

## RL

Task 1:

a)  $S = \{5, 4, 3, 2, 1\} \cup \{\text{Parked}, \text{Failed}\}$

A: P (try to park in current parking space)  
D (Drive on to next space)

$p$ : if  $s' = \text{Parked} \rightarrow p(s', 1 | s, T) = p$

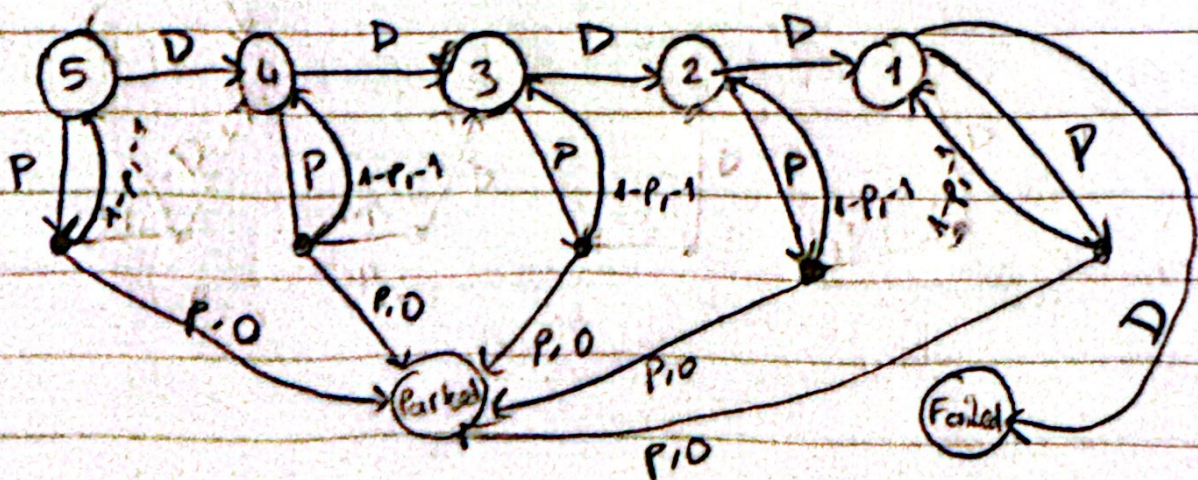
if  $s' = s-1 \rightarrow p(s', -1 | s, D) = 1$

$\rightarrow p(s', -1 | s, P) = 1-p$

if  $s=1$  and  $a=D \rightarrow p(\text{Failed}, -10 | 1, D) = 1$

$R: \{+1, -1, -10\}$

b)





c) Discount can be applied because we want to park as quickly as possible and therefore future rewards should be decreased by a discount factor.

Task 2) Since the next state is determined not only from the current state but also the previous state (last 2 states), Markov property is not fulfilled.

Task 3)

Discrete

Consider  $S = \{s_1, s_2\}$   
 $A = \{a_1, a_2\}$

In  $s_1$

→  $a_1$ : transitions to  $s_2$  with  $r=1$   
 $a_2$ : " " " "

In  $s_2$

→ both actions go back to  $s_1$  with a reward of 0

So  $J^*(s_1) = 1$ ,  $J^*(s_2) = 0$

and these don't change with different policies

$\pi^*(s_1) = a_1 = a_2$

$\pi^*(s_2) = a_1 = a_2$

optimal deterministic policy is not unique.