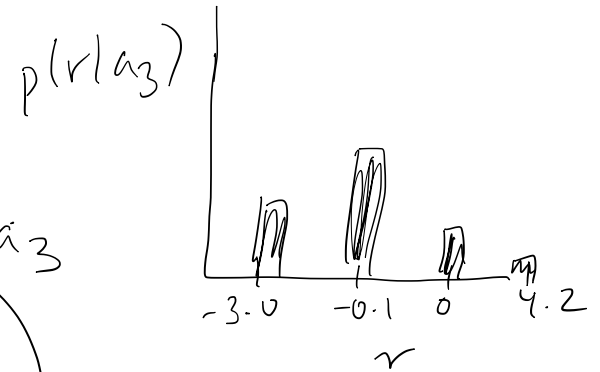
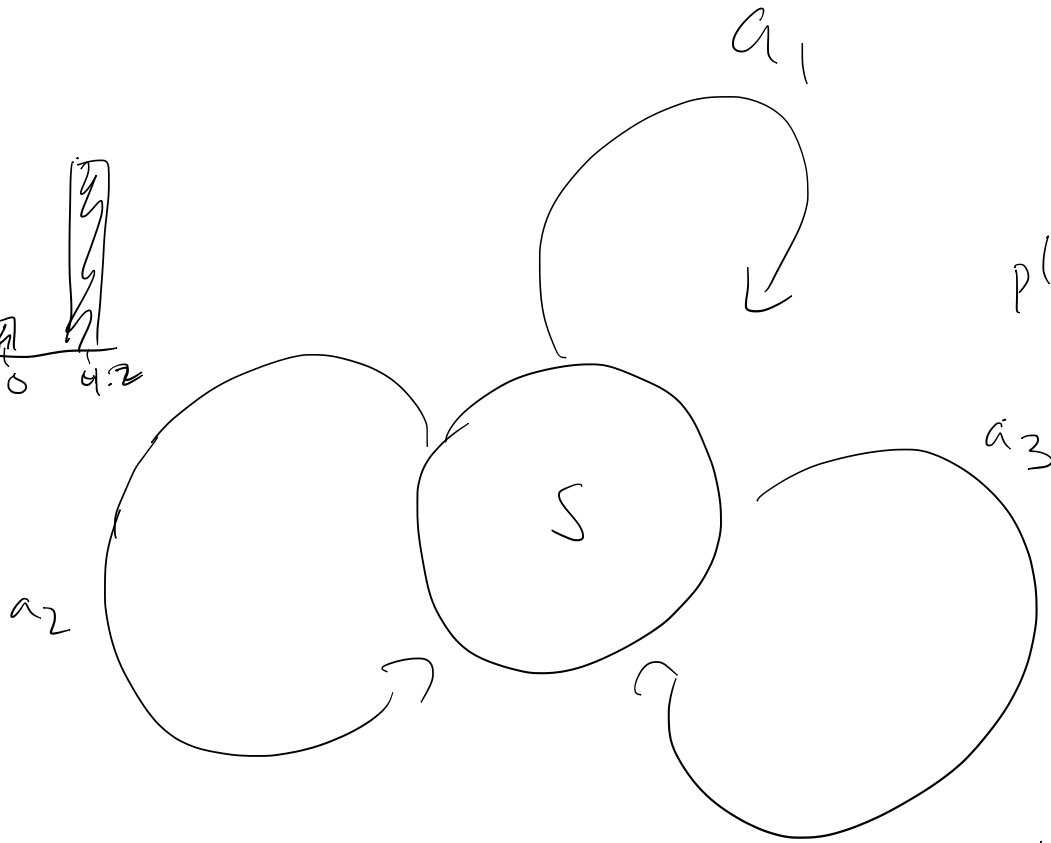
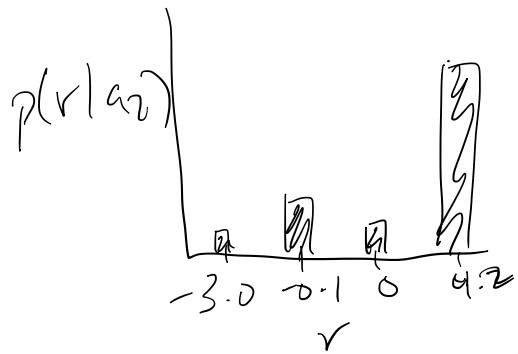


Bandit as MDP:



- Single state, agent always transitions back to  $S$  for every action
- Different distributions over rewards for each action