

Note set 4½

①

## Tic Tac Toe Reinforcement Learning

- Tic Tac Toe is the simplest game I could think of for us to learn how to make the computer learn by playing against it self.
- Several notable programs have learned how to play at near champion/champion level using similar techniques to what we will learn in this section. Some are

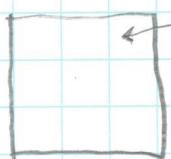
> TD gammon - Backgammon by Tesmo

> Alpha go

> etc.

Concerning Tic Tac Toe -

There are lots of states. A rough and incorrect number is  $3^9 = 19,683$



each square

↑  
ouch

can be empty,

X or O a → 9 squares  
so  $3 \rightarrow$  states

①

- ① In practice, there seems to be about 5,000 or so unique states, but that is still a large number — a couple of orders of magnitude more than the:
  - 6 State example
  - The cheese and pit game
  - ghost chase game
- ② Another problem here is that it is not immediately clear how to set up the
  - Matrix - action at current state leads to the next state
  - Why? For each state in the game, especially the early states, there are many actions available, leading to many next state. Each successive move, reduces this problem
  - Immediate - the goal state is loaded with immediate reward. In tic tac toe there are 8 ways to win the game, with many paths to get there resulting in a

(2)

multitude of goal states

$q$ -Table - a one to one match up with  $q$ Matrix that holds the delayed reward values associated with the moves at each state. Since we are unsure about defining  $q$ Matrix, these issues will impact how we set up  $q$ Table.

We will retain the functionality of  $q$ Matrix,  $q$ Table, and  $q$  Immediate by using a state dictionary

What is a state dictionary?

Each team, X and O will record each game.

After game is over the dictionary will be updated.

What is a dictionary?

In Python it is a data type that pairs a string with a value  
 $\text{myDictionary} = \{\}$

(3)

myDictionary["biff"] = 100

my Dictionary

{'biff': 100}

myDictionary["biff"]

100

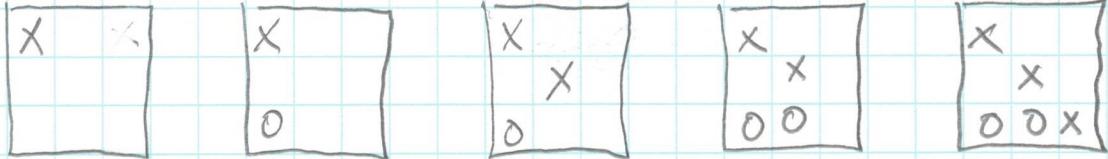
biff in myDictionary

True

- So our plan will be to put each game board into the dictionary as a string, with the value associated with being the reward value
- About reward - if a player wins a game it gets positive reward - for the game each state leading to the win gets delayed reward based on the discount factor
  - if a player loses, they get negative reward, similar to above
  - A cat game results in some negative reward for each player

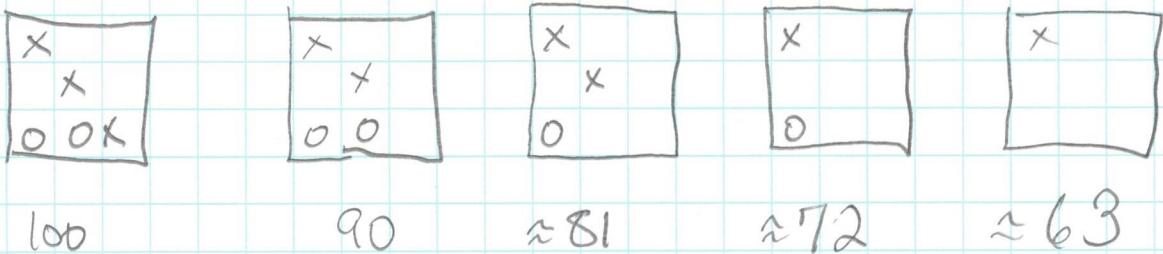
# Example

(4)



X wins!

> we will reverse the game and figure out the reward, discount = .9



\* no learn rate used here

Update is based on the standard of learning delayed reward equation

$$q\text{Table}[\text{state}, \text{action}] = q\text{Table}[\text{state}, \text{action}] + \ln * \left( \begin{array}{l} \text{immediate} \\ \text{reward} + \text{discount} * \\ q\text{Table}[\text{nextState}, \text{best action}] \\ - q\text{Table}[\text{state}, \text{action}] \end{array} \right)$$

- For tic tac toe only the last state (win state) gets immediate reward
- the  $q\text{Table}[\text{nextState}, \text{best action}]$  is equivalent to

we never see the game sequence, we are updating based on a game sequence that has already occurred.

game sequence at this point has already been determined - we end up with:

$$q\text{Table}[\text{state}, \text{action}] = q\text{Table}[\text{state}, \text{action}] + \ln * \left( \text{discount} * \text{CumReward} - q\text{Table}[\text{state}, \text{action}] \right)$$

# Pseudo Code

$x_{\text{states}} = \{\}$

$o_{\text{states}} = \{\}$

$x_{\text{states}}, o_{\text{states}} = \text{num TTTexVsOh (ngames,}$

15,000  
works well  
↓

$x_{\text{states}},$   
 $o_{\text{states}}$

$\text{num TTTexVsOh (ngames, } x_{\text{states}}, o_{\text{states}} \text{)}$

for ( $j = 0; j < \text{ngames}; j++$ ) {

    winner, game = tttExVsOh ( $x_{\text{states}},$   
 $o_{\text{states}}$ )

    if (winner is X) {

        updateDict ( $x_{\text{states}}, \text{game},$

            reward = 100, ln, discount)

        updateDict ( $o_{\text{states}}, \text{game},$

            reward = -100, ln, discount)

    else if (winner is O) {

        -- -- -

    }

    else if (cat game)

        updateDict ( $x_{\text{states}} -- \text{reward} = -20$   
            -- )

        updateDict ( $o_{\text{states}} -- \text{reward} = -20$   
            -- )

(6)

updateDict (states, game, reward, lr, discount)  
 $nGame = \text{reversed}(game)$

# For each board in the game, see if it is  
# in the dictionary

# if not, put it in, otherwise update its  
# reward state.

curReward = reward

for board in nGame :

bString = str(board)

# if the board is not in dictionary  
# put it in, initial reward is 0

if (bString in states) == False

states[bString] = 0

# update reward

states[bString] = states[bString] +  
 $lr * ((discount * curReward) - states[bString])$

curReward = states[bString]

# the next board back will get  
# curReward as its starting place.

(17)

TTTExVsOh (x states, 0 states)

board = initBoard()

turn = X

play = True

game = []

while (play) {

game.append(board)

moves = getMoves(board)

chance = getRand()

if (exploit(chance) == True)

if (turn is X's)

            find board + move is  
            x states with highest  
            reward

&gt; that is the best move

else if (turn is O's) {

            ...  
            ...  
            ...    else {  
        3  
        select a random move  
    }

board = board + move

wim, winner = checkForWin(board)

    if (wim)  
        play = False

return winner, game

8

How to do this.

- Write a simple person against computer tic tac toe first
  - computer picks randomly from available states
  - win / loss / cat are detected
- Next implement  
    + + + 2xVs Oh  
    then  
    num + + + 2xVs Oh and update Dict  
    check for win, get moves, etc  
    will be in simple tic tac toe
- Implement + + + ComputerVs Person  
    In this one computer uses dictionary to play person