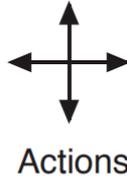
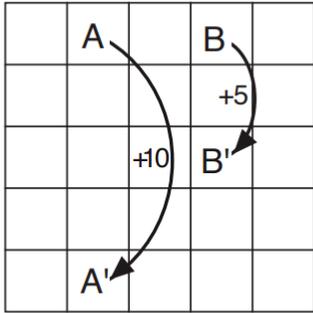


Worksheet C1M3

4. Consider the gridworld and value function in the figure below. Using your knowledge of the transition dynamics and the values (numbers in each grid cell), write down the policy corresponding to taking the greedy action with respect to the values in each state. Create a grid with the same dimension as the figure and draw an arrow in each square denoting the greedy action.



3.3	8.8	4.4	5.3	1.5
1.5	3.0	2.3	1.9	0.5
0.1	0.7	0.7	0.4	-0.4
-1.0	-0.4	-0.4	-0.6	-1.2
-1.9	-1.3	-1.2	-1.4	-2.0

5. Consider the continuing MDP shown on the bottom. The only decision to be made is that in the top state, where two actions are available, left and right. The numbers show the rewards that are received deterministically after each action.
- List and describe all the possible policies in this MDP.
 - Is the following policy valid for this MDP (i.e. does it fit our definition of a policy): Choose *left* for five steps, then *right* for five steps, then *left* for five steps, and so on? Explain your answer.
 - What policy is optimal if $\gamma = 0$? If $\gamma = 0.9$? If $\gamma = 0.5$?
 - For each possible policy, what is the value of state s ? Write down the numeric value to two decimal places. *Hint*: write down the return under each policy starting in state s (don't forget γ). Simplify the infinite sum, using the fact that many rewards are zero. Then plug in the rewards and γ and compute the number.

