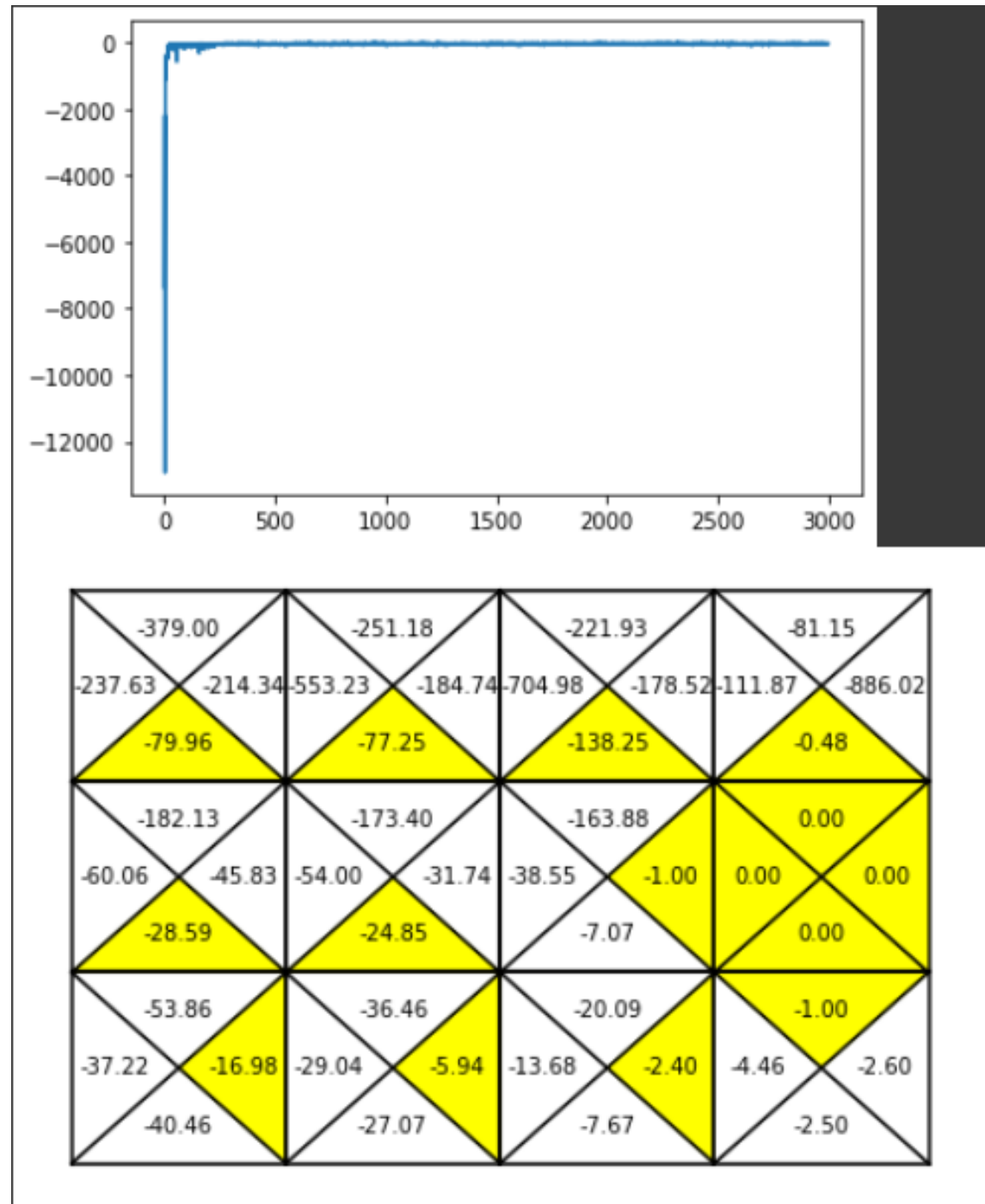


Experiments and Analysis (40%)

1. Plot the q _values in your result. (20%)

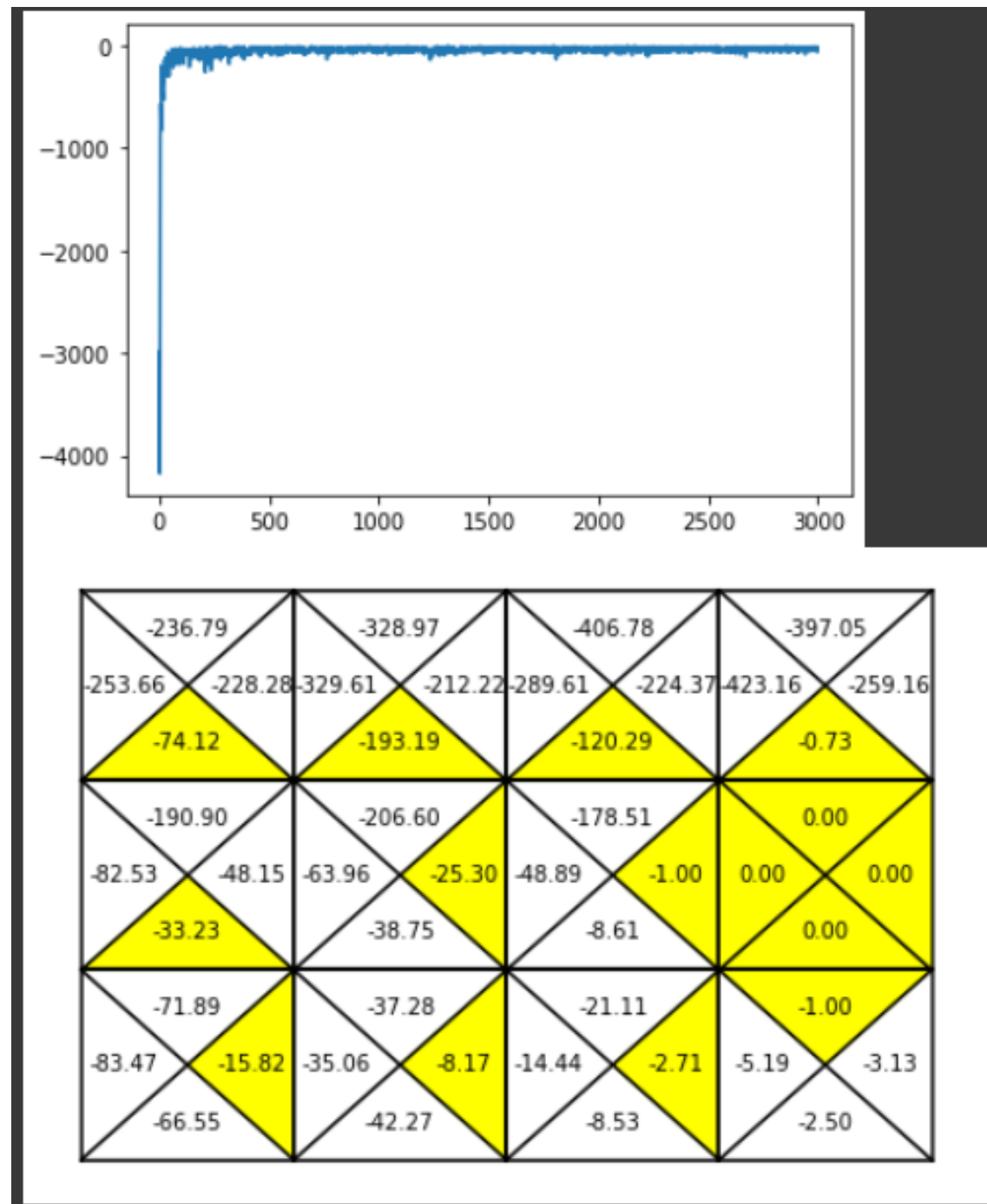
First-visit MC control:

State-action values and learning curve



Every-visit MC control:

State-action values and learning curve



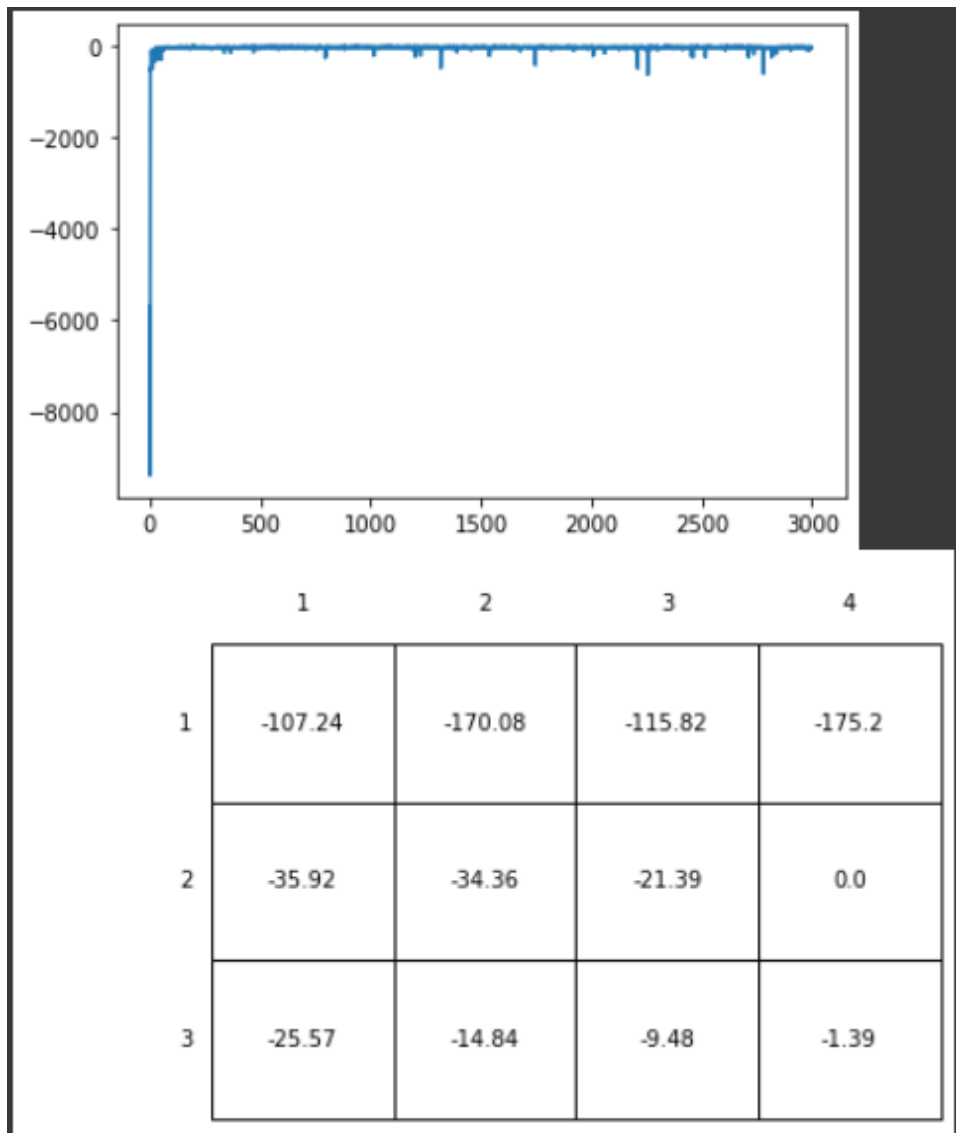
2. Whether q values are reasonable? Explain your result. (10%)

合理，可以從上圖中得知，他會選擇最遠的距離進行移動，因為他是蒙地卡羅法是使用 on-policy 的方法，而在沼澤地因為 reward 為 -100，所以他會全部往下走，離開沼澤地，所以 q value 向下的較大。

3. Transfer state-action values to state values and plot it. (10%)

First-visit MC predict:

State values and learning curve



Every-visit MC predict:

State values and learning curve

