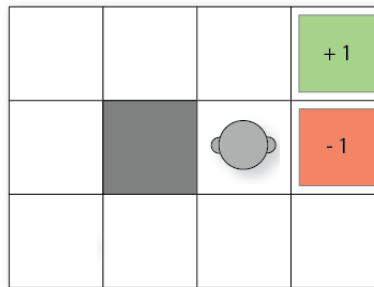


## Problem Set VIII: Monte-Carlo Tree Search

**Aim** The purpose of this workshop is to help you get a better understanding of Monte-Carlo Tree Search for solving MDPs in an online manner.

### Tasks

In this workshop, you will consider the example from the lectures of the agent that moves in a 2D grid world. Remember that if the agent tries to move in a particular direction, there is an 80% of success, and a 10% chance of it going to the left or right.



1. The agent is at cell (2,1), in which 2 is the x-coordinate and 1 the y-coordinate (both start from 0). It samples the following 5 iterations of MCTS, where E (East) goes right, W (West) goes left, N (North) goes up, and S (South) goes down:

Iter	Trace	
1	$N$	$simulate(succ) = -1$
2	$N \xrightarrow{succ} E$	$simulate(slip(S)) = -1$
3	$E$	$simulate(succ) = -1$
4	$W$	$simulate(succ) = 1$
5	$N \xrightarrow{succ} S$	$simulate(slip(W)) = 1$

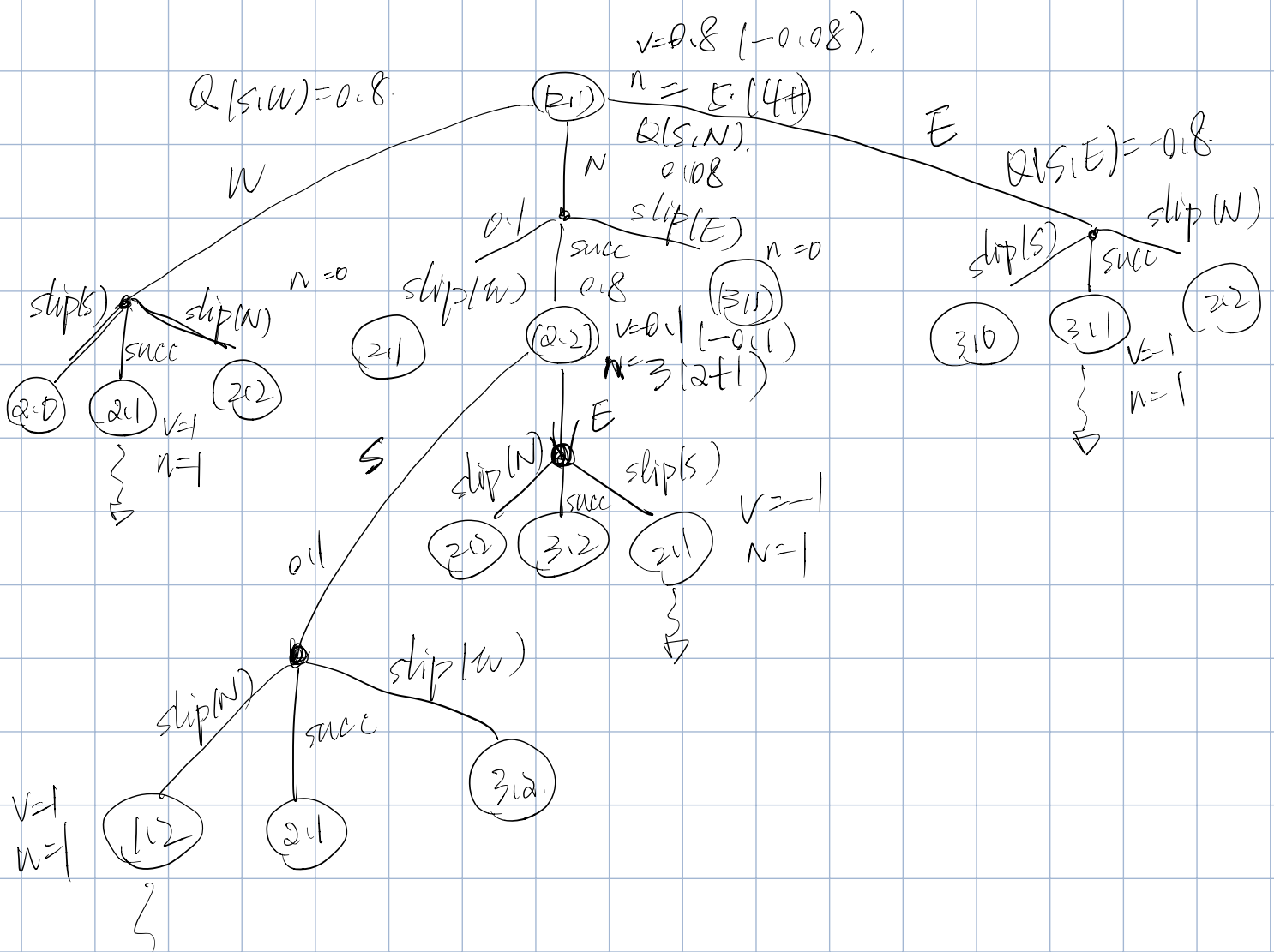
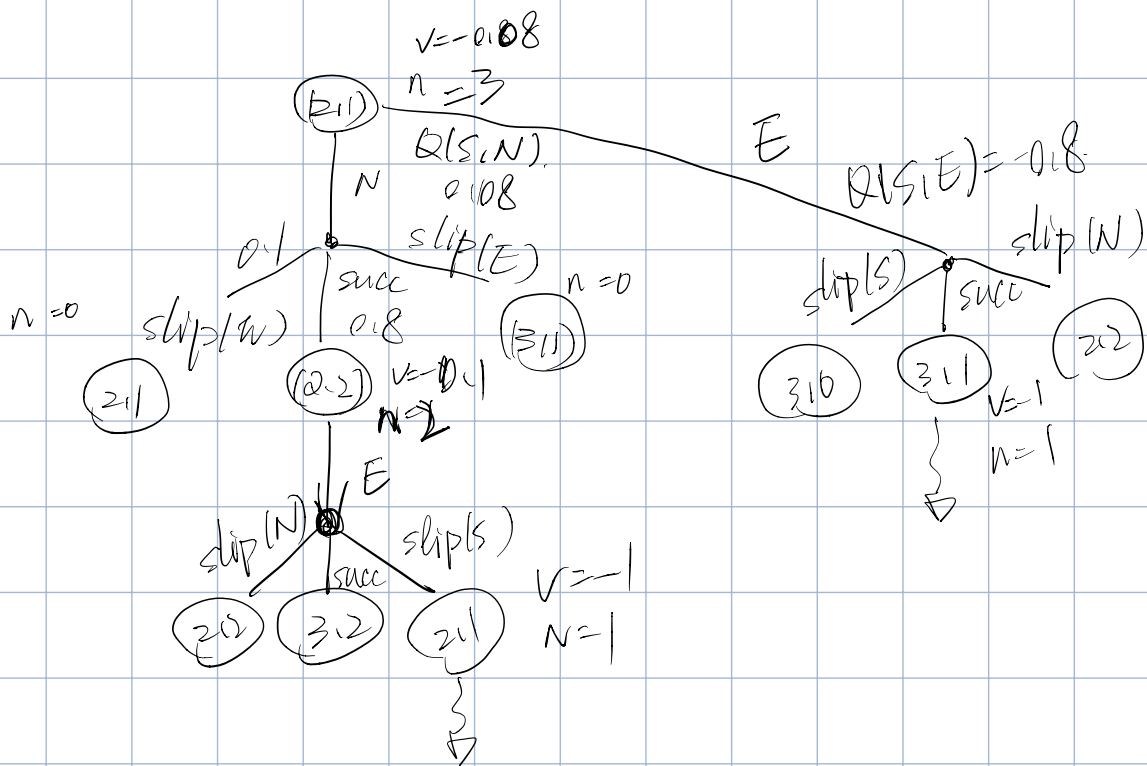
Here,  $N \xrightarrow{succ} E$  means that we select  $N$ , then select that the outcome of  $N$  was successful, and then select  $E$ . The notation  $N \xrightarrow{slip(E)} S$  means that we select  $N$  and the agent went East instead; then select South.

The notation  $simulate(X) = r$  means having just expanded a node, simulate from the outcome  $X$ , which will return reward  $r$ .

Draw the MCTS tree for this, assuming  $\lambda = 1.0$ . Label the lines on the tree with the actions & outcomes and label the nodes with the value  $V(s)$  for each state and  $N$  (the number of times this has been visited).

2. Based on your tree, calculate the action with the highest expected return.
3. Based on your tree, which of action, North, South, East, or West, would be more likely to be chosen if we use UCT to probabilistically select the next action? Show your working. Assume that  $C_p = \frac{1}{2}$ .





$$uct = R(s, a) + \gamma \sum_{N(s, a)} \sqrt{\frac{2 \ln(N(s, a))}{N(s, a)}}$$

$$uct(N) = 0.8 + \gamma \sqrt{\frac{2 \ln(5)}{3}}$$

$$uct(w) = 0.8 + \gamma \sqrt{\frac{2 \ln(5)}{1}}$$

$$uct(E) = 0.8 + \gamma \sqrt{\frac{2 \ln(5)}{1}}$$

$$uct(s) = 0 + \gamma \sqrt{\frac{2 \ln(5)}{0}} = \infty$$