

# On Combining Decisions from Multiple Expert Imitators for Performance

Jonathan Rubin and Ian Watson

Department of Computer Science

University of Auckland

Auckland, New Zealand

jrub001@aucklanduni.ac.nz

ian@cs.auckland.ac.nz

## Abstract

One approach for artificially intelligent agents wishing to maximise some performance metric in a given domain is to learn from a collection of training data that consists of actions or decisions made by some *expert*, in an attempt to *imitate* that expert's style. We refer to this type of agent as an *expert imitator*. In this paper we investigate whether performance can be improved by combining decisions from multiple expert imitators. In particular, we investigate two existing approaches for combining decisions. The first approach combines decisions by employing ensemble voting between multiple expert imitators. The second approach dynamically selects the best imitator to use at runtime given the performance of the imitators in the current environment. We investigate these approaches in the domain of computer poker. In particular, we create expert imitators for limit and no limit Texas Hold'em and determine whether their performance can be improved by combining their decisions using the two approaches listed above.

## 1 Introduction

Supervised machine learning algorithms typically construct models by processing a collection of training data. Novel instances/problems are later addressed by making generalised decisions based upon the model that was learned during training. One approach for artificially intelligent agents wishing to maximise some performance metric in a given domain is to learn from a collection of training data that consists of actions or decisions made by some expert, in an attempt to *imitate* that expert's style. The expert that provides the training data can be either human or artificial. We refer to this type of agent as an *expert imitator*. This approach has seen successful application in computer games [Ontañón *et al.*, 2009; Floyd and Esfandiari, 2009; Rubin and Watson, 2010], but has also been applied in other domains [Coates *et al.*, 2008].

Our research takes place in the domain of two-player Texas Hold'em poker – a game defined by simple rules, which however, provides a rich, dynamic environment for applying sophisticated strategies. Performance in the game of poker is known to have an intransitive relationship, i.e. while player

A beats player B and player B beats player C, it does not necessarily follow that A beats C. Given the intransitive nature of poker strategies, the successful combination of different styles of play has the potential to largely improve overall performance, against a range of different opponents.

By training on experts with different styles of play, we create multiple expert imitators in the domains of limit and no limit Texas Hold'em. We determine whether the performance of individual expert imitators can be improved by combining their decisions. We investigate two existing approaches for combining decisions.

The first approach combines decisions using ensemble voting [Dietterich, 2000], where a final decision is made by having each expert imitator vote on their preferred action. The second approach combines decisions by dynamically selecting a single expert imitator for each hand of play. This approach attempts to choose the expert imitator that achieves the greatest profit against a particular style of play in order to maximise overall profit against a range of opponents. This approach was originally applied to the game of limit Texas Hold'em in [Johanson *et al.*, 2007; Johanson, 2007] where the UCB1 [Auer *et al.*, 2002] allocation strategy was applied to dynamically select experts during play. Here we follow the same procedure described in [Johanson, 2007] for dynamically selecting *expert imitators*. Our own research extends this idea to the more complicated domain of no limit Texas Hold'em. Furthermore, our work focuses specifically on combining the decisions from multiple *expert imitators*. *Expert imitators* can be used to imitate the different playing styles of artificial or human experts, simply by observing and recording their betting decisions. A major benefit of using this approach is that producing new strategies simply requires updating the data that is used to train the system. By using expert imitators, a diverse set of strategies with various playing styles can be produced with negligible computational effort.

## 2 Related Work

While any learning algorithm that uses training data from a single expert to train a system can be considered an *expert imitator*, some learning algorithms are more suited to this particular purpose than others. For example, the *lazy learning* [Aha, 1997] of case-based reasoning is particularly well suited to expert imitation where expert observations can be

recorded and stored for later use at decision time. Depending on the particular domain, some CBR systems have demonstrated successful imitation after making only a handful of observations [Ontañón *et al.*, 2009]. In the games domain, [Floyd and Esfandiari, 2009] demonstrated successful imitation in the domain of simulated robot soccer to control a multi-agent soccer team. [Ontañón *et al.*, 2009] observes human experts playing a real time strategy game and uses case-based planning to generate unique strategies. Our own case-based poker playing system, *Sartre* [Rubin and Watson, 2010], has demonstrated strong performance at international computer poker competitions by training on hand history data from the strongest agent in the previous year’s competition. Using this approach a version of *Sartre* placed 3rd out of 15 competitors at the 2010 AAI competition<sup>1</sup>.

Many classification systems have made use of *ensemble* decision making to improve accuracy. Typically a collection of different learned hypotheses will vote for a particular category. The category with the most (possibly weighted) votes is the one that is selected. [Davidson, 2002] used an ensemble approach to predict future actions of an opponent in the domain of limit Texas Hold’em. The results showed that having multiple learning algorithms vote improved overall accuracy compared to just a single learning algorithm.

Rather than having different hypotheses vote on a solution, a more difficult problem involves dynamically selecting an appropriate hypothesis at runtime in order to maximise overall payoff. Johanson *et al.* [2007] created a collection of exploitive strategies in the game of limit Texas Hold’em that were designed to exploit particular opponents. They employed the UCB1 update policy together with showdown DIVAT [Billings and Kan, 2006] analysis to select the best exploitive player against a particular opponent during game play. In this work, we apply a similar idea to selecting *expert imitators* in the domain of limit and no limit Hold’em.

### 3 Expert Imitators

The expert imitators we describe in this work are lazy learners created using the case-based reasoning methodology [Kolodner, 1993]. Each expert imitator,  $I_j$ , is made up of a collection of cases,  $C_j = \{c_{j,1}, c_{j,2}, \dots, c_{j,n}\}$ , where  $j$  refers to the particular expert being imitated and  $\forall c \in C_j, c = (x, a)$ , where  $x$  is a feature vector that captures game state information and  $a$  is an action vector that specifies the probability of taking a certain action, given game state  $x$ .

Given a game state described by  $c_t$ , a decision is made by processing the collection of stored cases and maximising a global similarity metric to retrieve the most similar case to  $c_t$ .

$$c_{max} = \arg \max_{c_k} sim(c_k, c_t), \forall c_k \in C_j \quad (1)$$

where,  $sim(c_1, c_2)$  is some global similarity metric that determines the similarity between the feature vectors of the two cases  $c_1, c_2$ .

The retrieved case,  $c_{max}$ , suggests an action ( $a_{max}$ ) that  $I_j$  uses to approximate the decision that the original expert would have taken, given the current game state.

<sup>1</sup><http://www.computerpokercompetition.org/>

### 3.1 Action Vectors

The action vectors used to make a betting decision differ depending on whether the limit or no limit variation is being played.

#### Limit Hold’em

In the domain of limit Hold’em the action vector used by the expert imitator has the following format:

$$a = (f, c, r)$$

where each entry within the vector is a number between 0.0 and 1.0 that indicates the probability of taking a particular betting action. All entries within the vector must sum to exactly 1.0. The actions that are possible in limit Hold’em are either fold ( $f$ ), check/call ( $c$ ) or bet/raise ( $r$ ).

#### No Limit Hold’em

In limit Hold’em, if a player wishes to raise, the amount they are allowed to raise is set at a certain predefined level. On the other hand, in the no limit variation, players are allowed to bet or raise any amount, including all of the money or chips they possess (this is called going *all-in*). While this appears to be a simple rule change, the consequences it has on game play are quite dramatic. In order to construct an action vector for no limit Hold’em, a mapping is required to assign quantitative bet amounts into discrete categories. The no limit action vector used by the expert imitators in this paper is given by the following:

$$a = (f, c, q, h, i, p, d, v, t, a)$$

where each entry refers to different discrete betting categories as follows: fold ( $f$ ), call ( $c$ ), raise quarter pot ( $q$ ), raise half pot ( $h$ ), raise three quarter pot ( $i$ ), raise pot ( $p$ ), raise double pot ( $d$ ), raise five times pot ( $v$ ), raise ten times pot ( $t$ ), all-in ( $a$ ). Once again each entry corresponds to the probability of taking that particular action and all entries in the vector sum to 1.0.

## 4 Combining Decisions

Given an action vector, the final action an expert imitator selects can be determined by one of the following policies:

1. **Probabilistic** – Select an action probabilistically based on the values specified in the action vector. Or,
2. **Max Frequency** – Select the action within the triple that has the greatest probability value

### 4.1 Ensemble-based Decision Combination

Given an action vector  $a = (a_1, a_2, \dots, a_n)$  and a collection of expert imitators  $I_1, I_2, \dots, I_m$ , where each imitator applies the *max frequency* decision policy, we can derive a vector  $v = (v_1, v_2, \dots, v_n)$ , whose elements correspond to the total number of votes each action received.

The final action taken by the ensemble of imitators is given by selecting the action,  $a_j$ , that corresponds to the *strictly* maximum number of votes,  $v_j$ , if one exists. If a strictly maximum number of votes does not exist, then the action specified by imitator  $I_1$  is used, where  $I_1$  is an arbitrary imitator selected by the authors.

## 4.2 Dynamic strategy selection

### Multi-arm bandit problem

The multi-arm bandit problem is a well known problem in statistics and artificial intelligence. The problem consists of a gambler and a collection of slot machines (one-armed bandits). Each slot machine has a particular reward distribution associated with it, which is unknown to the gambler. The gambler seeks to maximise his or her overall profit by playing the right machines. The problem for the gambler is to determine which machine to play next based on the sequence of rewards he or she has witnessed so far.

The situation described above depicts a classic example of the **exploration/exploitation** trade off. The gambler needs to choose between selecting the machine that has so far offered the greatest reward (*exploitation*), versus selecting a different machine which may offer yet a greater reward than that witnessed so far (*exploration*).

### UCB1

The UCB1 [Auer *et al.*, 2002] algorithm offers a simple, computationally efficient solution to the multi-arm bandit problem. The UCB1 algorithm defines a policy for selecting the next machine to choose based on the concept of **regret**. In the multi-arm bandit problem, *regret* refers to the difference between the reward the gambler would have received had they always selected the optimal slot machine compared to the actual reward the gambler received. The UCB1 algorithm offers a logarithmic bound on regret and is given as follows:

1. Select each machine once
2. Select machine,  $j$ , that maximises the following equation:

$$\arg \max_j \bar{x}_j + \sqrt{\frac{2 \ln n}{n_j}} \quad (2)$$

where,

$\bar{x}_j$  is the average reward obtained from machine  $j$ ,  
 $n$  is the total number of trials so far, and  
 $n_j$  is the total number of times machine  $j$  has been selected.

The algorithm depicted above can be applied to dynamic strategy selection in the computer poker domain. Rather than selecting slot machines, the UCB1 algorithm can determine the order in which to select different expert imitators. In the computer poker domain,  $\bar{x}_j$  in equation (2) now refers to the average utility (profit or loss) achieved by expert imitator,  $I_j$ , while playing against the current opponent. Similarly,  $n_j$  refers to the total number of times  $I_j$  has been chosen to play against the current opponent. However, due to the inherent variance present in the game of Texas Hold'em, strategy selection based on profit alone can severely bias the results of the UCB1 allocation policy. Hence, efforts to reduce this inherent variance are required in order to stabilise the results.

### Variance Reduction

As in [Johanson, 2007], we employ the use of DIVAT analysis [Billings and Kan, 2006] in order to improve the average

utility values,  $\bar{x}_j$ , used by the UCB1 allocation strategy. DIVAT (ignorant value assessment tool) is a perfect information variance reduction tool developed by the University of Alberta Computer Poker Research Group<sup>2</sup>. The basic idea behind DIVAT is to evaluate a hand based on the expected value (EV) of a player's decisions, not the actual outcome of those decisions. DIVAT achieves this by comparing the EV of the player's decisions against the EV of some baseline strategy, see equation 3.

$$EV(\text{Actual Actions}) - EV(\text{Baseline Actions}) \quad (3)$$

Equation 3 attempts to factor out the effects of luck as both the actual strategy and the baseline strategy experience the same *lucky* or *unlucky* sequence of events. For example, an imitator that benefits by a statistically unlikely event no longer makes a difference as the baseline strategy also benefits by the same improbable event. What actually matters is any difference in EV the imitator is able to achieve by varying the actions that it takes. When the imitator is able to achieve a greater EV than the baseline, it is rewarded. Alternatively, if the imitator's actions result in a lower EV then it is punished with a negative outcome. If the EV of both strategies is the same the outcome is 0.

Any type of strategy can be used as a baseline. Typically, in limit Hold'em a bet-for-value baseline strategy is adopted. A bet-for-value strategy makes no attempt to disguise the strength of its hand. By establishing game-theoretic equilibrium-based thresholds, a sequence of baseline betting decisions can be constructed, based on hand strength alone. Strong hands will be bet and/or raised, whereas weak hands will be checked or folded.

The EV, from equation 3, is approximated by a concept known as *pot equity* or *all-in equity*<sup>3</sup>. To determine *all-in equity* in a two-player Hold'em match, both player's hidden cards need to be known. A player's all-in equity is given by multiplying the current amount in the pot by the probability that the player will win the hand. On the last round the probability that the player will win is either 0 or 1. For intermediate rounds the probability that the player will win is determined by enumerating all possible combinations of the remaining community cards and counting the number of times the player's hand wins, loses and ties. All-in equity is computationally efficient as it makes the basic assumption that a player's equity is given by the *current* pot size and does not attempt to predict what the final pot size might have been had there been future betting.

Finally, DIVAT requires perfect information about both players' hidden cards. When either player folds this information is not available. Therefore, to calculate the utility values used by the UCB1 allocation strategy, DIVAT analysis is employed only when a showdown occurs. If a fold occurs the actual monetary outcome of the hand is used instead.

<sup>2</sup><http://webdocs.cs.ualberta.ca/games/poker/index.html>

<sup>3</sup>A better EV calculation involves the concept of *roll-out equity*. Here we use all-in equity as it is less computationally expensive.

Table 1: Expert Imitator Results against Fell Omen 2 and AlistairBot

	Fell Omen 2		AlistairBot		Average	
<b>Dynamic</b>	<b>0.01342</b>	$\pm 0.006$	0.68848	$\pm 0.009$	<b>0.35095</b>	$\pm 0.0075$
<b>Ensemble</b>	0.00830	$\pm 0.010$	0.67348	$\pm 0.010$	0.34089	$\pm 0.0100$
<b>Rockhopper-max</b>	-0.01445	$\pm 0.009$	<b>0.69504</b>	$\pm 0.009$	0.34030	$\pm 0.0090$
<b>PULPO-max</b>	-0.00768	$\pm 0.006$	0.66053	$\pm 0.024$	0.32643	$\pm 0.0150$
<b>Sartre-max</b>	0.00825	$\pm 0.014$	0.63898	$\pm 0.019$	0.32362	$\pm 0.0165$
<b>PULPO-prob</b>	-0.00355	$\pm 0.011$	0.63385	$\pm 0.018$	0.31515	$\pm 0.0145$
<b>Sartre-prob</b>	-0.01816	$\pm 0.006$	0.64535	$\pm 0.015$	0.31360	$\pm 0.0105$
<b>Rockhopper-prob</b>	-0.02943	$\pm 0.011$	0.63870	$\pm 0.020$	0.30464	$\pm 0.0155$

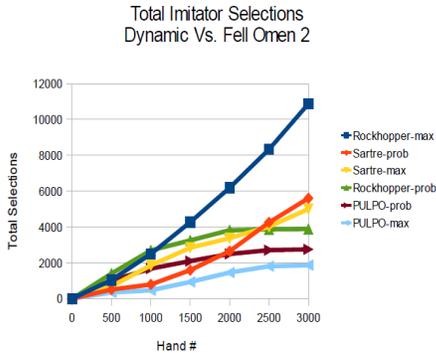


Figure 1: Shows the total number of times each expert imitator was selected to challenge Fell Omen 2

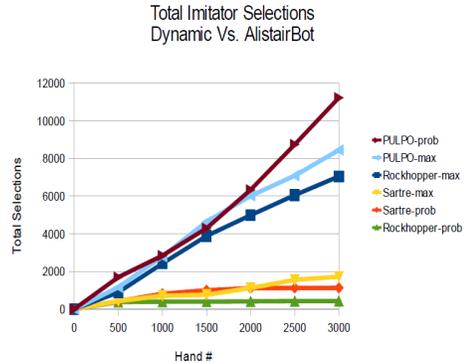


Figure 2: Shows the total number of times each expert imitator was selected to challenge AlistairBot

## 5 Experimental Results

### 5.1 Methodology

We provide experimental results in the domains of both limit and no limit Texas Hold'em. In each domain 3 basic expert imitators were constructed. All expert imitators were trained on hand history data from the annual computer poker competition. Each expert imitator was trained on separate data determined as follows:

1. The first expert imitator was trained on decisions made by the agent that won the bankroll division of the 2010 AAAI computer poker competition.
2. The second expert imitator was trained on decisions made by the agent that won the equilibrium division of the 2010 AAAI computer competition.
3. The final expert imitator was our own entry into the 2010 AAAI computer poker competition.

Both *max frequency* and *probabilistic* decision re-use policies were investigated resulting in a total of 6 basic expert imitators.

From these 6 expert imitators, we also constructed two players that combined decisions in the following ways:

1. An *ensemble-based* player was constructed where each of the basic *max frequency* expert imitators voted for a particular action and the action with the most votes was

selected. The decision of one designated expert imitator, selected by the authors, was used in case all expert imitators voted for a separate action.

2. A *dynamic-selection* based player was constructed that dynamically selected at runtime the best basic expert imitator to use against the current opponent. Selection was based on the UCB1 allocation strategy together with showdown DIVAT analysis. In the limit variation, a bet-for-value baseline strategy was adopted during DIVAT analysis. In no limit a more simplistic *always-call* strategy was used as the baseline.

These 8 players were then each challenged against two different types of computerised players. All matches played were duplicate matches. A single duplicate match involves playing 3000 hands, after which the players' memories are wiped and the 3000 hands are played again, but in the reverse direction, i.e. the cards that were initially given to player A are instead given to player B and vice-versa. This way both players get to play both sets of cards and this reduces the variance that is involved in simply playing a set of hands in one direction only. For each of the 16 match-ups, 5 duplicate matches were played, for an overall total of 480,000 hands.

### 5.2 Limit Results

In the limit domain each expert imitator challenged the following competitors:

Table 2: Expert Imitator Results (No Limit) against MCTSBot and SimpleBot

	MCTSBot		SimpleBot		Average	
<b>Hyperborean-max</b>	<b>1.7781</b>	$\pm 0.193$	<b>0.7935</b>	$\pm 0.084$	<b>1.2858</b>	$\pm 0.139$
<b>Dynamic</b>	1.3332	$\pm 0.146$	0.5928	$\pm 0.058$	0.9630	$\pm 0.102$
<b>Ensemble</b>	1.3138	$\pm 0.047$	0.5453	$\pm 0.058$	0.9295	$\pm 0.053$
<b>Hyperborean-prob</b>	1.1036	$\pm 0.165$	0.5075	$\pm 0.093$	0.8055	$\pm 0.129$
<b>Sartre-max</b>	0.9313	$\pm 0.117$	0.5248	$\pm 0.032$	0.7281	$\pm 0.075$
<b>Sartre-prob</b>	0.4524	$\pm 0.106$	0.3933	$\pm 0.073$	0.4228	$\pm 0.090$
<b>Tartanian-max</b>	0.1033	$\pm 0.451$	0.3450	$\pm 0.087$	0.2242	$\pm 0.269$
<b>Tartanian-prob</b>	-0.0518	$\pm 0.127$	0.4221	$\pm 0.032$	0.1852	$\pm 0.080$

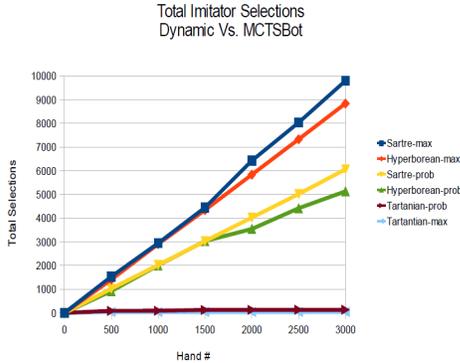


Figure 3: Shows the total number of times each no limit expert imitator was selected to challenge MCTSBot

1. **Fell Omen 2** – A solid Nash equilibrium-based agent commonly used as a benchmark for testing limit Hold'em agents.
2. **AlistairBot** – An exploitive agent that uses Monte-Carlo simulation to determine the decision with the best EV against the current opponent.

Table 1 presents the results of the 8 expert imitators against the above two competitors. The results are presented in small bets per hand (sb/h), where the total number of small bets won or lost are divided by the number of hands played. Each original expert imitator begins with the name of the original expert used to train the system followed by the decision policy used, i.e. either *max frequency* or *probabilistic*. The Sartre agent was our entry into the 2010 computer poker competition. Figures 1 and 2 show the total number of times each expert imitator was selected via the UCB1 allocation policy (used by Dynamic) to play against each of the competitors.

### 5.3 Limit Discussion

Overall, the results presented in Table 1 support the idea that combining decisions from multiple expert imitators can improve performance over a single expert imitator alone. Against Fell Omen 2, Dynamic does better than any other expert imitator. Against AlistairBot, the expert imitator trained on the Rockhopper agent, using a *max frequency* betting policy, fares the best. Combining the results against the two opponents and taking the average sees Dynamic achieve

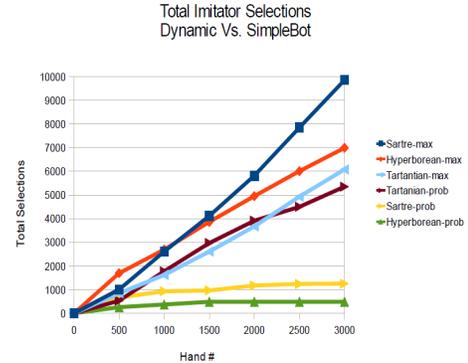


Figure 4: Shows the total number of times each no limit expert imitator was selected to challenge SimpleBot

the best average outcome, followed by Ensemble. However, care must be taken in interpreting these results as there is overlap between the standard deviations.

Also depicted in the results are plots showing the total number of times each original expert imitator was selected via the UCB1 allocation strategy to play against Fell Omen 2 (Figure 1) and AlistairBot (Figure 2). Table 1 suggests that the expert imitator Sartre-max did the best out of all the original imitators against Fell Omen 2, achieving a profit of +0.00825 sb/h. In particular, Sartre-max was the only original imitator to achieve a slight profit against Fell Omen 2, whereas the other 5 achieved slight losses. It was therefore thought by the authors that this agent would be selected the most by Dynamic when challenging Fell Omen 2, however this appears not to be the case. Instead, Figure 1 shows that Rockhopper-max was selected the most overall, followed by Sartre-prob and Sartre-max in third place. A similar scenario is depicted in Figure 2, where PULPO-prob is selected by far the most, whereas the obvious choice (Rockhopper-max, as it achieves the greatest profit of +0.69504 sb/h against AlistairBot) is only 3rd on the list. These discrepancies are likely due to a combination of variability in the data and possibly exaggerated utility values within the UCB1 allocation strategy.

### 5.4 No Limit Results

In the no limit domain, each expert imitator challenged the following opponents:

1. **MCTSBot** – an exploitive agent that uses Monte-Carlo Tree Search [Van den Broeck *et al.*, 2009].
2. **SimpleBot** – a no limit rule-based agent.

The *opentestbed*<sup>4</sup> project was used to gather results as the above no limit agents were made publicly available within this framework. Table 2 presents the no limit results. The results are in big blinds per hand. Once again Figures 3 and 4 show the total number of times an expert imitator was selected via the UCB1 allocation strategy to play against the above two competitors.

## 5.5 No Limit Discussion

While both Dynamic and Ensemble achieve high placings in the no limit results (2nd and 3rd, respectively). Table 2 suggests that they were not able to improve upon the performance of a single expert imitator, i.e. Hyperborean-max. Here Hyperborean-max performs the best against both MCTSBot and SimpleBot, suggesting that an alternative choice of opponent may be required in order to receive any benefit by combining decisions from multiple expert imitators. In both Figures 3 and 4, Sartre-max was selected the most by Dynamic, followed by Hyperborean-max.

The altered betting structure from limit to no limit had a few follow on effects when it came to dynamically selecting expert imitators. First, the amount of money invested by each player is typically larger than in limit poker. This can result in larger swings in the utility values used by the UCB1 allocation strategy. Furthermore, in no limit a lot more hands end when one player folds to another player’s bet, compared to limit where a lot more hands proceed all the way to a showdown. Showdown DIVAT (as its name suggests) can only be applied when a showdown is reached. The increased proportion of fold actions observed in no limit poker results in less variance reduction taking place on the utility values used by UCB1 and hence could have resulted in possibly selecting an inappropriate expert imitator a lot more often than was necessary.

## 6 Conclusion

In conclusion, we have evaluated a series of expert imitators in the domains of limit and no limit Texas Hold’em and investigated two separate approaches for combining their decisions in an attempt to improve performance. Overall, combining decisions either via ensemble voting or dynamic selection at runtime appears to improve performance. In the limit variation, combining decisions was able to produce better results than relying on one expert alone. This is likely due to the intransitive nature of poker strategies. In the no limit variation, dynamic selection did better than most single expert imitators alone, but was not the overall best strategy to use. This could be due to the fact that we simply have not tested against a diverse enough set of opponents and therefore, no benefit was received by selecting different expert imitators, or that further improvements to the variance reduction techniques used in the no limit domain are required to stabilise the results.

<sup>4</sup><http://code.google.com/p/opentestbed/>

## References

- [Aha, 1997] David W. Aha. Editorial. *Artificial Intelligence Review*, 11(1-5):7–10, 1997.
- [Auer *et al.*, 2002] Peter Auer, Nicolò Cesa-Bianchi, and Paul Fischer. Finite-time analysis of the multiarmed bandit problem. *Machine Learning*, 47(2-3):235–256, 2002.
- [Billings and Kan, 2006] Darse Billings and Morgan Kan. A tool for the direct assessment of poker decisions. *The International Association of Computer Games Journal*, 2006.
- [Coates *et al.*, 2008] Adam Coates, Pieter Abbeel, and Andrew Y. Ng. Learning for control from multiple demonstrations. In *Machine Learning, Proceedings of the Twenty-Fifth International Conference (ICML 2008)*, pages 144–151, 2008.
- [Davidson, 2002] Aaron Davidson. Opponent modeling in poker: Learning and acting in a hostile and uncertain environment. Master’s thesis, University of Alberta, 2002.
- [Dietterich, 2000] Thomas G. Dietterich. Ensemble methods in machine learning. In *Multiple Classifier Systems, First International Workshop, MCS 2000*, pages 1–15, 2000.
- [Floyd and Esfandiari, 2009] Michael W. Floyd and Babak Esfandiari. An active approach to automatic case generation. In *Case-Based Reasoning Research and Development, 8th International Conference on Case-Based Reasoning, ICCBR 2009*, pages 150–164, 2009.
- [Johanson *et al.*, 2007] Michael Johanson, Martin Zinkevich, and Michael H. Bowling. Computing robust counter-strategies. In *Advances in Neural Information Processing Systems 20, Proceedings of the Twenty-First Annual Conference on Neural Information Processing Systems*, 2007.
- [Johanson, 2007] Michael Bradley Johanson. Robust strategies and counter-strategies: Building a champion level computer poker player. Master’s thesis, University of Alberta, 2007.
- [Kolodner, 1993] Janet Kolodner. *Case-based reasoning*. Morgan Kaufmann Publishers, San Mateo, CA, 1993.
- [Ontañón *et al.*, 2009] Santiago Ontañón, Kane Bonnette, Prafulla Mahindrakar, Marco A. Gómez-Martín, Katie Long, Jainarayan Radhakrishnan, Rushabh Shah, and Ashwin Ram. Learning from human demonstrations for real-time case-based planning. In *IJCAI-09 Workshop on Learning Structural Knowledge From Observations*, 2009.
- [Rubin and Watson, 2010] Jonathan Rubin and Ian Watson. Similarity-based retrieval and solution re-use policies in the game of texas hold’em. In *Case-Based Reasoning, Research and Development, 18th International Conference on Case-Based Reasoning, ICCBR 2010*, pages 465–479, 2010.
- [Van den Broeck *et al.*, 2009] Guy Van den Broeck, Kurt Driessens, and Jan Ramon. Monte-Carlo tree search in poker using expected reward distributions. In *Advances in Machine Learning, First Asian Conference on Machine Learning, ACML 2009*, pages 367–381, 2009.