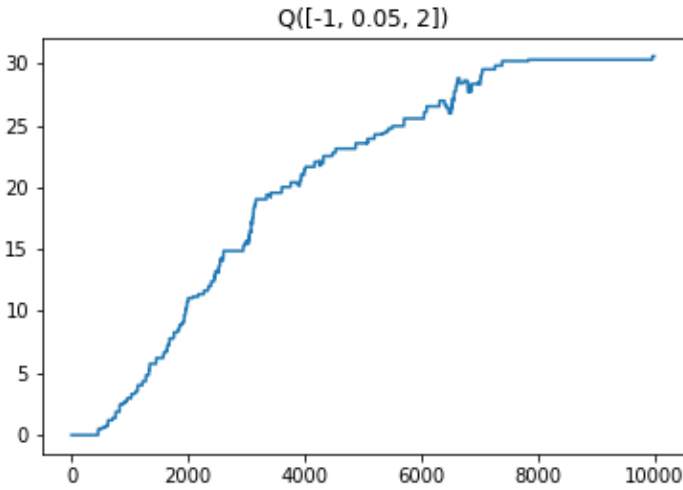
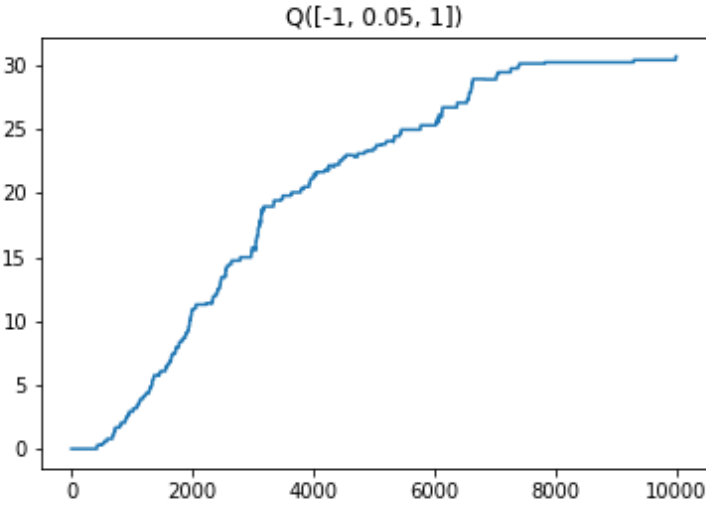
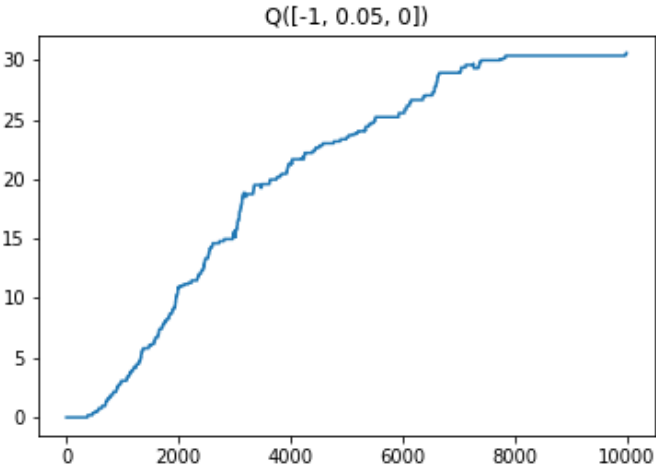
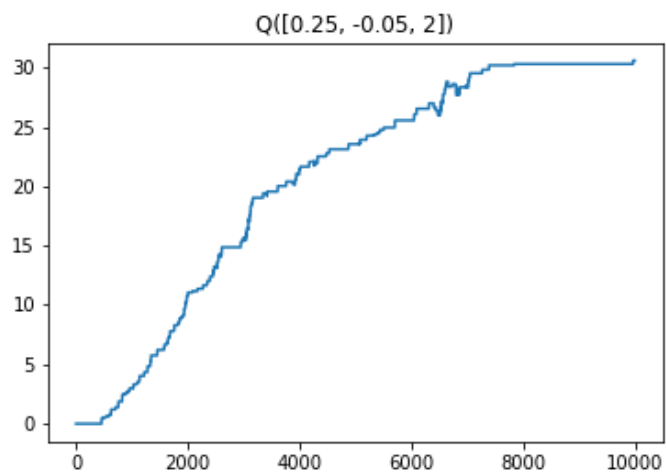
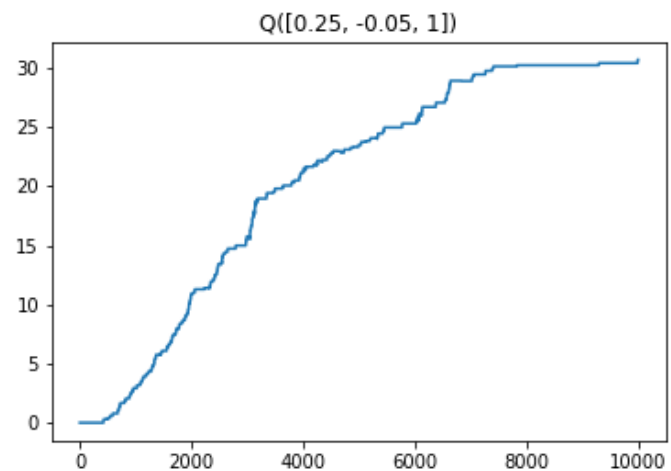
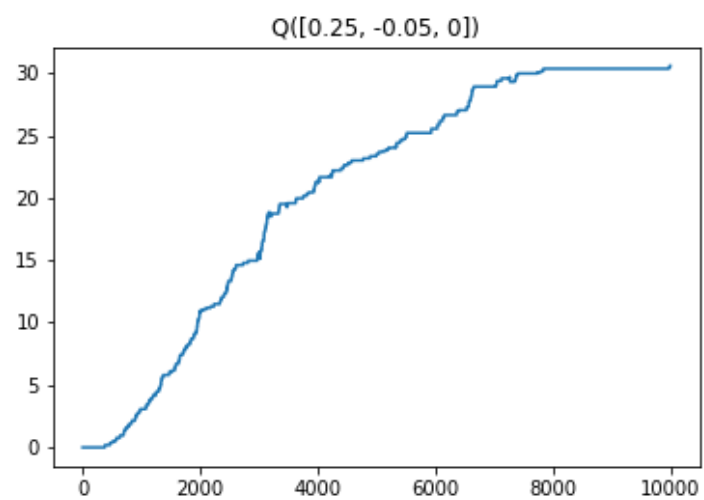


State 2:

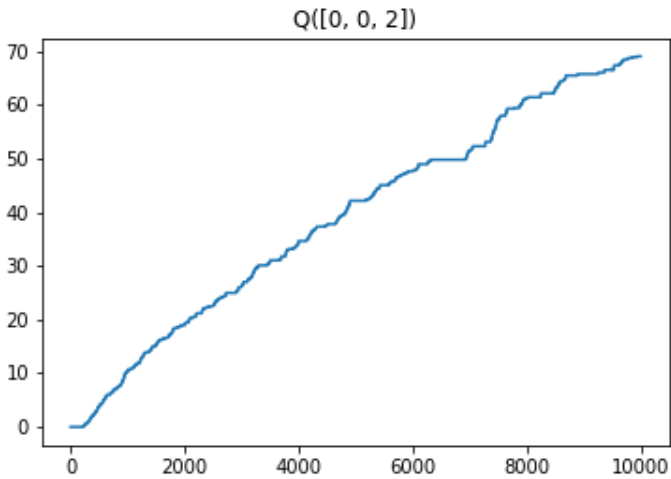
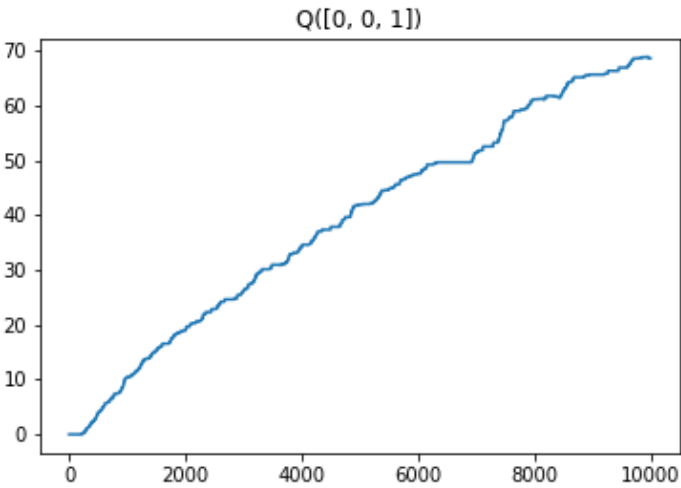
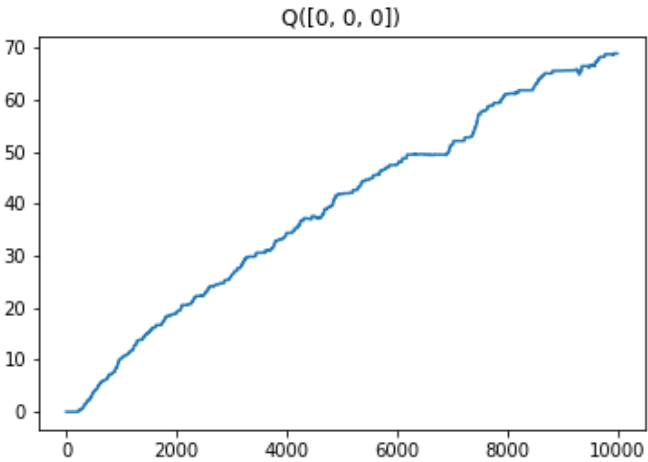


State 3:

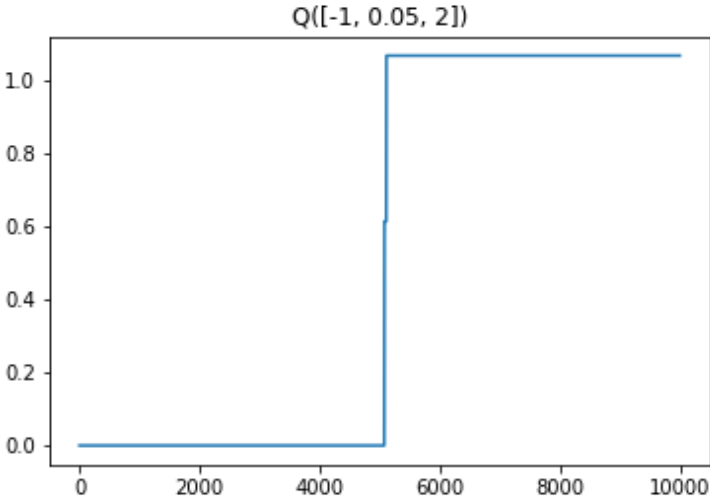
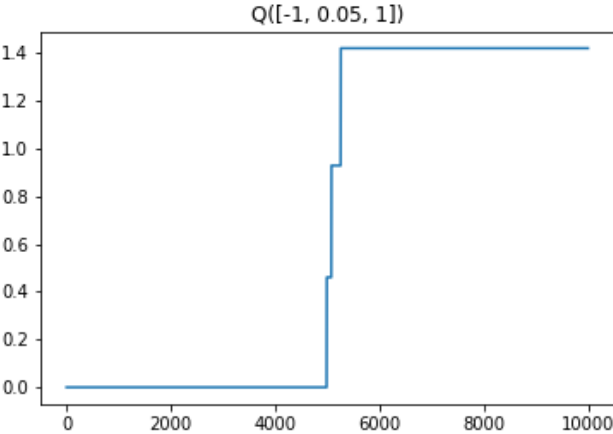
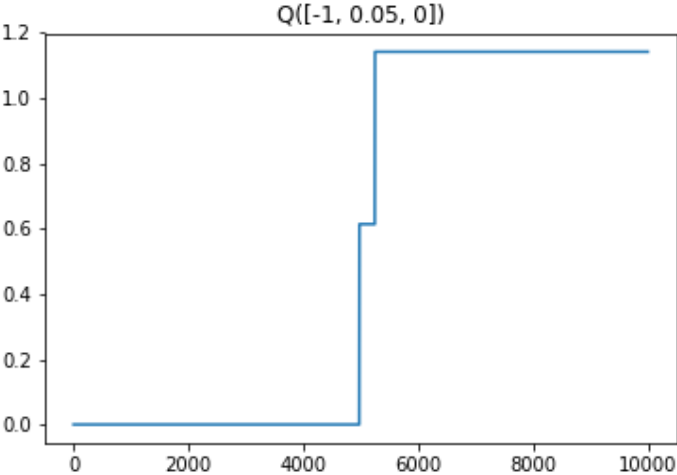


Off-policy TD with Q learning:

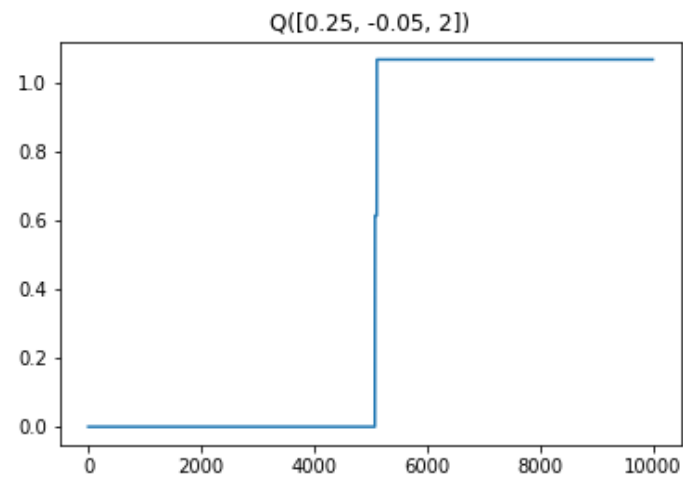
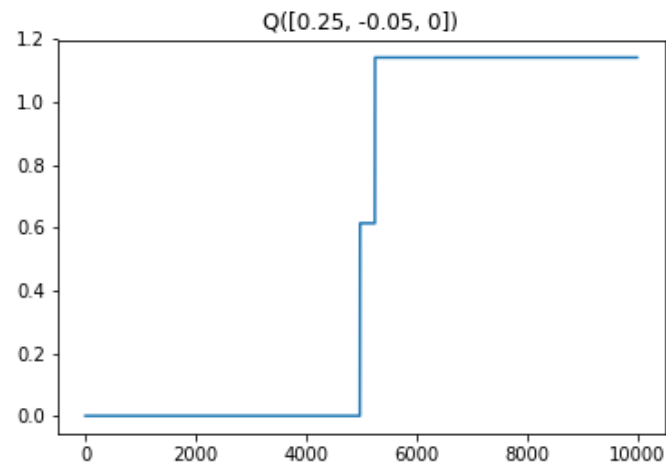
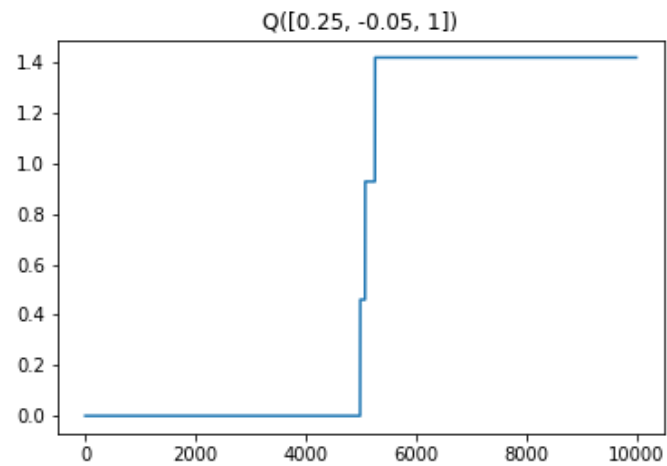
State 1:



State 2:



State 3:



Optimized Policy: (x-axis is velocity [from -0.7 to 0.7], y-axis is position [from -1.2 to 0.6] )

