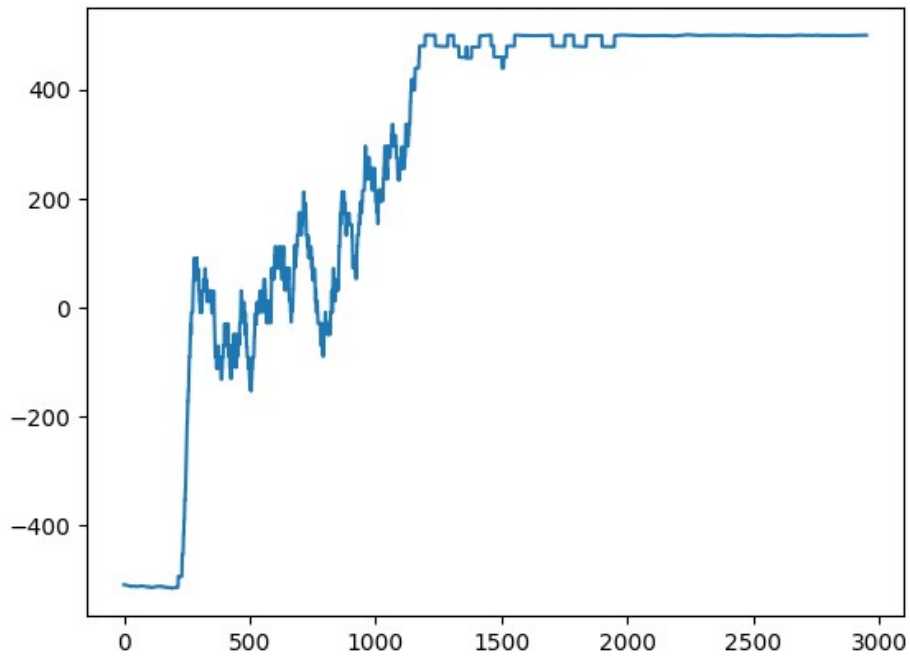


## RL Programming assignment 1 – Tabular Q learning – Moshe Beutel 037580792

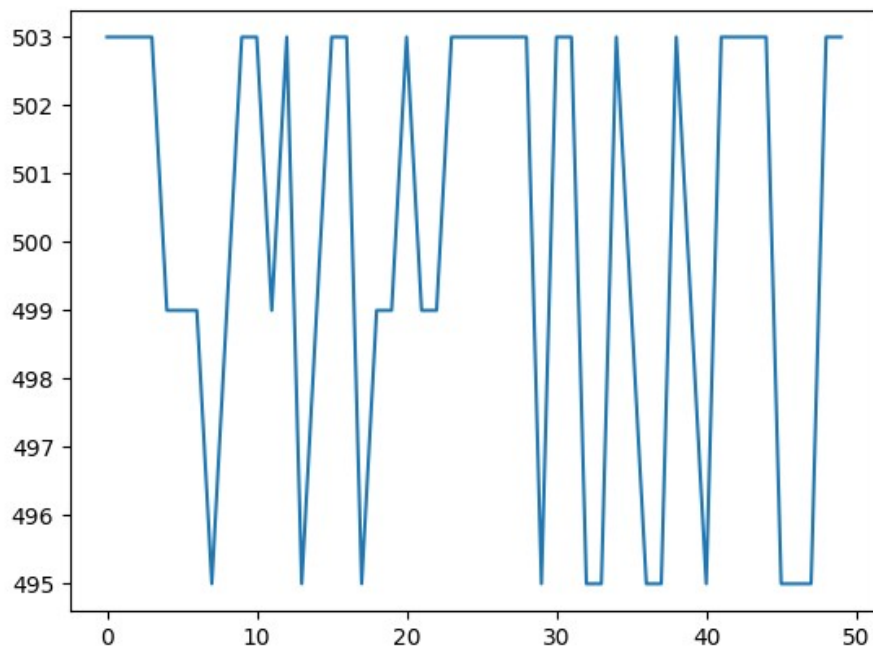
Questions 3+4

**epsilon = 0.0**

Train rewards 50 ep. avg.alpha\_0.2\_gamma\_0.8\_epsilon\_0.0

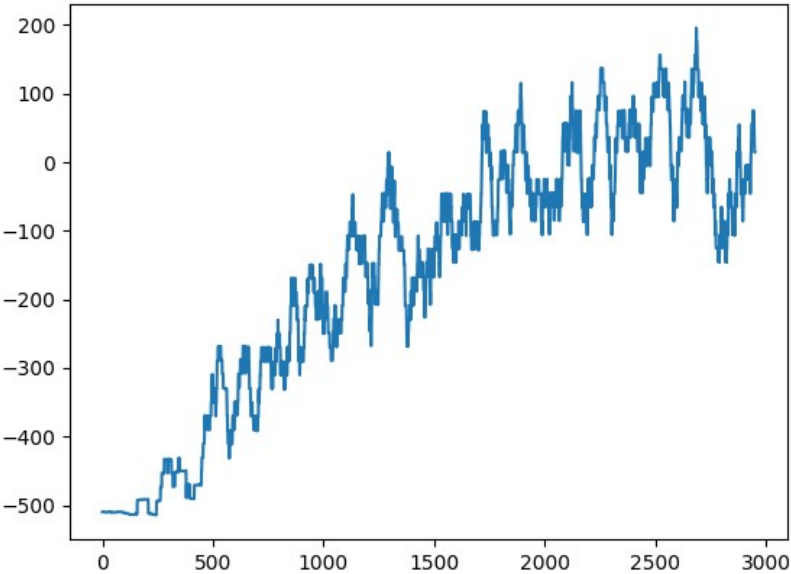


Test rewards after train with alpha\_0.2\_gamma\_0.8\_epsilon\_0.0 win rate 1.0

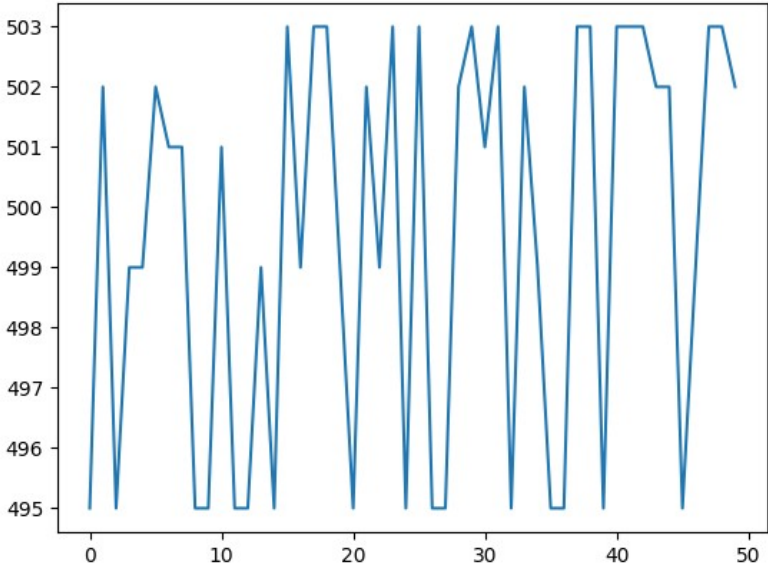


**epsilon=0.1**

Train rewards 50 ep. avg.alpha\_0.2\_gamma\_0.8\_epsilon\_0.1

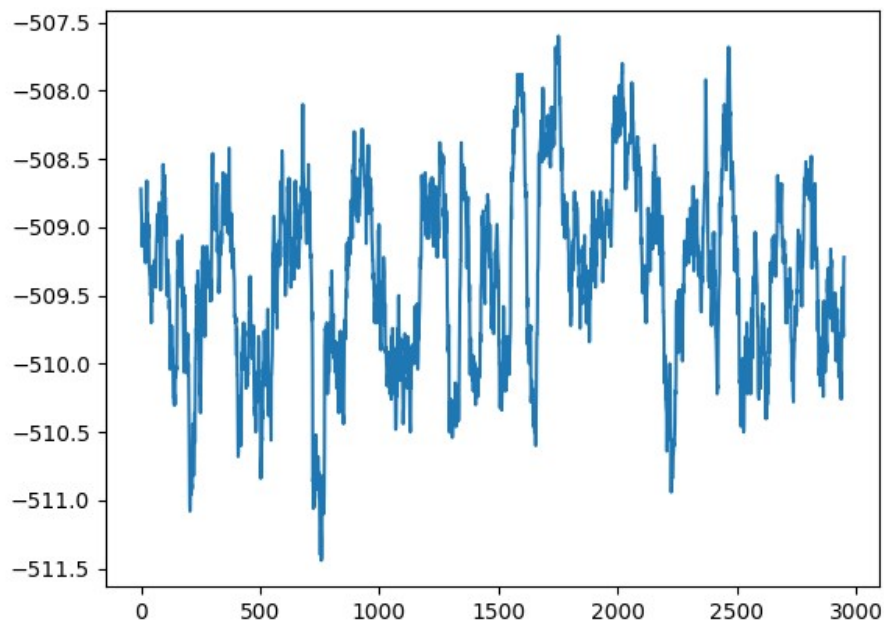


Test rewards after train with alpha\_0.2\_gamma\_0.8\_epsilon\_0.1 win rate 1.0

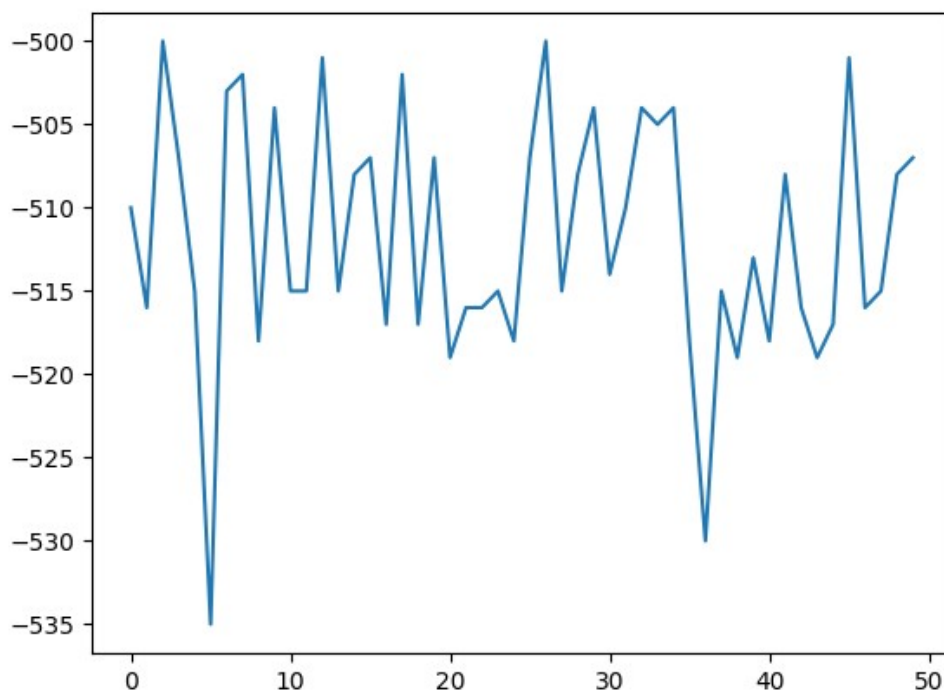


**epsilon=1.0**

Train rewards 50 ep. avg.alpha\_0.2\_gamma\_0.8\_epsilon\_1.0

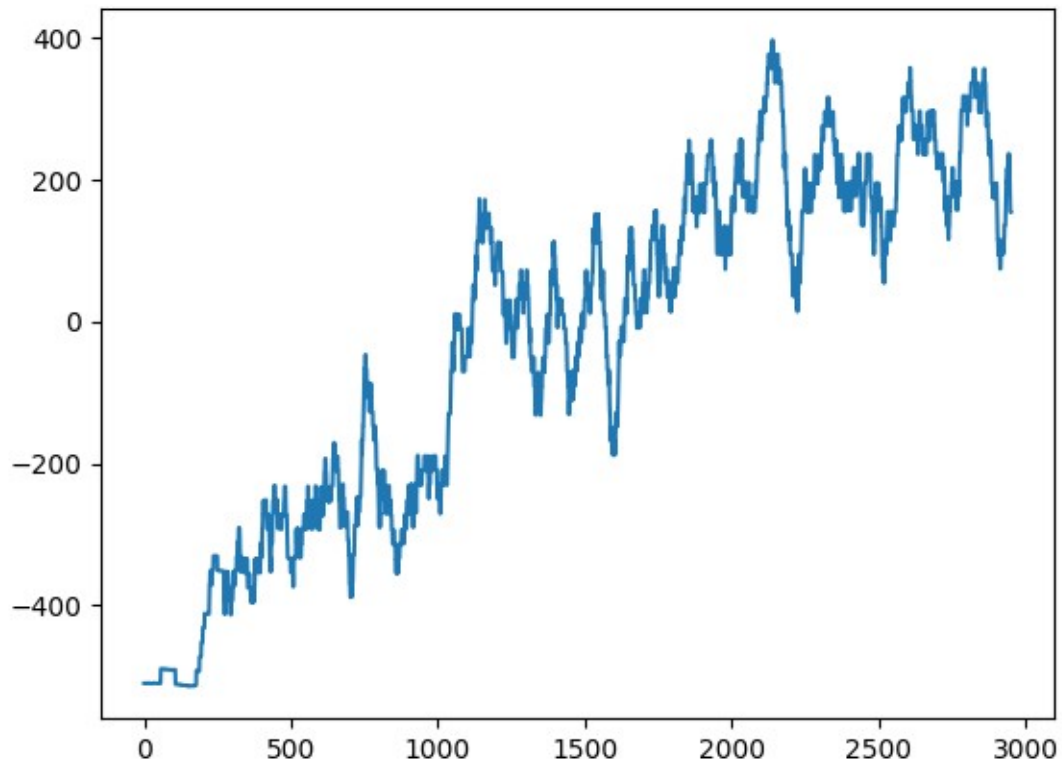


Test rewards after train with alpha\_0.2\_gamma\_0.8\_epsilon\_1.0 win rate 0.0

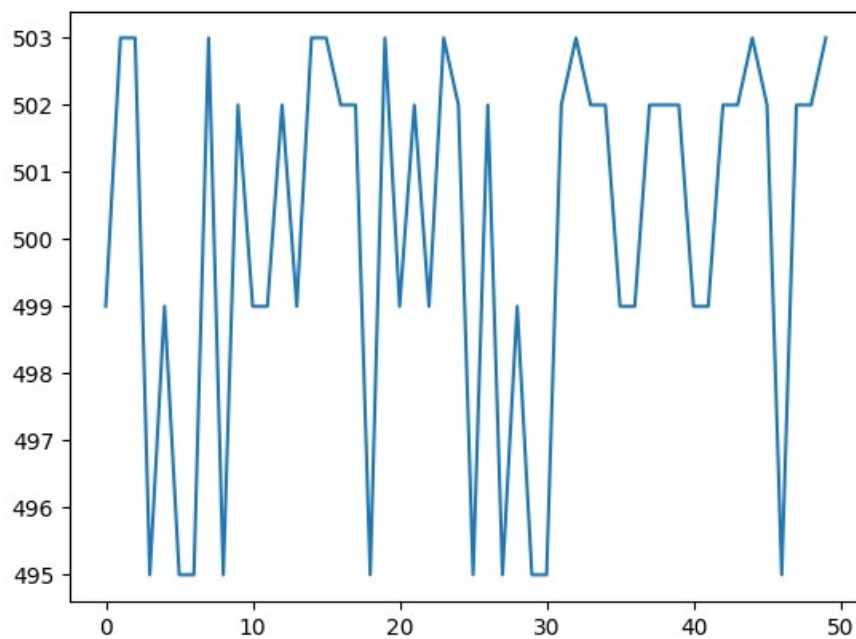


**alpha=0.1**

Train rewards 50 ep. avg.alpha\_0.1\_gamma\_0.8\_epsilon\_0.05

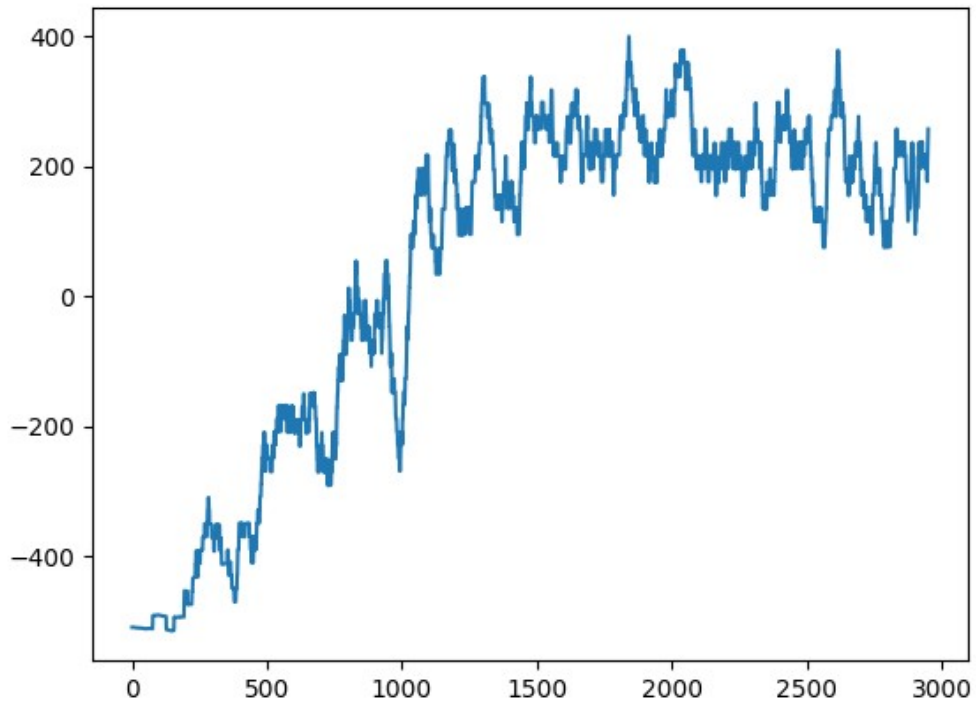


Test rewards after train with alpha\_0.1\_gamma\_0.8\_epsilon\_0.05 win rate 1.0

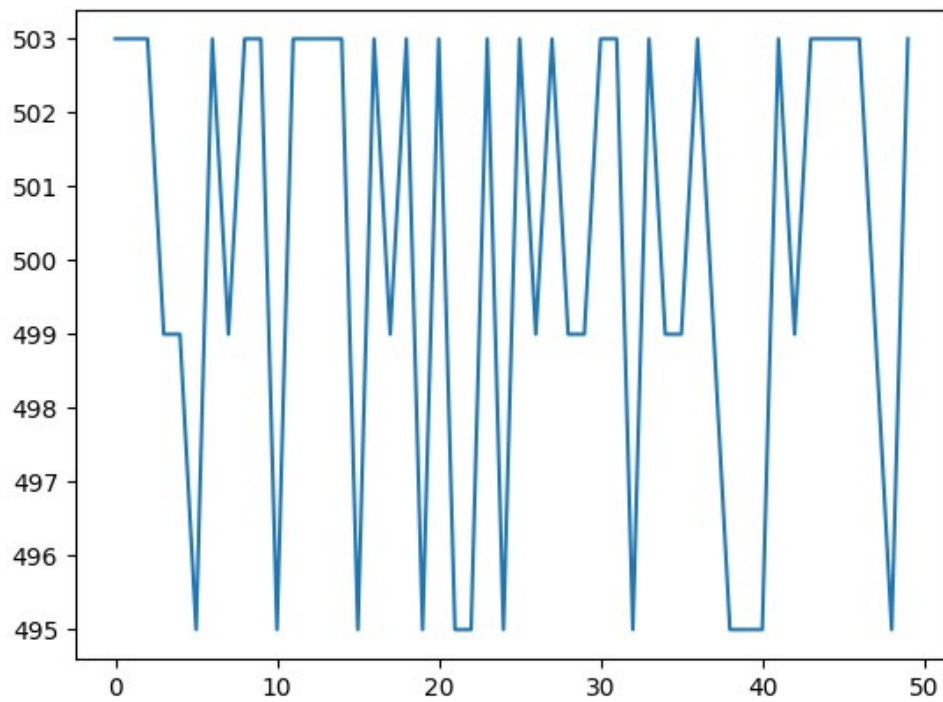


**alpha=0.2**

Train rewards 50 ep. avg.alpha\_0.2\_gamma\_0.8\_epsilon\_0.05

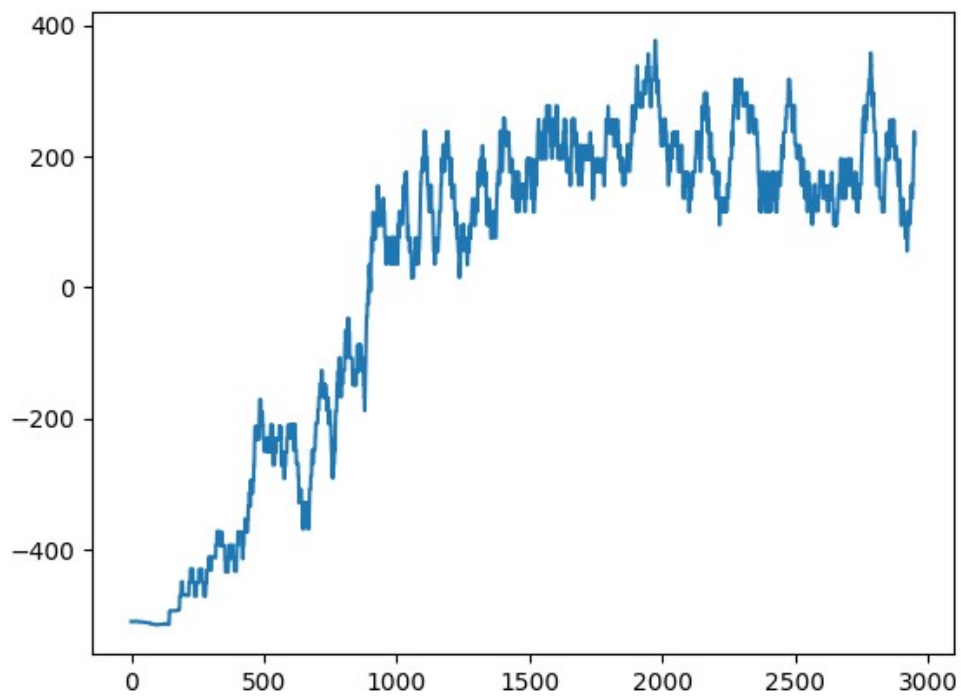


Test rewards after train with alpha\_0.2\_gamma\_0.8\_epsilon\_0.05 win rate 1.0

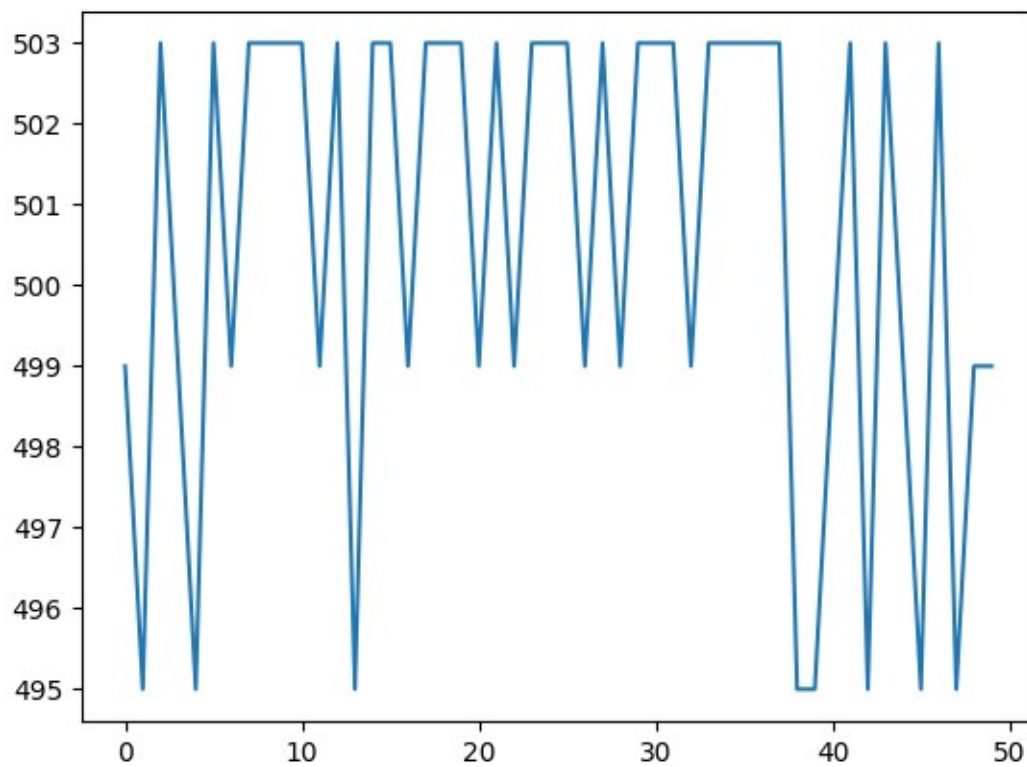


**alpha=0.3**

Train rewards 50 ep. avg.alpha\_0.3\_gamma\_0.8\_epsilon\_0.05



Test rewards after train with alpha\_0.3\_gamma\_0.8\_epsilon\_0.05 win rate 1.0



**Question 5.** The agent with no exploration works well because of the negative rewards. The Q table is initialized by default to 0 (I even initialized with 500 to increase exploration) so the max in the q update encourage the agent to try those 0 valued unseen state-actions.

**Questions 6+7.** The results on the medium layout were about 88%. The number of unseen states (or seen just a few times) is very large:

\*\*\*\*\*

518299 states in Q table

110509 states were visited less than 5 times

75590 states were visited once

\*\*\*\*\*

I made a few changes that raised the test average reward 100%:

When storing the state I zeroed the score and the direction (direction.STOP). The states that only differ in score or direction were unified

This made the q-table a lot smaller and learning more efficient.

\*\*\*\*\*

9148 states in Q table

1357 states were visited less than 5 times

847 states were visited once

\*\*\*\*\*