i.e. how often the opponent re-raises a pre-flop raise, and fold to 3bet; and **cbet**, i.e. how often an opponent bets the flop after being the last to raise pre-flop and fold to cbet; and various other statistics that are now a staple for online player's HUD's.

The notion is: if statistics are specific to a particular situation, they will surpass generic statistics in predicting the opponent's action accurately. For example, in the case where the opponent was last to raise pre-flop and is first to act on the flop, a statistic for cbet percentage will be more accurate in determining how the opponent will act than the generic fold, check, bet flop percentages. Online players understand this and populate their HUDs with statistics for situations which occur frequently and are frequently played incorrectly. Much of the research on opponent modeling, however, uses generic frequency statistics, instead. The reason for this is that much of the opponent modeling research has been conducted on fully opponent based strategies, which must use the statistics when making every action decision. Researchers usually keep

frequency statistics general, thereby reducing the complexity of determining which statistics to use in order to impact the action decision, and the weight that each statistic will be given. If a fully opponent-based strategy were to use situation specific frequency statistics, the logic to determine which statistics to use and their impact on the action decisions would have to be largely altered. One would have to consider which of the specific frequency statistics to use and how to weight their impacts on the action decision for every possible game state. Although the use of specific frequency statistics could increase the accuracy of the opponent models used by fully opponent based strategies, it would make the agents much more complex, and would be difficult and time consuming to implement.

Consequently, our statistical exploitation module is based around the use of specific frequency statistics. Each exploitation uses a situation specific statistic along with some generic statistics to determine if the opponent is playing that situation incorrectly and if it can profitably exploit this. The module provides a partially opponent based strategy, using frequency statistics to determine an action only if an exploit is available for the situation. Unlike a fully opponent-based strategy it does not have to use frequency statistics for every situation; rather, it plays the base strategy in situations where an exploit does not exist or an exploit does not apply. The workings of the statistical exploitation module will be presented in more detail in the next section.

VI. MODULE SYSTEM DESIGN

We have created a module that we have added to SartreNL, presented in [8], that exploits players through the use of statistical exploits. This produces an agent which has a significantly increased win rate, without increasing its exploitability. This was accomplished by creating a statistical model that records detailed frequency statistics of an opponent in many contexts, and a number of exploits that provide highly profitable actions in the situation they apply to. The module also includes an opponent exploiter that tracks the exploits that apply to a given opponent model at any given time and provides the underlying agent with actions from the exploits when a situation arises in which one of the exploits applies. While we have used SartreNL as the underlying agent, this module could be added to any other underlying agent by making minor changes to the opponent exploiter and the chosen underlying agent.

Fig. 1 depicts the manner in which the individual parts work together to provide the addition of a partially opponent based strategy to SartreNL. Partially opponent based because the actions determined by the agent/module combination are not always dependent upon the module's opponent model; they are only dependent on the module's opponent model when an exploit applies to both the current game state and the module's opponent model. Our opponent model is based on frequency statistics and thus faces the two problems discussed

in section V. Without enough observations of a situation, the action distribution statistics can be largely inaccurate. This problem has been overcome by our agent by playing a partially opponent based strategy, if the frequency statistics are not significant the agent plays the base strategy. Our agent has not yet overcome the problem of opponents often altering their strategy throughout the game, invalidating the statistics that have been observed. In section IX we discuss a way in which we aim to solve this problem.

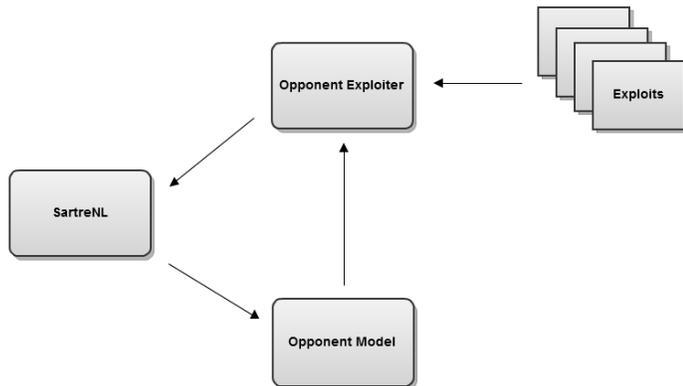


Fig. 1. Model of exploitation system design.

SartreNL is used as the underlying poker agent because it plays an approximate Nash equilibrium strategy, providing the trait of being difficult to exploit. SartreNL uses the opponent imitator described in [9] to play the style of the player whom its case-base was trained on. In the case of SartreNL the case-base was trained on the Hyperborean agent’s hand histories from the 2011 AAAI Heads-up No-Limit Texas Hold’em competition. This agent plays an approximate Nash equilibrium strategy, so the approximate Nash equilibrium strategy should emerge in SartreNL’s play.

A. Opponent model

The opponent model is a collection of counters and variables representing the frequency statistics. After each hand is played out, the counters are updated and then the statistics are recalculated based on the updated counters. In the poker framework used in the Annual Computer Poker Competitions (ACPC), messages are passed to and from the agents. The messages sent to the agents include all of the contextual information from the game and the messages sent from the agents contain the actions they wish to take. The bulk of the code for the opponent model consists of methods that update the counters from the information found in the messages sent from the server. The statistics that are used in the model are often context based and in order to update the counters one has to check numerous conditions for each.

B. Exploit specifications

Exploits are rule modules that adhere to the generic exploit template. The generic exploit template includes two methods: `appliesToStats` and `getAction`. `appliesToStats` is the same for

every exploit, it takes a stats model object and returns whether the given exploit applies to the stats model. An exploit can only apply if the prerequisite statistics have been observed more times than a given threshold, which is expert defined for each statistic in each exploit. The `getAction` method is given the game context and returns an action if the exploit applies to the specified context. For `getAction` to be called `appliesToStats` must have been called previously and must have returned true.

All in pre-flop exploit example

```

//If the opponent has been raised all in at least 11 times
if(model.whereRaisedAllinPre > 10){
  //if they have played at least 31 hands in position
  if(model.vpipInPosition > 30){
    if(model.vpipD > 30 && model.vpipD < 50){
      if(model.folded2AllinPre < 10)
        //lowest EV = + 3224 ~32BB
        rankThresholdD = 25;
      else if(model.folded2AllinPre < 30)
        //lowest EV = +3346 ~33BB
        rankThresholdD = 20;
    }else if ...
  }
}

```

The profitability of each exploit is determined by calculating the exploit’s expected value. In order to determine the profitability we must first calculate the equity which is the chance that a hand or range of hands (the possible hands the opponent could have) will win the pot. We used Poker Stove [10] to calculate the equities that were used in our expected value calculations. Exploits have hand rank thresholds to determine which action to take. The value of these thresholds depends on the opponent’s frequency statistics. The threshold determines the range of hands with which the agent will take the action. In the all in pre-flop example there are two thresholds for the hand rank, one for when the opponent is the dealer (`rankThresholdD`) and one for when they are not (`rankThresholdND`). This is because opponents often play vastly different hand ranges from different positions. The exploit’s hand range is determined by a hand ranking list and the value given by the various `rankThresholds`. Each hand in the list has a number associated with it, it is ordered from best to worst, AA being 1. The exploit’s hand range consists of all cards with a value less than or equal to the `rankThreshold` that applies. The opponent’s range is estimated, using the frequency statistics from the opponent model. The estimation always takes from the top of the opponent’s range ensuring that the opponent’s range is never stronger than the estimated one. Yet, it may be weaker, which works in our favor. SartreNL’s hand range, based on the thresholds and the opponent’s hand range, are used to calculate the equities when calculating the exploitation’s expected value. Expected Value (EV) is the long term expected outcome of a given hand or situation and is calculated using the following equation:

$$EV = [Our\ Equity] * [What\ we\ win] - [Opponent's\ equity] * [What\ we\ lose]$$

Through the use of this calculation we can ensure that our exploits are profitable. This guarantees that our use of an exploit is never exploitable by the opponent, given that our opponent model is accurate. To prove this we will walk through the all in pre-flop example, explaining how it is decided if the exploit applies or not and when the exploit is used once it applies.

The snippet from the all in pre-flop exploit comes from the all in pre-flop exploits appliesToStats method. First it is checked whether the opponent has been raised all in pre-flop at least 11 times. 11 is an expert defined threshold chosen on the assumption that, since going all in is such a big decision, players will not be doing it for deceptive purposes and the statistic will be accurate within a small number of trials. Next the exploit ensures the opponent has voluntarily put money into the pot enough times such that another expert defined threshold is met. If this is the case, it then determines the dealer rank threshold (rankThresholdD) based on the opponent's rate of voluntarily putting money into the pot in position (vpipD) and the opponent's rate of folding to all in raises pre-flop (fold2AllinPre). The value of vpipD tells the exploit the range of hands the opponent plays in position. The portion of hands the opponent folds to an all in raise pre-flop is removed from the bottom of the opponent's hand range.

A number of different thresholds have been created that take effect depending on the frequency statistics of the opponent. For each threshold the minimum expected value of going all in pre-flop was calculated by using the smallest range of hands for the opponent that still applies. For example the minimum EV for the rankThresholdD 25 was calculated as follows:

Our hand range was set to the top 25 hands. The opponent's hand range was set to the top 27% of hands, since he plays at least 30% of all hands and folds 10% of these to a pre-flop all in. Using PokerStove [10] our equity was 57.811% and our opponents was 42.189%, giving us the calculation below.

$$EV = 0.57811 * 20100 - 0.42189 * 19900 = 3224.4$$

The calculations for determining the non-dealer rank threshold are similar, but use voluntarily put money in pot out of position and update rankThresholdND. When creating the thresholds for the all in pre-flop exploit we made sure each had an expected value of at least 30 big blinds, since the risk when going all in is your entire stack.

The getAction method of the all in pre-flop exploit determines the position of the opponent. It then checks the hole card rank given in the context against the threshold that applies, if the hole cards are within the acceptable range

getAction returns the string "a" for "all in action" otherwise it returns null.

C. Opponent exploiter

The opponent exploiter is the module that provides the interaction between SartreNL, the statistical model, and the exploits. The opponent exploiter has a statistical model associated with it and has four lists of exploits: pre-flop, flop, turn and river. Each list represents all exploits that apply to the associated statistical model for each of the four betting rounds. These lists are populated through the findApplicableExploits method which considers each exploit, calling its appliesToStats method, and supplying the statistical model the opponent exploiter is associated with as the parameter. If an exploit returns true, it is added to the list of exploits for the betting round it applies to. The opponent exploiter also includes a method for each betting round that takes the game state information and calls the getAction method for each exploit in the list for the given betting round. If an exploit returns a non-null value, this is passed along to the agent that subsequently uses the action.

VII. EXPERIMENTAL RESULTS

A. Methodology

We required several opponents to challenge in order to evaluate the results of using our statistical exploitation module. Optimally we would evaluate statistical exploitation against a variety of competencies, ranging from easily exploitable to un-exploitable. The participants in the ACPC represent a good variety of computer players. While it is not possible to challenge the agents submitted to the competition directly, due to them not being publicly available, the hand history information is available for each agent that participated. Expert imitator case-bases were created for several of the no limit Texas Hold'em participants from the 2011 ACPC that imitate and generalize the opponent's style of play from their hand histories. These case-bases were created to be used by the expert imitation based framework described in [9], training each expert imitator on the decisions made in the competition by each of the chosen agents. Agents imitating an opponent through the use of the expert imitator always play a static strategy since the case-base used to determine their actions does not change. For these experiments this is fortunate because the statistical exploitation module is currently only able to handle static opponents. A strong poker agent will be able to deal with non-static strategies, and we will address this in section VIII. The agents chosen cover a variety of exploitability, ranging from *highly exploitable* to *difficult to exploit*. Table 1 shows the authors' views on the exploitability of the various agents.

Table 1. Exploitability of the opposition.

Player	Exploitability
POMPEIA	highly exploitable
Kappa	highly exploitable
Hugh	highly-moderately exploitable
Lucky7	moderately exploitable
Hyperborean-iro	difficult to exploit
Hyperborean-tbr	difficult to exploit

Six opponents were challenged against SartreNL without exploits and SartreNL with the addition of exploits (henceforth SartreNLExp). Each of the six agents played two seeded duplicate matches against both SartreNL and SartreNLExp. A duplicate match consists of 20,000 hands in total. 10,000 hands are initially played, the players then switch seats and the same 10,000 hands are played again. This way each of the players receives the cards that their opponents received before. The duplicate match style was used to decrease the variance that is normally involved in poker. To decrease the overall variance further, the same seed value was used for each of the duplicate matches played between each of the variants of SartreNL and the various opponents, making the matches identical in terms of the cards for each hand.

Overall 24 duplicate matches were played; both, SartreNL and SartreNLExp played two duplicate matches against each of the opponents. To determine the effectiveness of the addition of the partially opponent based strategy based on statistical exploits on SartreNL the duplicate matches were split into two subsets: Run 1 and Run 2. Run 1 consisted of the first duplicate match SartreNL played against each of the opponents and the first duplicate match SartreNLExp played against each opponent. Run 2 consisted of the second duplicate match that SartreNL and SartreNLExp ran against each of the opponents. The agents played each run without any knowledge of the opponent they were facing.

In each run the match SartreNL played against a particular opponent is used as the base-line. The difference in performance between SartreNL and SartreNLExp can then be

taken as the effect of the exploits. Some of the agents have some randomness associated with their strategies, so although the same hands and community cards came up in all matches this does not mean the agents chose the same action each time. Due to this it is likely that the scores fluctuated between matches. Furthermore an exploit was not applied to every hand and SartreNLExp chooses the actions as SartreNL normally would for hands where exploits are not applied to. Therefore, the situation could have occurred where exploitations were used and the overall score for SartreNLExp was lower than SartreNL's score. This would not have been due to the exploits losing money, since exploits always have a large positive expected value, it would have been caused by the randomness in SartreNL's action selection process. This means that although the exploits had positive effects, SartreNLExp chose less profitable actions in the hands in which exploits were not used, causing SartreNLExp's overall score to be lower than SartreNL's.

To combat these problems we only compared the resulting scores of the hands in which exploits were used. There exists a duplicate match between SartreNLExp and each opponent and a corresponding duplicate match between SartreNL and each opponent in each run. The scores for the exploited hands were found for both the SartreNL and SartreNLExp match by going through the log files and tallying up the result for each of the hands in which exploits were used. The difference between these two scores shows the impact of the exploits much more accurately than the difference between the overall scores.

B. Results

Table 2 presents the results against the set of chosen opponents from the 2011 ACPC competition. The opponent in the match is given as the column heading, which is further split into the two Runs. The table is split into two sections due to spatial limitations. The *SartreNL* row depicts the outcome of SartreNL's matches. The *SartreNLExp* row depicts the outcome of SartreNLExp's matches. The match outcomes are depicted in milli-big blinds per hand, and only take into account hands in which exploits were used by SartreNLExp. Milli-big blinds records the average number of big blinds won per hand, multiplied by 1000. The *number of hands* row shows

Table 2. Results

Bot Name	POMPEIA		Kappa		Hugh	
	Run 1	Run 2	Run 1	Run 2	Run 1	Run 2
Number of Hands	2090	2090	4160	4160	4703	5122
SartreNL	2,111	2,205	29,425	29,425	-1,258	-1,184
SartreNLExp	48,900	48,900	47,200	47,200	1,628	1,168
Difference	46,789	46,695	17,775	17,775	2,885	2,351
Bot Name	Lucky7		Hyperborean-iro		Hyperborean-tbr	
	Run 1	Run 2	Run 1	Run 2	Run 1	Run 2
Number of Hands	6392	6353	1265	1265	1357	1357
SartreNL	224	220	179	58	159	-37
SartreNLExp	1,008	503	972	895	855	885
Difference	784	283	793	837	696	921

the number of hands in the match in which exploits were used by SartreNLExp. The *difference* row depicts the difference in win rate in milli-big blinds per hand between SartreNLExp and SartreNL, clearly indicating the effect that the use of exploits had on the hands.

All of the results shown in Table 2 indicate a large improvement in performance when statistical exploits are used. The results suggest that the statistical exploits module is able to appropriately determine and exploit statistical anomalies found in the play of opponents. There are significant increases not only against the *highly exploitable* opponents, but also against the *difficult to exploit* opponents, which suggest that using an approximate Nash equilibrium strategy as the base strategy has the desired effect. SartreNLExp is able to remain difficult to exploit while concurrently exploiting its opposition.

VIII. CONCLUSION

In conclusion, we have presented an approach for exploiting statistical anomalies in the game of No Limit Texas Hold'em. Rather than create a fully opponent based agent we have created a module which can be added to a base strategy to create a partially opponent based agent. This allowed us to easily overcome one of the major problems faced by frequency statistic based opponent models, allowing us to play the base strategy as we wait for our frequency statistics to become accurate depictions of the opponent's strategy. We added the statistical exploitation module to an approximate Nash equilibrium base strategy to create an agent which was both difficult to exploit and able to exploit opponents. The module was able to safely exploit static opponents, i.e. it was able to exploit opponents without making itself any more exploitable than the base strategy. Our experimental results show that the use of statistical exploitations significantly increased the win rate of the agent in the hands in which exploits were used without increasing exploitability.

The statistical exploitation module's tradeoff comes not in the form of increasing the agent's exploitability to allow for the ability to exploit opponents, as in Polaris [1]. Instead we limit the type of opponent exploitation to only statistical exploitation, allowing the module to exploit opponents without increasing the exploitability of the base strategy. This limitation means the statistical exploitation module is not able to fully exploit an opponent like a best response would. However, ensuring the agent is as difficult for opponents to exploit as possible is more effective than allowing it the ability to maximally exploit some opponents.

IX. FUTURE WORK

There are several improvements that could be made to increase the performance of the current system and increase

the scalability. The most obvious improvement is the addition of more exploits. The more exploits the system includes, the better it is at exploiting opponents. So far the system does not have a large number of exploits. Further, the frequency statistics model currently does not use decaying history and is unable to remember opponents. Decaying history is an important upgrade as the system is vulnerable to any non-static strategies without it, as it cannot react quickly to strategy alterations after it has built up a view of its opponent.

In future work, the module will therefore be improved by the re-implementation of the statistical model to store hand histories. This will improve the performance and scalability of the statistical model, making it easy to include and compute further frequency statistics and allowing the ability to recognize previously played opponents. Decaying history is easier to implement for a system that stored hand histories than for the current system that has only counters. If only counters are available, it is difficult to determine how to alter the counter to implement decaying history. If hand histories were stored, it would be simple to determine what the last M hands were and how they affected each counter. Stored hand histories would also allow the module to have statistic-specific histories. In this way, state specific statistics will be based on the last k occurrences of the state that the statistic is associated with, instead of just having the last M hands played.

REFERENCES

- [1] M. B. Johanson, "Robust Strategies and Counter-Strategies: Building a Champion Level Computer Poker Player," M.S. thesis, Dept. CompSci., Univ. of Alberta, Edmonton, Canada, 2007.
- [2] S. Markovitch, R. Reger. "Learning and Exploiting Relative Weaknesses of Opponent Agents." *Autonomous Agents and Multi-Agent Systems*, vol. 10, pp. 103-130, 2005.
- [3] T. C. Schauenberg, "Opponent Modelling and Search in Poker," M.S. thesis, Dept. CompSci., Univ. of Alberta, Edmonton, Canada, 2006.
- [4] D. Billings, D. Papp, J. Schaeffer, D. Szafron "Opponent Modeling in Poker," *Proc. of the 15th Nat. Conf. of the American Association for Artificial Intelligence, Wsconsin*, 1998, pp 493-499 .
- [5] D. Billings, L. Pena, J. Schaeffer, D. Szafron, "Using Probabilistic Knowledge and Simulation to play Poker," *Proc. of the 16th Nat. Conf. of the American Association for Artificial Intelligence, Florida*, 1999, pp 697-703.
- [6] A. Davidson, "Opponent Modeling in Poker: Learning and Acting in a Hostile and Uncertain Environment," M.S. thesis, Dept. CompSci., Univ. of Alberta, Edmonton, Canada, 2002.
- [7] J. Rubin, I. Watson, 1st Initial. , "Opponent Type Adaptation for Case-Based Strategies in Adversarial Games," *ICCBR 2012*, Vol. , no. , 357-368, 2012.
- [8] J. Rubin, I. Watson., "Successful Performance via Decision Generalisation in No Limit Texas Hold'em," *Int. Conf. on Case-based Reasoning*, 2011, pp 467-481.
- [9] J. Rubin, I. Watson., "Similarity-based retrieval and solution re-use policies in the game of texas hold'em," *18th Int. Conf. on Case-Based Reasoning* , 2010, pp 465-479.
- [10] *PokerStove* [Online]. Available: <http://www.pokerstove.com/blog/> [Accessed: Oct 12, 2012]