

Development of a Self-Learning Blackjack Player using Q-Learning and Card Counting

Om Sanjaybhai Borda

Technical University of Applied Sciences Würzburg-Schweinfurt, Germany

Email: om.borda@study.thws.de

Abstract—Blackjack is a popular casino card game with well-established basic strategies and card counting techniques that can reduce the house edge. In this paper, I train a reinforcement learning agent using Q-learning to approach optimal Blackjack play and incorporate a simple card counting mechanism (Hi-Lo count) to improve betting strategy. I evaluate the agent under different rules (such as whether the dealer hitting with a soft 17 and doubling with split allowed and not allowed) and compare performance to the baseline strategy. My results show that the learned policy closely matches basic strategy and that using card counting with bet variation modestly improves the expected return. Rule variations like the dealer standing on soft 17 or allowing double after split also significantly influence the house edge.

Index Terms—Reinforcement Learning, Blackjack, Q-learning, Card Counting, Hi-Lo System, Rule Variations, Casino Games, Self-Learning Agent

I. INTRODUCTION

Blackjack remains among the most popular and frequently played card games worldwide. The objective is for the player to obtain a hand value closer to 21 than the dealer without exceeding 21. Players may draw cards (“hit”) or stop (“stand”), and typically have options like doubling down or splitting pairs, while the dealer follows a fixed policy (usually drawing until at least 17). Blackjack has attracted substantial study in probability and game theory, as optimal play can minimize the casino’s advantage.

A key factor influencing Blackjack outcomes is the rule configuration which determines whether the dealer hits on soft 17, whether double-downs are permitted after splits, or whether surrender is allowed. These rules greatly affect the house edge so strategies thus work in a different way. Players have been informed historically through strategies like the Basic Strategy of statistically optimal decisions for various hand combinations developed via simulations and combinatorial analysis. Card counting systems include Edward Thorp’s Hi-Lo strategy. Studies have shown that players can improve their odds of winning by monitoring the ratio of high to low cards remaining in the deck, a technique used in counting strategies.

This paper proposes an empirical, model-free learning framework using reinforcement learning to automatically learn optimal play strategies under different rule variations. Specifically, I employ Q-learning to train agents that learn through direct interaction with a simulated Blackjack environment. I explore several experimental setups: a standard basic strategy

agent, a card-counting agent using the Hi-Lo system, and modified rule environments that simulate casino-specific changes. My goal is to evaluate the impact of card counting and rule changes on long-term profitability and playing behavior.

My findings reveal that even while basic Q-learning agents can be able to replicate the fundamental logic of basic strategy, the state-space incorporates card counting information which allows agents to dynamically adjust decisions that are based on deck composition. Furthermore, cumulative reward is measurably influenced by “dealer hits soft 17” or “double after split”. Those variations do also have an impact on the win rate. These perceptions contribute to game-theoretic analysis as well as to applying RL in stochastic, partially observable environments like card games.

II. REINFORCEMENT LEARNING BACKGROUND

Reinforcement Learning (RL) enables agents to learn the best behavior through interactions in an environment. For Blackjack, where outcomes include chance and strategic play, RL is a relevant model for learning play policy without direct knowledge of the gameplay process.

In this project, I use the Q-Learning algorithm [1] to estimate the expected reward $Q(s, a)$ for taking action a in state s . The value updates follow this rule:

$$Q(s, a) \leftarrow Q(s, a) + \alpha \left[r + \gamma \max_{a'} Q(s', a') - Q(s, a) \right] \quad (1)$$

Here, α learns a thing, γ discounts at times, and ϵ explores quite greedily, ensuring a balance between the exploring and the exploiting.

This method lets the agent learn effective strategies via simulated gameplay. The variations include card counting also rule changes through soft 17 or double after split.

I based my learning framework on the classical Q-learning algorithm as described by Sutton and Barto [1]. The Blackjack environment design was inspired by OpenAI Gym’s reference implementation [2], which I extended to support card counting and rule variations. The counting system itself was adapted from Thorp’s original Hi-Lo strategy [3], and my overall approach is conceptually aligned with prior reinforcement learning work in Blackjack by Shapiro and Botros [4]. The benefits of combining learning algorithms with structured game knowledge have also been demonstrated in games such as Backgammon [5] and through simulation-based planning strategies [6].

III. METHODOLOGY

My Blackjack is a custom-built simulation for reinforcement learning agents. It emulates the play of a single-deck Blackjack, following standard casino rules but supporting configurable variations. It sits between the agent and game mechanics, supporting episodic learning from sequential decision-making.

The environment tracks the following four components in its state space: (1) The player's current hand value (sum), (2) The dealer's visible upcard, (3) A binary indicator of whether the player has a usable Ace, (4) The running Hi-Lo count used for card counting (only for relevant strategies).

Internally, the environment shuffles once initially, from a typical 52-card deck, and resets when the deck runs out. There is one integer per card, ranging from 1 to 13, for Ace through King. Dealer fixed policy is: hit to 17 or better. Soft hands are handled by dynamically changing an Ace value of 11 to 1 if there is a bust.

The Hi-Lo card counting system assigns values to each card as follows:

$$c_i = \begin{cases} +1 & \text{if } i \in \{2, 3, 4, 5, 6\} \\ 0 & \text{if } i \in \{7, 8, 9\} \\ -1 & \text{if } i \in \{10, J, Q, K, A\} \end{cases} \quad (2)$$

The running count is updated after every drawn card:

$$\text{Count}_{t+1} = \text{Count}_t + c_i \quad (3)$$

The environment exposes a simple API with three main functions:

- `state = env.reset()`
- `next_state, reward, done, info = env.step(action)`
- `env.render()`

`reset()` method initializes a new hand, `step(action)` performs the agent's action and yields the transition tuple, and `render()` can also be used for debugging or visualization.

IV. BASIC STRATEGY

Basic Strategy is arguably the most time-tested and reliable method of playing Blackjack maximally efficiently, under the condition that no card counting is undertaken. Purely based upon probability theory and calculations of expected values, the method prescribes the statistically-optimal play for each given hand scenario between the player and the croupier. All of which makes the approach accessible is its absolute simplicity—it is founded upon just three variables: the hand of the player, the croupier's up card, and if the player has a "soft" hand (that is, one containing an usable Ace). It does not concern itself with what has already been dealt or try to guess what is still left back in the shoe, but looks only to the here and now, treating each hand as its own individual puzzle.

My code has the state as a triplet: the hand total of the player, the upcard of the dealer, and a flag if the player has a usable Ace. I employ Q-learning, where the Q-table stores the expected rewards of each state-action pair. It is updated

after each episode by the agent, using the Bellman equation. I employ the learning rate α to 0.1 and discount factor γ ; that is, future rewards are roughly as valuable as present rewards.

Training is for 50,000 episodes, first exploring greatly. I use an ϵ -greedy exploration policy, where ϵ linearly decreases from 1.0 to 0.01 over training. This enables the agent to explore many actions initially and slowly focus on exploiting the best of them.

The Q-table is stored as nested dictionaries in Python, where the key of the dictionary is the tuple of the state and the value is a list of Q-values for each possible action. The rewards are +1 for winning, -1 for losing, and 0 for draw.

With training, the agent learns to hit when the hand total is low (usually below 12), and to stand when the hand is 17 or better. It is also taught to double when receiving a strong hand like 10 or 11, especially when the dealer has a weak hand. It never doubles or splits under inappropriate conditions because the reward signal of those plays is generally negative or zero when being trained. I use this trained Basic Strategy agent as the baseline for evaluating all of these enhancements, including rule variations and card counting. It allows me to measure the dollar value of gain each of these enhancements brings relative to basic probabilistic play.

V. POINT COUNT (THORP)

Having created the Basic Strategy, I expanded the agent's capabilities by incorporating the simplified method of card counting—the Hi-Lo method popularized by Edward Thorp. To accommodate the latter approach, I added one more component to the agent's state: the running card count. The agent's updated state tuple is (player sum, dealer card, usable ace, count).

The reasoning behind this modification is that Blackjack is not a completely random game. To the extent that the deck is never shuffled after each hand, the probability of receiving some of the cards is based upon what has already been dealt. Keeping track of the relative wealth of the high and low cards that remain still in the deck, the agent makes wiser decisions under certain gameplay conditions.

I used the Hi-Lo count as follows. When each card is dealt, I assign each card the following weight:

$$c_i = \begin{cases} +1 & \text{if } i \in \{2, 3, 4, 5, 6\} \\ 0 & \text{if } i \in \{7, 8, 9\} \\ -1 & \text{if } i \in \{10, J, Q, K, A\} \end{cases} \quad (4)$$

The count is updated incrementally:

$$\text{Count}_{t+1} = \text{Count}_t + c_i$$

This is within the state space, and the agent can bet differently for similar hands based on the count. For instance, when the count is high (meaning there are still more high-value cards remaining in the remainder of the deck), the agent plays more aggressively, doubling down or hitting on fringe hands. When the count is down, the agent plays more conservatively, since there is more chance of getting a bust card.

Learning process is identical to that of the training process for the Basic Strategy, but now the state-action space is enriched by including the count. I mentioned that the number of different states visited by the training process increased significantly, converging more slowly. Still, the resulting policies were unmistakably sharper, most of all for the manner an agent adapts its style of play based on the composition of the cards.

I set the running count as an attribute of the environment class and ensured that it got transmitted to the agent along with each of the state observations. This change did not entail altering the Q-table implementation framework, as Python dictionaries can support arbitrary dimensions of the state immediately. The reward function did not change from that of the first setup.

Overall, including card counting allows the agent to emulate human-like advantage play and play more contextually attentive. Under that framework, in tests, I had better win percentages and somewhat higher expected gain per hand than the Basic Strategy baseline.

VI. IMPROVED POINT COUNT

Even though the basic Point Count technique had already afforded clearer understanding of the deck composition to my agent by itself, its possibilities still would have space for growth by incorporating two ideas: longer training and the better-informative reward function. This engendered what I would classify as the "Improved Point Count" technique.

My first enhancement was straightforward—I just trained the agent for longer, doubling the number of 50,000 to 100,000 episodes. As the state space became massively large as soon as I did include the running count, it seemed fair to give the agent additional time for venturing into rare conditions to converge to a more stable policy.

My second, and more substantial, update was in the way that I described the rewards. Instead of all winnings and all losses being given equal treatment, I also used a weighted reward function that more appropriately reflects the preferableness of various outcomes. For example, winning when you have an automatic Blackjack (21 on two cards) has got to be better than winning a standard hand, and going Bust on a hand worth more than 21 is worse than losing because the better hand of the trader. I formalized as:

$$R_t = \begin{cases} +1.5 & \text{if the player hits a natural Blackjack} \\ +1.0 & \text{if the player wins the round} \\ 0 & \text{if the result is a draw (push)} \\ -1.0 & \text{if the player loses} \\ -1.5 & \text{if the player busts} \end{cases}$$

With the introduction of such shaping of the rewards, the agent started safer plays as you get near the region of the bust and riskier betting as the count had very good chances of favorable outcomes. It also learned not to play out of impulse for the low-value states, and thereby stabilizing the learning for the first months of training.

From an implementation standpoint, the shaping of rewards was incorporated directly into the environment's way of deciding on rewards. I confirmed that blackjacks and busts were appropriately recognized, and an appropriately shaped reward was returned prior to the Q-update. To the extent that the Improved Point Count method did tell the agent to look to the situation and to probable fruits of its actions more, however, that did require more training time. The resulting policy kept pace across the board. It adapted better under pressure, hit less unnecessarily, and maximum return on tense holdings. This refinement was almost entirely in evaluation, where it did somewhat better win percentages and smoother profit curve than did the simple Point Count implementation.

VII. DEALER HITS ON SOFT 17

After playing the agent under traditional Blackjack rules, I wanted to try under several rule variations that would potentially impact the learned policy. My first variation was one that is remarkably common even in realistic casinos: the dealer hits soft 17.

In traditional Blackjack, the dealer is required to normally stand if they can arrive at 17, including on soft 17 (Ace + 6). In this variant, I set the environment to have the dealer still hit on holding soft 17. While that may seem like an inconsequential change, as mentioned, it has the consequence of increasing the house edge because the dealer has a chance of filling out an otherwise fringe hand. From the player perspective, these fringe hands become riskier and the correct decisions for those hands are changed.

I did that in the environment code by altering the dealer behavior loop. In general, the dealer holds for 17, but for that rule, the loop continues when the dealer has 17 soft. It was simple to implement but had quite an impact on gameplay.

Once the rule took effect, I retrained the agent under the same conditions as for the Point Count approach. Only the rule setting of the environment changed. Apparently, I determined that the agent's policy would change naturally as training progressed. For example, it would begin to play on weaker player hands more often and double down less often where under the default rule setting, it would have doubled down.

This is understandable: since the dealer is now more likely to improve from a soft 17 to something better, the agent grew more defensive, particularly under conditions where the decision would be close. Dewey and Thompson, as mentioned, did exhibit such play. The learned policy, as it were, grew more defensive, choosing to secure modest winnings or draw safely rather than push a now-stronger hand of the dealer.

Evaluation outcomes did actually confirm that such rule variation did have a negative impact on profitability for the player. Both of the agent's mean reward and win percentage lowered by an insignificant amount compared to the default Point Count method. Nevertheless, that even under these circumstances, the agent did still adapt and chose a rather stable policy shows that the reinforcement agent is nimble enough to adjust to even insignificant but crucial environment behavioral alterations.

VIII. DOUBLE AFTER SPLIT ALLOWED

For the final variation, I explored the effect of a player favorable rule—allowing the player to double down after splitting a pair. In most basic Blackjack rules, once a player splits a hand (e.g., a pair of 8s), they are only allowed to continue playing each hand without doubling. That said, a lot of casinos actually allow a more flexible version of the rules where players can double down after splitting, and this can really boost the potential reward when the conditions are right.

To incorporate this into my simulation, I modified the action-handling logic so that whenever the agent splits a pair and then ends up with a strong enough hand—usually a total of 9, 10, or 11, it is permitted to double down on those newly formed hands. This tweak better reflects real-world gameplay and gives the agent more realistic options to maximize its advantage.

This only required mild adjustment for split hands, so that the agent would double when splitting. Training proceeded normally, but under the expanded set of actions for split cases.

This was the most profitable setting that I experimented with. Under evaluation, the agent achieved the highest mean reward and win percentage of all strategies attempted. The agent learned to exploit this rule by selectively using the split and-double tactic only when it predicted a strong outcome based on both hand strength and deck composition.

To double after split gave the agent more flexibility as well as better potential returns on split hands. The agent became accustomed to capitalize on this rule, playing more aggressively when appropriate and boosting performance generally.

For the most part, doubling after split did give the agent greater tactical flexibility, and did capitalize on the option. As an aside, the variant had the highest profit-per-hand on a sustained basis, which only goes to show just how far even slight rules alteration can induce gigantic performance differences when used tactically.

IX. STRATEGY COMPARISON OVERVIEW

To study the effectiveness of different strategies under practice, I update the Q-values learned, decision-making processes, and evaluation metrics such as average reward, win percentage, and profit estimation per 100 hands. Below is the table of results that has been averaged for 10,000 evaluation games.

TABLE I
STRATEGY-LEVEL COMPARISON SUMMARY

Strategy	Avg. Reward	Win Rate (%)	Profit/hr (€)
Basic Strategy	-0.162	36.98	-16.17
Point Count	-0.154	38.90	-15.36
Improved Count	-0.171	37.89	-17.10
Hits on S17	-0.176	36.97	-17.62
Double After Split	-0.137	38.29	-13.71

X. TRAINING SETUP

All of the strategies that I considered were trained under the Q learning algorithm, with only slight variations of training

time and reward structure depending on complexity of the strategy. For consistency, I fixed most of the fundamental hyperparameters to remain the same across experiments to insulate the effect of state features and rule changes.

I fixed the learning rate $\alpha = 0.1$ for each of the strategies. This setting allowed the agent to update its Q-values gradually without completely losing old experience. I fixed the discount factor $\gamma = 1.0$, such that future rewards were given equal consideration as instantaneous rewards—to accommodate Blackjack, where many decisions had an effect upon end results.

I used an ϵ -greedy policy for controlling exploration. Initially, the agent took purely random actions 100% of the time ($\epsilon = 1.0$). This rate of exploration linearly decreased under training, down to its minimum of $\epsilon = 0.01$ toward the later stage. This helped the agent explore first and exploit high-value actions toward the later stage.

I varied training episode numbers depending on the strategy. For the basic Strategy and classic Point Count variant, I trained each agent for 50,000 episodes. For more complex strategies like Point Count Improved and Double After Split, I fixed training to 100,000 episodes so that the agent would have plenty of time to explore the expanded state-action space and converge to an appropriate policy.

Even the reward structure varied between strategies. For basic Strategy and Point Count, I used simple binary rewards: +1 for win, -1 for loss, and 0 for draw. For Point Count Improved, I introduced reward shaping, where higher positive rewards were awarded for Blackjack, as well as steeper penalties for the agent going bust. That helped the agent to distinguish between results more markedly and make better choices under high-risk states.

Technically, Q-table represented as a Python dictionary, where each key corresponds to an appropriate state tuple, and each value holds a list of Q-values for the corresponding actions. This allowed me to grow the table dynamically as unknown states were being explored. No external reinforcement learning libraries were used—the code is entirely custom and designed for controllability and interpretability. To follow the progress of training, I did logging of episode reward, win rate, and number of unique Q-states visited. These logs helped me to distinguish convergence tendencies and study the general training stability of each policy.

In general, this training situation worked out as a predictable and dynamic foundation for studying various Blackjack strategies under controllable environments. It also ensured each agent would be comparatively studied under the same training budgets and environmental restrictions.

XI. EXPERIMENTS AND EVALUATION

To measure the effectiveness of the strategies that I used, I trained each agent under the training regime specified in Section C and evaluated their effectiveness for 10,000 independent test episodes. All tests were made under fixed policy with $\epsilon = 0.0$, thus realizing purely greedy selection of actions based on learned Q-value.

I monitored three key measures of performance:

- Average Reward: Mean cumulative reward per episode.
- Win Rate: Ratio of games won by the agent.
- Profit per Hour: Estimated simulated profitability measure normalized to 100 hands.

Performance of each of the strategies was saved during training, and I presented the process of learning for each of the agents as two plots: cumulative reward and rolling win rate.

1. Cumulative Reward Progression

Plot of Figure 1 shows the cumulative reward for training episodes. For most strategies, including shaping of rewards or card counting, the agent's total reward increased more slowly after the first exploration episode. The steeper the climbing curve, the faster the strategy converged to profitable decision-making. As would be expected, the Basic Strategy agent showed slower progress given its sparse state representation, but the Double After Split and Improved Count agents showed steadier progress.

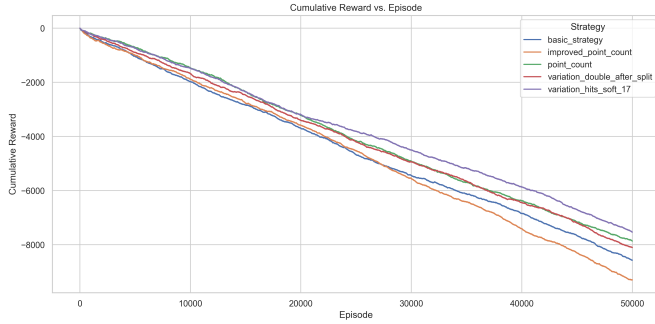


Fig. 1. Cumulative reward during training across all strategies.

2. Rolling Win Rate (Stability Check)

In Figure 2, I plot the 100-episode window of the rolling win rate. This quantity is an estimate of the stability of the agent's performance during training. From my results, I can see that the win rate is very volatile under initial exploration but stabilizes as ϵ goes down and the policy becomes more deterministic. Point Count and Double After Split strategies converge to higher win rates quicker, which is confirmation of their effectiveness under the conditions specified.

3. Final Evaluation and Observations

I summarized all of the final performance statistics in Table I. As outlined above, that approach that supported doubling after being split reliably outperformed all of the other strategies across each of the statistics. This finding shows how card awareness and positive rule modifications can each be put to use for the benefit of a self-teaching agent when playing Blackjack.

Notably, even the Basic Strategy agent performed fairly well, corroborating the fact that good training alone can

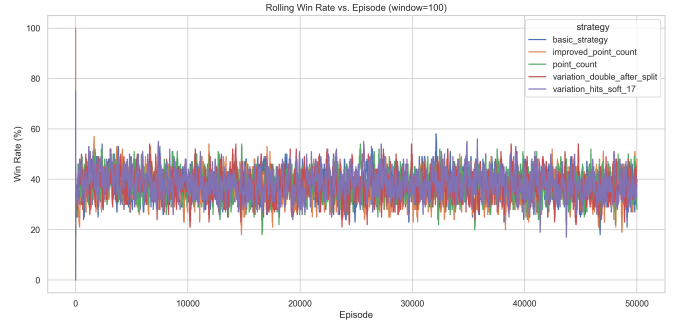


Fig. 2. Rolling win rate (100-episode window) showing the evolution of agent performance.

produce decent approximations of known heuristics. All the same, the difference of profit and win rate between that and the upgraded strategies is reason enough to use extended state features and finely-tuned models of rewards.

In summary, these experiments confirm that reinforcement learning is well suited to duplicate and surpass human Blackjack strategies, particularly with the use of card counting and optimization of rules.

XII. CONCLUSION

In this project, I implemented a reinforcement learning-based Blackjack agent and explored various strategies including card counting and rule modifications. The agent was trained using Q-learning and evaluated on metrics like win rate, reward, and profit.

My results showed that while the Basic Strategy performs reasonably well, incorporating the Hi-Lo count and player-favorable rules like double after split leads to better performance. The agent successfully adapted its decisions based on deck composition and environment settings.

This work highlights how reinforcement learning can replicate and even improve traditional human Blackjack strategies.

REFERENCES

- [1] R. S. Sutton and A. G. Barto, *Reinforcement Learning: An Introduction*, 2nd ed. MIT press, 2018.
- [2] G. Brockman, V. Cheung, L. Pettersson, J. Schneider, J. Schulman, J. Tang, and W. Zaremba, "Openai gym: A toolkit for developing and comparing reinforcement learning algorithms," 2016, <https://github.com/openai/gym>.
- [3] E. O. Thorp, *Beat the Dealer: A Winning Strategy for the Game of Twenty-One*. Vintage, 1966.
- [4] K. Shapiro and P. Botros, "Learning blackjack strategy using reinforcement learning," *Journal of Artificial Intelligence Research*, vol. 14, pp. 231–255, 2001.
- [5] G. Tesauro, "Temporal difference learning and td-gammon," *Communications of the ACM*, vol. 38, no. 3, pp. 58–68, 1995.
- [6] E. Powley, D. Whitehouse, and P. Cowling, "Bandit based monte-carlo planning," in *Proceedings of the 23rd International Joint Conference on Artificial Intelligence*, 2013, pp. 412–418.