

## Goals

We want to maximize the long term monetary gain to produce a stable source of return.

## Deliverables

Produce a simulation on never before tested data and produce a graph of the returns long term (multiple years). It should produce a strong return (minimum 5% target 10%, with a stability comparable to the stock market).

## Formalization

$g_t$  is the money in the pot,  $g_0$  is the initial investment,  $r_t$  is the percent return between  $g_t$  and  $g_{t+1}$  such that  $g_{t+1} = r_t g_t$ .

$$\max \lim_{t \rightarrow \infty} g_t$$

We can then define the amount we invest in each team, where we have teams  $A$  and  $B$ . An investment in team  $A$  would be represented as  $A$  and a investment in no team is  $\emptyset$ . Thus  $A_t + B_t + \emptyset_t = g_t$ .

Let's assume we have a model that provides us with both a probability of team  $A$  winning,  $P[W_t = A]$ , and its confidence on that probability,  $\epsilon_t$ .

Below is the Hoeffding bound

$$P[|v - \mu| > \epsilon] \leq 2e^{-2\epsilon^2 N}$$

There is also a given gambling odds  $O_A$  and  $O_B$  such that  $O_A = 1 - O_B$ .

$$g_{t+1} = \emptyset_t +$$

## Implementation

### Architecture

1. **Probability estimation:** Generates  $P[W_t = A]$  and  $\epsilon$ . Should create a game environment similar to a sport where I can measure and re-run to manually calculate  $P[W_t = A]$  and calculate the correctness of  $\epsilon$  to test is the model is constructed properly prior to using real data
2. **Gambler model:** Uses  $P[W_t = A]$  and  $\epsilon$  generated by the probability estimation model to find the optimal gambling strategy. Note that this model should be re-trained for each sport that is played as the sampling distribution will vary from sport to sport

## Milestones

1. Train gambler reinforcement model with sim that gives  $P[W_t = A]$  and  $\epsilon$  see if we can profit maximize in this environment, this way we don't need any data yet its just a game agnostic model that can be applied in any gambling situation
2. Create a sports environment with a similar setup to baseball and provides a similar data
3. Write and test model until epsilon is correct within a certain bound to be defined
4. Create data API for use in training the probability estimation model
5. Train the model
6. Use end to end and test
7. Write front end and ship

## Concerns to consider

There will be a tradeoff between achieving a good return and losing a large amount of pot. How accurate are betting odds in terms of expected value.

## Updates

After running the gambler model that just takes in  $g$ ,  $O$ , team\_a, team\_b and date, it returned an average earnings with an episode length of 300 of  $\mu = 9.94404e-1$  and  $\sigma = 6.81531e-3$ . This was run with the TD3 algorithm with  $\lambda_{loss} = -1.1$  and  $\lambda_{gain} = 1$ .

PPO with  $\lambda_{loss} = -1$  had final result of  $\mu = 2.19557e-13$  and  $std = 2.15727e-12$ , very poor. Ran A2C with the same  $\lambda$  performed the same.

Changing the  $\lambda$  value to -1.3 to see if an increase will help

None of the above worked. This is because the discount factor  $\gamma = 0.99$  which throws off the model completely as it thinks that its current actions affect the future rewards. When I changed episode length to 1 the final results were  $\mu = 1.68395$   $\sigma = 1.13996$ . I have not tried testing large length episodes so I am unsure if this actually produces stable results. However, after simulating it with just random variables I get that it sometimes returns very large amounts but most of the time it returns 0. Out of 50 trials of episode length of 50 only 2 returned any money (it was in the billions though). Question: can I divide both  $\mu$  and  $\sigma$  by the same amount to bolster the stability?

I will now try removing any possibility of cheating by removing all duplicates. Additionally, I will save the model so I can do future testing.