# CA425
# Artificial Intelligence

Expressing Naughts & Crosses as a
Reinforcement Learning problem

Name: Ciaran McNally

Student ID: 58079978

# Game Introduction

Naughts-and Crosses has Nine positions in total, lined up in a 3x3 format. There are eight win states to this game.  A win occurs when either Player lines up 3 of their own symbols in a row. The winning positions are 3 going across the state-space, 3 going down the state-space and then the 2 diagonals. Each of the two players marks their symbol in turn. The Symbols commonly used are 'X' and 'O'.

Before I began diving in and implementing a Q-learning solution to Naughts-and-Crosses, I decided to investigate if there were any inherent advantages or disadvantages in the game itself. I also thought this would be a useful way of testing that my code works. I made two bots play against each other 10 Million times. These Naive Players simply picked a random free position in alternating goes. The data I generated can be seen below.

| Runs | Player 1 wins | Player 2 Wins | Draw |
|---|---|---|---|
| 1 | 576362 | 296251 | 127387 |
| 2 | 575465 | 297807 | 126728 |
| 3 | 575888 | 297817 | 126295 |
| 4 | 575864 | 296962 | 127174 |
| 5 | 576341 | 296781 | 126878 |
| 6 | 576302 | 296406 | 127292 |
| 7 | 576433 | 296532 | 127035 |
| 8 | 575421 | 297397 | 127182 |
| 9 | 576189 | 297245 | 126566 |
| 10 | 575956 | 296560 | 127484 |
| Total | 5760221 | 2969758 | 1270021 |
| % of 10 million | 57.60221 | 29.69758 | 12.70021 |

This data would suggest that there is quite a significant advantage to the Player who moves first. I figured this out by setting it so Player 1 always moved first. This means I must set my game to choose the player who moves first at random to ensure my results in further subsequent stages aren't off. I may also look into how this advantage could affect the learning rate of my Q-learning player at later stages.

# Game Code - Game.java

```java
1.  import java.util.Random;
2.  import java.util.Arrays;
3.
4.  //This class represents the Noughts and Crosses Game Environment.
5.  public class Game{
6.      //Stores game board
7.      private int[] state = new int[9];
8.      public int turn = 1;
9.
10.     public void newBoard(){
11.         //reset state of Game Board
12.         for(int i=0;i<state.length;i++){
13.             state[i] = 0;
14.         }
15.         //randomize who takes first turn
16.         Random r = new Random();
17.         boolean whoFirst = r.nextBoolean();
18.         if(whoFirst)    turn = 1;
19.         else    turn = 2;
20.     }
21.
22.     public String getState(){
23.         String s = Arrays.toString(state);
24.         String theState = s.replaceAll("[^\\d]","");
25.         return theState;
26.     }
27.
28.     //Print out board - Debugging purposes...
29.     public void printState(String state){
30.         char[] sp = state.toCharArray();
31.         System.out.println(sp[0]+" "+sp[1]+" "+sp[2]);
32.         System.out.println(sp[3]+" "+sp[4]+" "+sp[5]);
33.         System.out.println(sp[6]+" "+sp[7]+" "+sp[8]+"\n");
```

```java
34.    }
35.
36.    //return next state string
37.    public String getNextState(int who, int action){
38.        int[] nextState = state;
39.        nextState[action] = who;
40.        String s = Arrays.toString(nextState);
41.        String theState = s.replaceAll("[^\\d]","");
42.        return theState;
43.    }
44.
45.    public void makeMove(int who, int move){
46.        state[move] = who;
47.        if(who == 1)    turn=2;
48.        else    turn=1;
49.    }
50.
51.    public int whoWins(){
52.        if((state[0]==state[1] & state[1]==state[2]) & (state[0]!=0))
53.            return state[0]; //top row
54.        else if((state[3]==state[4] & state[4]==state[5]) & (state[3]!=0))
55.            return state[3]; //middle row
56.        else if((state[6]==state[7] & state[7]==state[8]) & (state[6]!=0))
57.            return state[6]; //bottom row
58.        else if((state[0]==state[3] & state[3]==state[6]) & (state[0]!=0))
59.            return state[0]; //left side down
60.        else if((state[1]==state[4] & state[4]==state[7]) & (state[1]!=0))
61.            return state[1]; //middle down
62.        else if((state[2]==state[5] & state[5]==state[8]) & (state[2]!=0))
63.            return state[2]; //right side down
64.        else if((state[0]==state[4] & state[4]==state[8]) & (state[0]!=0))
65.            return state[0]; //left to right diag down
66.        else if((state[2]==state[4] & state[4]==state[6]) & (state[2]!=0))
67.            return state[2]; //right to left diag down
```

```
68.        else if(isDraw())
69.            return 3;
70.        else return 0;
71.    }

72.

73.    public boolean isDraw(){
74.        int count = 0;
75.        for(int i=0; i < state.length; i++){
76.            if(state[i]==0)
77.                count++;
78.        }
79.        if(count > 0)   return false;
80.        else     return true;
81.    }

82.

83.    public boolean isValidAction(int who, int move){
84.        if((state[move]==0) & (who == turn))
85.            return true;
86.        else return false;
87.    }
88. }
```

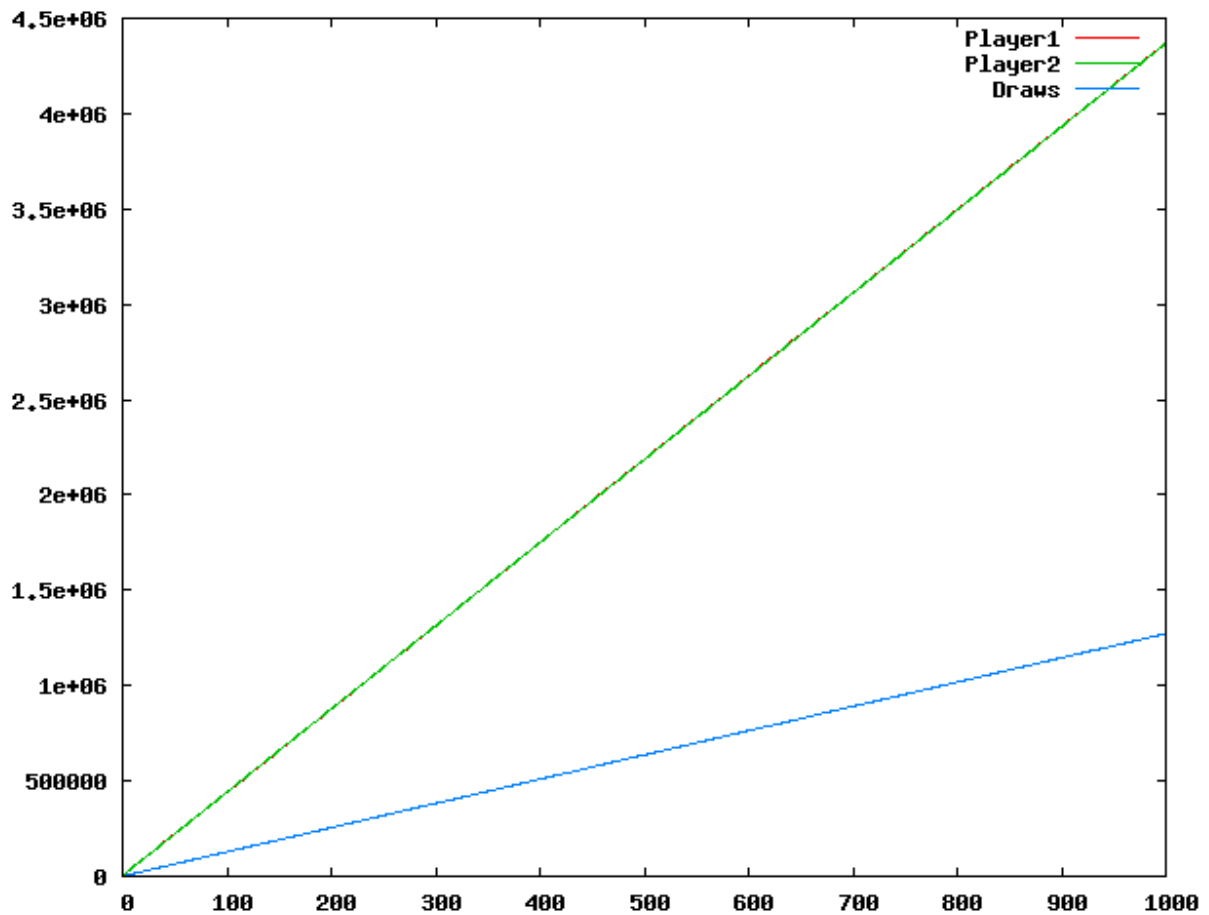# Random Player code - RPlayer.java

```
1.  import java.util.Random;

2.

3.  public class RPlayer{
4.      int currentMove;
5.      String currentState;
6.      public void nextAction(String state){
7.          Random r = new Random();
8.          int random = r.nextInt(9);
9.          currentMove = random;
10.         currentState = state;
11.     }
12. }
```

# Random Player vs Random Player

In the first round of tests I tried a random Player vs a random Player and received very expected results. Below you will see that both players win approximately an equal number of games. I ran this test run 10 million times taking a snapshot of the win ratio every 10,000 runs and then plotted the result. Below (Image 1.0) you will notice that both players line up as expected.
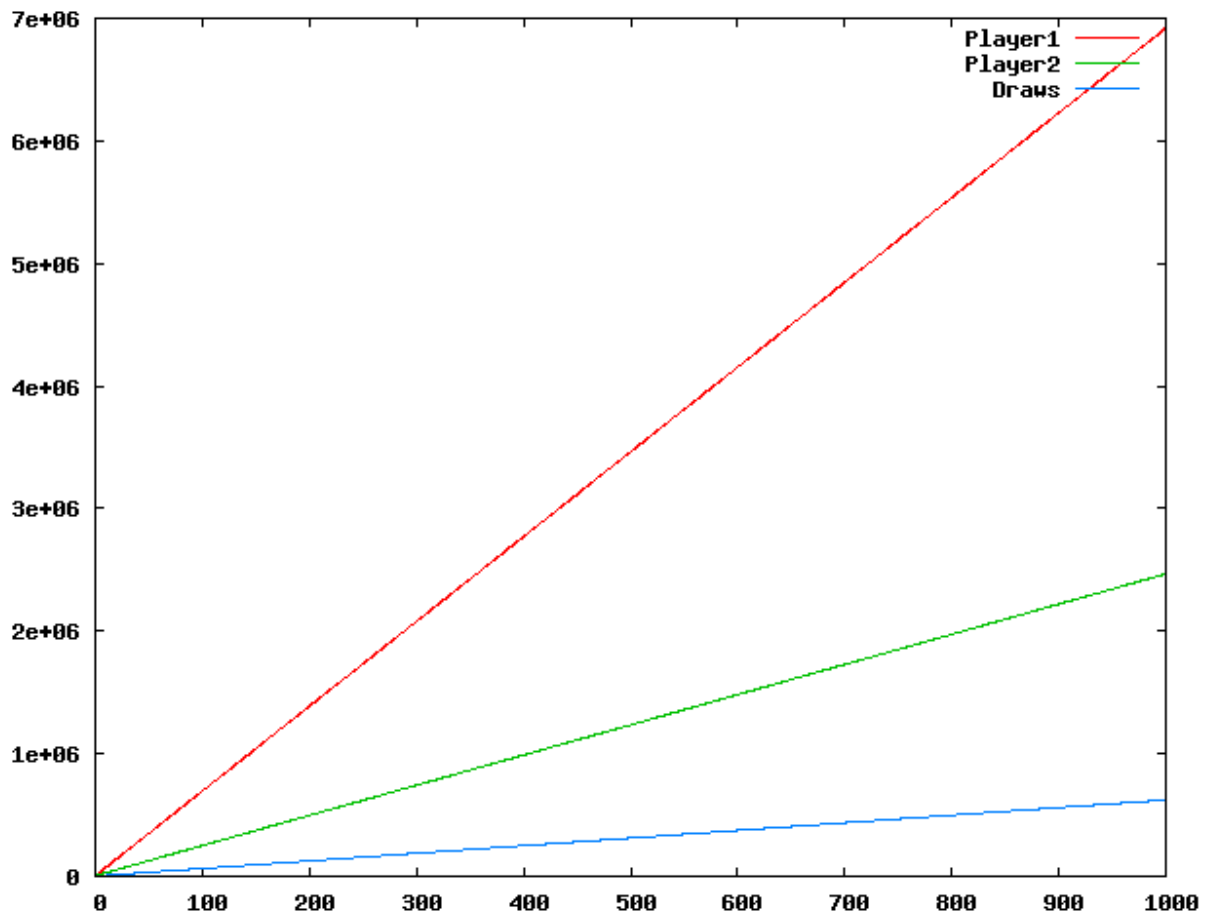
**Image 1.0**

# Q-Learning Player vs Random Player

For the first test run of my regular Q-Learning Player I used a value of 0.8 for gamma. In this implementation I didn't run the reward function after player 2's moves. I ran this 10 million times taking the win ratio on 10,000 run intervals. I then plotted this data using Gnuplot. The results can be seen below in Image 1.1.
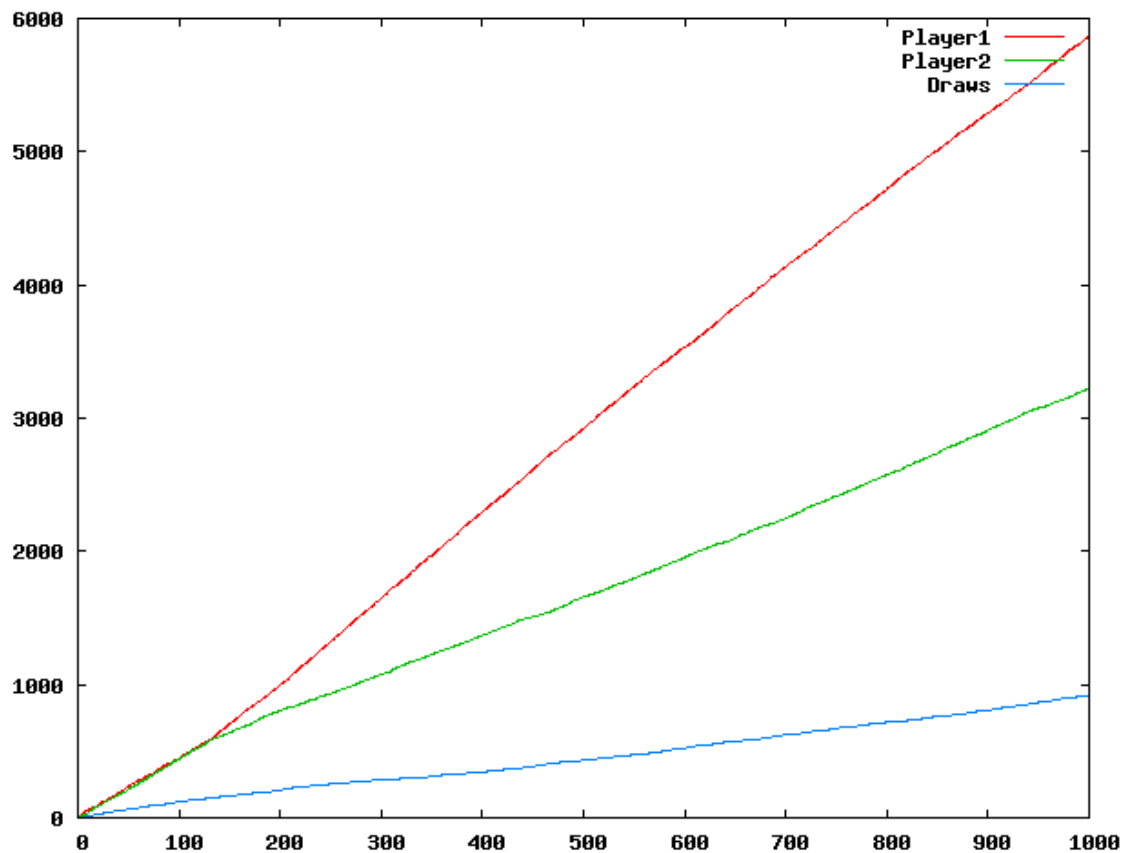
**Image 1.1**



We can see that there is a consistent rate of improvement to our Q-learning bot. The Q-learning Player is player 1, our random Player is player 2. If we change our gamma value I noticed there isn't a huge effect on the outcome of the results. It just takes a little bit longer for the bot to start taking better moves.

Below (Image 1.2) you will find a microscopic snapshot of the first 10000 runs. I took a snapshot every 100 runs and plotted this data. I would have expected to find more of a curve but there is a point where the Q-player has enough data to start winning more and more frequently and we can see the results of Player 1 and Player 2 begin to separate.

**Image 1.2**



You can notice how at almost 2000 runs our bots are taking the equivalent of mostly random actions, then when our bot has enough data it starts to greatly improve.
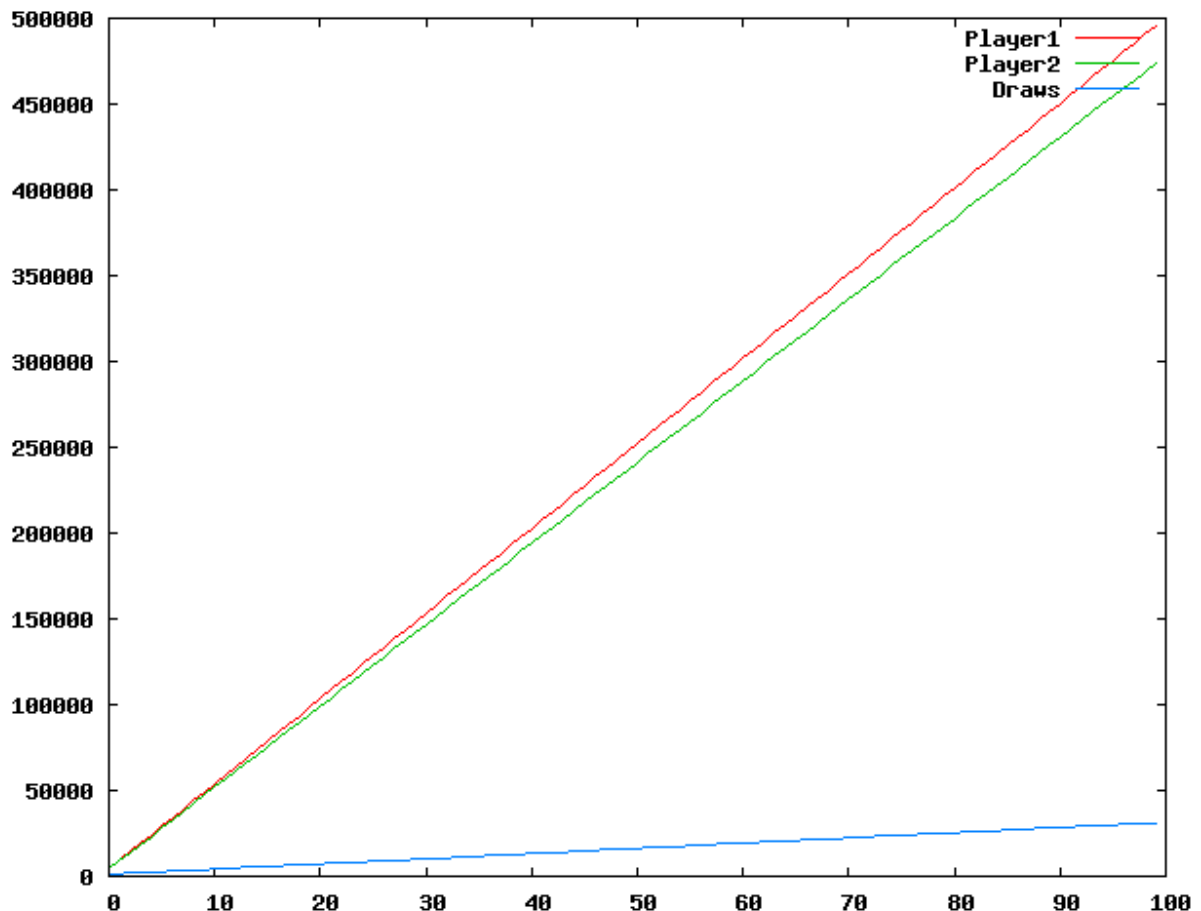
# Q-Learning Player vs Q-Learning Player

The data I generated for this was unexpected, I assumed both bots would learn to play against each other and this would eventually result in a lot more draws but this wasn't the case. This could possibly be an inconsistency in my implementation. Below (Image 1.3) you will find the outcome of my testing. I ran this test multiple times and what I found interesting is that either Player1 or 2 could take a slight lead.

You can see that after a point Player 1 has a slight lead that is then maintained. The same thing happened a few times with Player 2 on different test runs. It seems that once the Q-learning player takes a lead it maintains it, possibly it may have learned something that gives it an advantage over the other Q-Player.
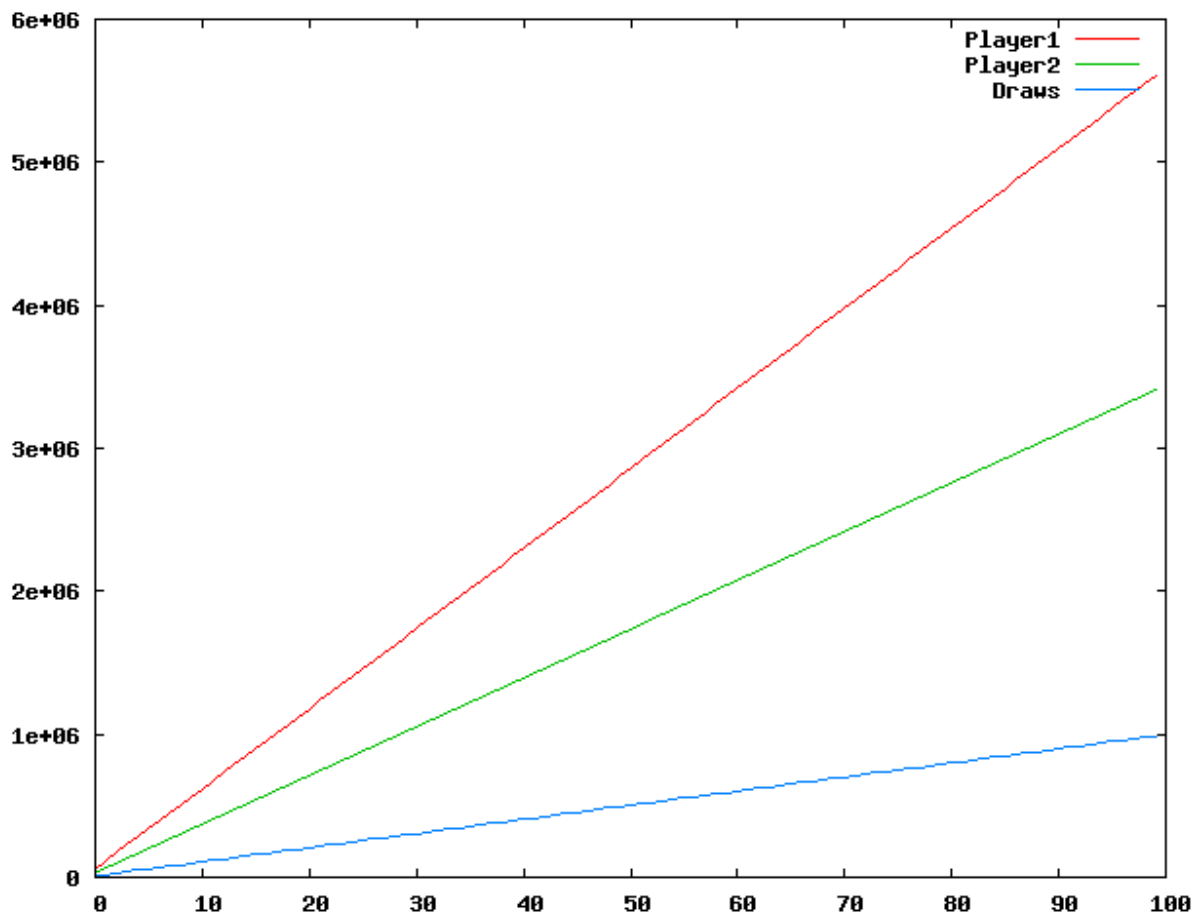
**Image 1.3**



As I was getting some unexpected results I decided to modify my code and now my Q-learning algorithm updates the reward for every state, as opposed to only observing states based on the Players own moves. This resulted in me getting some different results so I reinvestigated.

# Round2: Q-Learning Player vs Random Player

This time around my Q-learning algorithm updates itself on every state. This has resulted in different results. I modified my code so that I could look at different elements a bit clearer. I ran the same tests as before and did indeed get different data. I also added the ability to dump the Q-values in a json format. I would also like to implement a way of importing this data.

Below in the image 2.0 you will notice the Q-learning agent doesn't perform as well as in my recent testing. It does however demonstrate that the agent does beat the Random bot by a significant amount. This data below was 10 million runs, having a snapshot of the win ratio taken every 100,000 runs. We can see similarities in our old data when comparing the beginning stages. This Agent used a gamma value of 8.0 again.
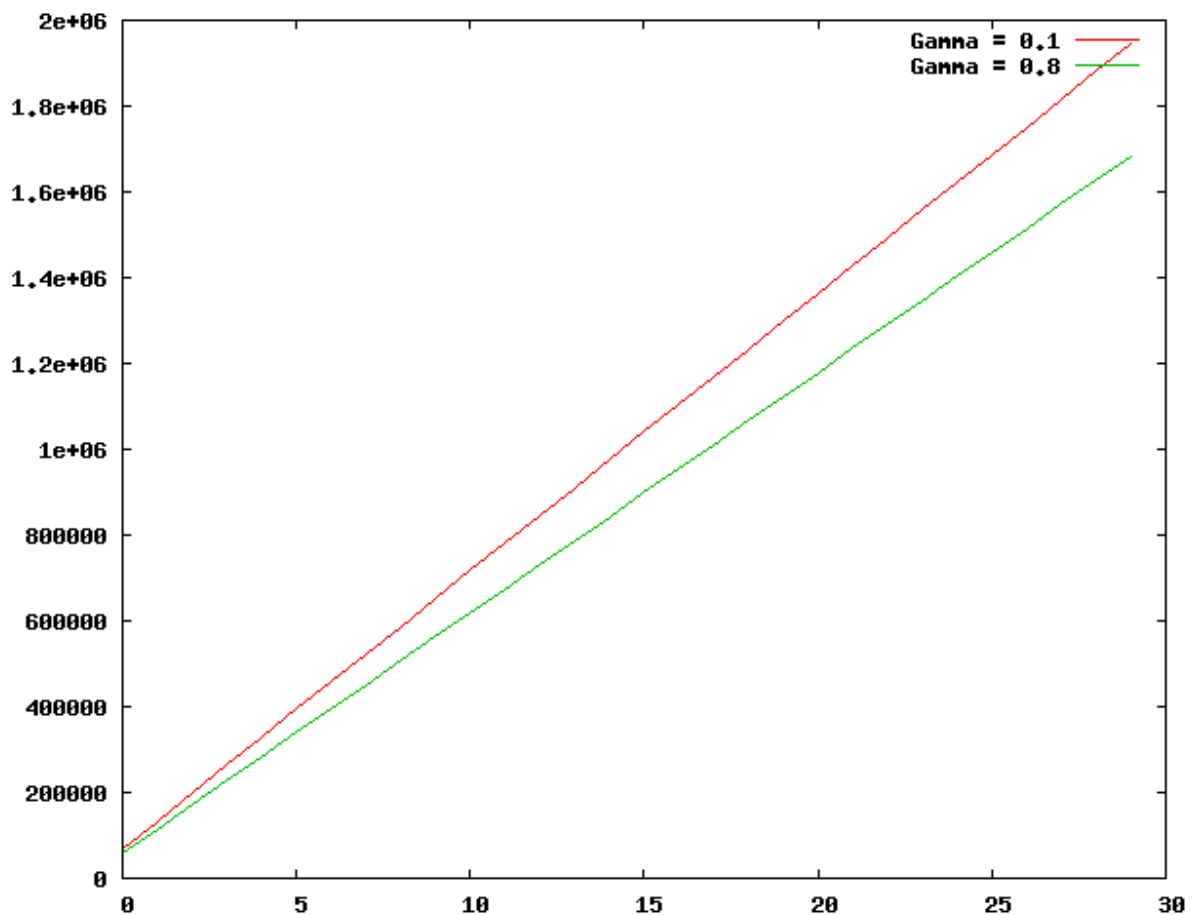
**Image 2.0**

Next I want to investigate what happens if I modify the discount factor, gamma. I ran the same tests as above but lowered gamma to 0.1. This had a very noticeable impact even in the early stages and our Q-learning agent performed much better. A lower discount award should make our agent aim for more short term targets, whereas a higher value will make it aim for a higher long term award. The results can be seen quite clearly below.

In Image 2.1 we can see that our short term reward Q-learning agent performs better from the offset. This image is a smaller snapshot of an overall trend. This was generated using 3 million runs with each different gamma setting, then I took a snapshot of the win ratio every 100,000 runs on interval.

**Image 2.1**



On the next page in image 2.2 you can see the overall results for our testing with a discount factor of 0.1. We can clearly see that this outperforms the results demonstrated in Image 2.0.

**Image 2.2**



I have successfully allowed the learning agent to store all of the Q-values in a json format. I can then continue from anywhere with these values and load them in on request. My only reason for not implementing the Running average element to this assignment was time constraints. It also might have been a little tricky with the way I have my code designed.

# Q-Learning Data / Mind Data code - Q.java

```java
import java.util.*;
import org.json.simple.*;
import java.io.*;
public class Q{
    /*
     * This Stores state as a key, in String format
     * then double Q values for each action in state.
     */
    HashMap<String, double[]> qMap = new HashMap<String, double[]>();
    String mindImport = "";

    public void importMind(String mind){
        mindImport = mind;
        boolean isData = true;
        String mindData = "";
        try{
            mindData = strFromFile(mindImport);
        }
        catch(NoSuchElementException e){
            System.out.println("Nothing in file... no data imported!");
            isData = false;
        }
        if(isData){
            reuseMind(mindData);
        }
    }
    /*
     * Lookup a state and find the best action to take
     * if more than single action is the same
     * pick random from bunch
     */
    public int getBestAction(String state){
        double[] actions = new double[9];
```

```java
34.        if(qMap.containsKey(state)){
35.            actions = qMap.get(state); //return all action values
36.        }
37.        double highestQ = -5000;
38.        ArrayList<Integer> next = new ArrayList<Integer>(); //possible choices
39.        int choice = 0;
40.        for(int i=0; i < actions.length; i++){
41.            if(actions[i] >= highestQ){
42.                highestQ = actions[i];
43.                choice = i;
44.            }
45.        }
46.        for(int i=0; i < actions.length; i++){
47.            if(actions[i] == actions[choice]){
48.                next.add(i);
49.            }
50.        }
51.
52.        //return a random choice if more than one
53.        //or single choice if only one
54.        if(next.size()==1){
55.            return next.get(0);
56.        }
57.        else if(next.size() > 1){
58.            Random r = new Random();
59.            int random = r.nextInt(next.size());
60.            return next.get(random);
61.        }
62.        else{
63.            //System.out.println(next.size());
64.            Random r = new Random();
65.            int random = r.nextInt(9);
66.            return random;
67.        }
```

```java
68.    }
69.
70.    /*
71.     * Get Q value for particular state/action
72.     */
73.    public double getQ(String state, int action){
74.        double[] actions = new double[9];
75.        if(qMap.containsKey(state)){
76.            actions = qMap.get(state);
77.        }
78.        return actions[action];
79.    }
80.
81.
82.    /*
83.     * Update Q value for particular state/action with reward
84.     */
85.    public void update_Q(String state, int action, double reward){
86.        double[] actions = new double[9];
87.        if(qMap.containsKey(state)){
88.            actions = qMap.get(state);
89.        }
90.        actions[action] = reward;
91.        qMap.put(state, actions);
92.    }
93.
94.    public void update_Q_reuse(String state, double[] actionsData){
95.        qMap.put(state, actionsData);
96.    }
97.
98.    /*
99.        This populates our Hashmap with the Q values from a file
100.       */
101.       public void reuseMind(String data){
```

```java
102.          Object obj2 = JSONValue.parse(data);
103.          JSONObject arr = (JSONObject) obj2;
104.          parseJson(arr);
105.      }
106.
107.      /*
108.       * Parse our Json object to pull array and object info
109.       */
110.      public void parseJson(JSONObject arr){
111.          Set<Object> set = arr.keySet();
112.          Iterator<Object> iterator = set.iterator();
113.          double[] thisk = new double[9];
114.          String state = "";
115.          while (iterator.hasNext()){
116.              Object obj = iterator.next();
117.              if (arr.get(obj) instanceof JSONArray){
118.                  //State and actions into Q
119.                  state = obj.toString();
120.                  thisk = getArray(arr.get(obj));
121.                  update_Q_reuse(state, thisk);
122.              }
123.              else{
124.                  if (arr.get(obj) instanceof JSONObject) {
125.                      parseJson((JSONObject) arr.get(obj));
126.                  }
127.                  else{
128.                      System.out.println(obj.toString());
129.                  }
130.              }
131.          }
132.      }
133.
134.      /*
135.       * Parse a JSONArray into a double array
```

```java
136.     */
137.    public double[] getArray(Object obj){
138.     JSONArray jsonArr = (JSONArray) obj;
139.     double[] vals = new double[9];
140.     for (int i=0; i < jsonArr.size(); i++) {
141.         vals[i] = (double)jsonArr.get(i);
142.     }
143.    return vals;
144.
145.    //Print out board - Debugging purposes...
146.    public void printState(String state){
147.        char[] sp = state.toCharArray();
148.        System.out.println(sp[0]+" "+sp[1]+" "+sp[2]);
149.        System.out.println(sp[3]+" "+sp[4]+" "+sp[5]);
150.        System.out.println(sp[6]+" "+sp[7]+" "+sp[8]+"\n");
151.    }
152.
153.    public void dump(){
154.        Set<String> enumk = qMap.keySet();
155.        Iterator<String> iter = enumk.iterator();
156.        while(iter.hasNext()) {
157.            String key = iter.next();
158.            double[] val = qMap.get(key);
159.            System.out.print(key+":  ");
160.            for(int i=0; i<val.length; i++){
161.                System.out.print("a"+i+":"+val[i]+",");
162.            }
163.            System.out.println();
164.        }
165.    }
166.
167.    public void dumpJSON(){
168.        JSONObject obj = new JSONObject();
169.        JSONArray arr = new JSONArray();
```

```java
            double[] val = new double[9];
            Set<Map.Entry<String, double[]>> set = qMap.entrySet();
            String key = "";
            for (Map.Entry<String, double[]> me : set) {
                key = me.getKey();
                val = me.getValue();
                for(int i=0; i<val.length; i++){
                    arr.add(val[i]);
                }
                obj.put(key, arr);
                arr = new JSONArray();
            }

            String jsonString = JSONValue.toJSONString(obj);
            strToFile(mindImport, jsonString);
    }

    //write string data to a file
     public void strToFile(String filename, String data){
        File file = new File(filename);
        try{
            FileWriter dataf = new FileWriter(file);
            dataf.write(data);
            dataf.close();
        }
        catch(FileNotFoundException e){
            System.out.println("\nFile not found!");
        }
        catch(IOException e){
            System.out.println("\nIO Error!");
        }
    }

    //Reads string from a file
```

```
204.        public String strFromFile(String filename){
205.            String mydata = "";
206.            File file = new File(filename);
207.            try{
208.                Scanner data = new Scanner(file);
209.                mydata = data.nextLine();
210.                data.close();
211.            }
212.            catch(FileNotFoundException e){
213.                System.out.println("\nFile not found!");
214.            }
215.            return mydata;
216.        }
217.    }
```

# Q-Learning Player code - Player.java

```
1.  import java.util.Random;
2.
3.  public class Player{
4.      int currentMove;
5.      String currentState;
6.      int lastMove;
7.      String lastState;
8.      Q q = new Q();
9.
10.     Player(String mind){
11.         q.importMind(mind);
12.     }
13.     public void nextAction(String state){
14.         currentMove = q.getBestAction(state);
15.         currentState = state;
16.     }
17.
18.     public void storeInLast(){
```

```java
19.         lastMove = currentMove;
20.         lastState = currentState;
21.     }
22.
23.     //Q(y,b)
24.     public double retQNextAction(String state){
25.         return q.getQ(state, q.getBestAction(state));
26.     }
27.
28.     public double retQ(String state, int action){
29.         return q.getQ(state, action);
30.     }
31.
32.     public void rewards(String state, int action, double reward){
33.         double oldQ = retQ(state, action);
34.         double newQ = oldQ + reward;
35.         //System.out.println("State: "+state+", Action: "+action+", Value: "+newQ);
36.         q.update_Q(state, action, newQ);
37.     }
38.
39.     public void directReward(String state, int action, double reward){
40.         double oldQ = retQ(state, action);
41.         q.update_Q(state, action, oldQ+reward);
42.     }
43.
44.     public void dump(){
45.         q.dump();
46.     }
47.
48.     public void dumpJSON(){
49.         q.dumpJSON();
50.     }
51.
52. }
```

# Running the game/tests.

I have wrote a simple bash script that can be used for compiling and running my code called Run.sh. I am using a JSON parsing library from the following location.
https://code.google.com/p/json-simple/
Simply include this in the directory to run/compile the code.
https://code.google.com/p/json-simple/downloads/detail?name=json-simple-1.1.1.jar&can=2&q=

## Run.sh Script

```
1.  |#!/bin/bash

2.  rm *.class

3.  javac -cp "./json-simple-1.1.1.jar:." RunStore2.java

4.  java -cp "./json-simple-1.1.1.jar:." RunStore2 $*
```

The RunStore2.java file contains the most recent version of my code. It simply accepts an int as a command line parameter. This is the amount of times you'd like to run the code.
For example (./Run.sh 1000) would run the code 1000 times and print the win/lose rations for the players. It would also by defaults store the Q values to the QData.json file. This file must be located in the run directory (it can be empty.)

## Game running code - RunStore2.java

```
1.  //Running our game

2.  public class RunStore2{

3.      public static int countRuns = 0;

4.      public static Player p1 = new Player("QData.json");

5.      public static RPlayer p2 = new RPlayer();

6.      public static void main(String[] args){

7.          Game game = new Game();

8.          int[] results = new int[4];

9.

10.         for(int i=0; i<Integer.parseInt(args[0]); i++){

11.             results[run(game)]++;

12.             countRuns += 1;

13.             if(countRuns % 100000==0){

14.                 System.out.println("1:"+results[1]+", 2:"+results[2]+",
```

```java
                draw:"+results[3]);
15.             }
16.         }
17.         for(int i=1; i<results.length; i++){
18.             if(i==3){
19.                 System.out.println("Draws: "+results[i]);
20.             }
21.             else{
22.                 System.out.println("Player "+i+": "+results[i]);
23.             }
24.         }
25.         p1.dumpJSON();
26.         //System.out.println("-------");
27.         //p1.dump();
28.     }
29.
30.     public static int run(Game game){
31.         double reward = 0.0;
32.         double yb = 0.0;
33.         double gamma = 0.5;
34.         String nextState = "";
35.         String x = "";
36.         int a = -1;
37.         String y = "";
38.         int b = -1;
39.         game.newBoard();
40.         while(game.whoWins() <= 0){
41.             if(game.turn == 1){
42.                 //Fill in our current move & state
43.                 x = game.getState();
44.                 p1.nextAction(game.getState());
45.                 while(!game.isValidAction(1,p1.currentMove)){
46.                     p1.nextAction(game.getState());
47.                 }
```

```
48.            a = p1.currentMove;

49.            y = game.getNextState(1, p1.currentMove);

50.            yb = p1.retQNextAction(y);

51.

52.            game.makeMove(1, p1.currentMove);

53.         }

54.

55.         //2nd player takes move

56.         else{

57.            x = game.getState();

58.            p2.nextAction(game.getState());

59.            while(!game.isValidAction(2, p2.currentMove)){

60.                p2.nextAction(game.getState());

61.            }

62.            a = p2.currentMove;

63.            y = game.getNextState(2, p2.currentMove);

64.            yb = p1.retQNextAction(y);

65.            game.makeMove(2, p2.currentMove);

66.         }

67.

68.         reward = gamma * yb;

69.         p1.rewards(x, a, reward);

70.      }

71.      if(game.whoWins() == 1){

72.         reward = 100.0;

73.         p1.directReward(game.getState(),a,reward);

74.         reward = 0.0;

75.      }

76.      else if(game.whoWins() == 2){

77.         reward = -100.0;

78.         p1.directReward(game.getState(),a,reward);

79.         reward = 0.0;

80.      }

81.      return game.whoWins();
```

```
82.     }
83.
84. }
```