

1. *Self-Play.* Suppose, instead of playing against a random opponent, the reinforcement learning algorithm described above played against itself, with both sides learning. What do you think would happen in this case? Would it learn a different policy for selecting moves?

Yes, I expect it to learn a different policy: an optimal policy against an optimal policy. (Apart from the exploratory moves.)

This can be shown by induction. For each state that necessarily leads to the end of the game, the greedy action will be the one that leads to a victory (if there exists such an action) and thus the value of this state after enough time steps will be really close to 1, if there exists a winning action and really close to 0 otherwise. Suppose that the values of states leading to the end of the game in maximum k steps are all really close to 0 or 1 reflecting who wins the game if both players play perfectly. For a state that leads to the end of the game in maximum $k + 1$ steps, all the actions lead to states leading to an end in maximum k steps, hence after enough time steps, the greedy action in this state will be optimal and the value will be close to 0 or 1.