

Bike MDP

* state :- state at time 't' denoted via ' S_t '
 $S_t = [b_t^1 \ b_t^2]$

$S_t \rightarrow$ Denotes No. of bikes available at the end of the day at particular time 't'

$b_t^i \rightarrow$ No. of bicycle at station 'i' at time 't'

* Action state $A([b^1, b^2])$

Considering the Notation -

- Transferring bicycle from stations $1 \rightarrow 2 = +ve$
- Transferring bicycle from stations $2 \rightarrow 1 = -ve$

eg1. Consider the state = $[2, 6]$

- Here we can transfer 5 bicycles from station 1 to 2
- Also we can transfer 2 bicycles from stations 2 to 1

so Action state = $[-5, \dots, 2]$

eg2. Consider the state $[10, 18]$

- Here we can transfer a maximum of 2 bicycles from station $1 \rightarrow 2$ as for every station maximum capacity is limited to 20.
- Also we can transfer 5 (max transfer limit) bicycles from $2 \rightarrow 1$

so Action state = $[-5, \dots, 2]$

from $2 \rightarrow 1$ ← → from $1 \rightarrow 2$

So Generally,
Action state $A([b^1, b^2]) =$

$$\{ -\min(\min(5, b_t^1), 20 - b_t^1), \min(\min(5, b_t^2), 20 - b_t^2) \}$$

* Value Function:- It maps states to real values.

$$V_{\pi} : S \rightarrow \mathbb{R}$$

$$|S| = 21 \times 21 \text{ (Size of state)}$$

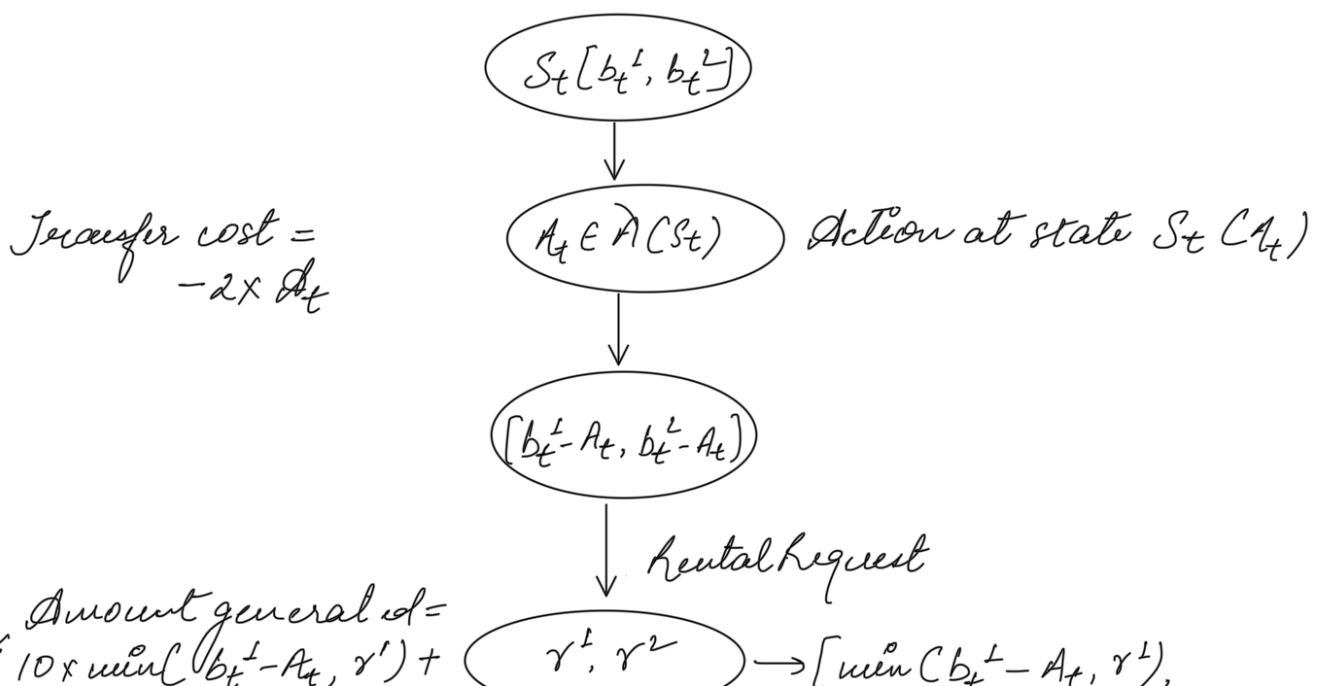
* Action Value Function -

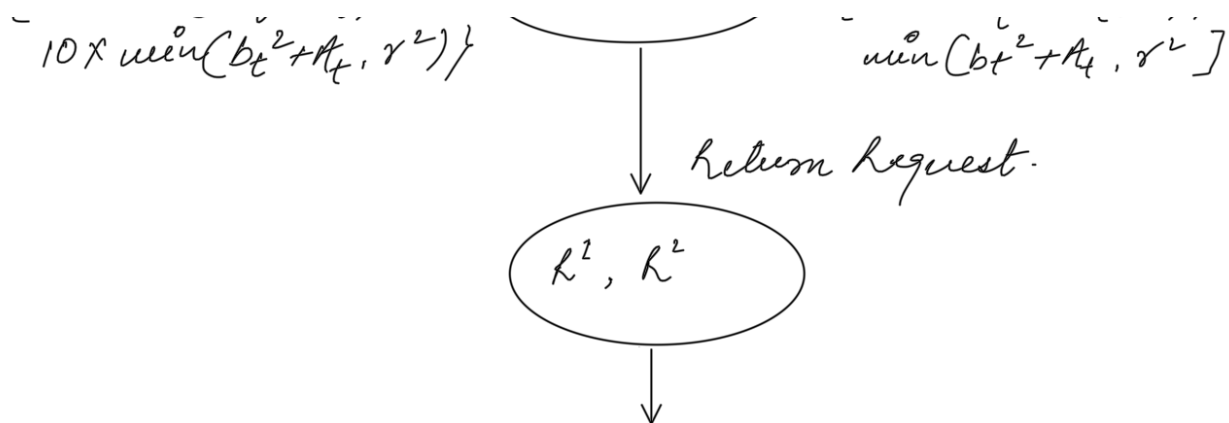
$$Q_{\pi} : S \times A \rightarrow \mathbb{R}$$

* Transition Probability -

$$P(S_{t+1}, R_{t+1} | S_t, A_t)$$

It means that at the state S_t taking action A_t we reach a new state S_{t+1} and received reward R_{t+1}





$$S_{t+1} = [b_{t+1}^1 = \min(b_t^1 - A_t - \min(b_t^1 - A_t, r^1) + R^1, 20) \\ b_{t+1}^2 = \min(b_t^2 + A_t - \min(b_t^2 + A_t, r^2) + R^2, 20)]$$

$$P(S_{t+1} | S_t, A_t) = P(r^1) P(r^2) P(R^1) P(R^2)$$

→ Probability with which we reach from state S_t to S_{t+1} taking action A_t

Probability of receiving rental request r^1 & r^2

$$P(r^1, r^2) = P(r^1) P(r^2) \rightarrow \text{Independent.}$$

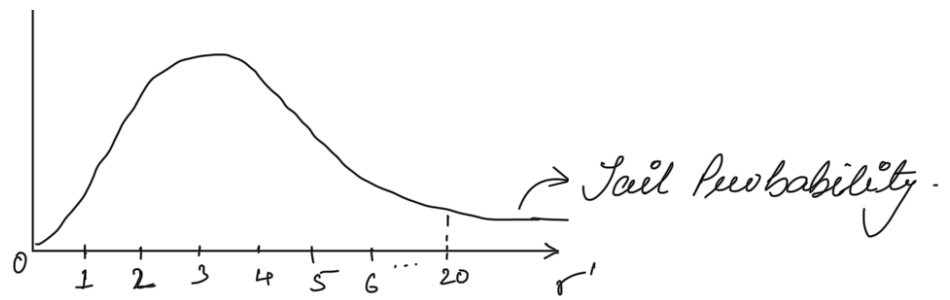
$$= \frac{1_1 r^1 \cdot e^{-r_1}}{r^1!} \cdot \frac{1_2 r^2 \cdot e^{-r_2}}{r^2!}$$

Similarly probability of return request will be same as that of rental request with different λ values.

* Reward at S_{t+1}

$$R_{t+1} = \left\{ -2 \times A_t + 10 \times \min(b_t^1 - A_t, r^1) + 10 \times \min(b_t^2 + A_t, r^2) \right\}$$

$$P(r^1) \uparrow$$



To get a meaningful distribution we divide the tail probability to each values from 0...20.