

APS 1080 LECTURE 2

REINFORCEMENT LEARNING

If observability is missing and if what you see could be in different states

Markov Decision Processes require observability.



$$A = \{A^{(1)}, \dots, A^{(n)}\}$$

The goal

Goal

Reward signal does not come intrinsically from the environment.

-1 → for one step missed
-100 → loss
100 → win

Agents mechanism to generate an action A based on S is called a policy.

$$\pi(A|S)$$

Policy of an action given a state.

n - 1 is to design a distribution for the policy

the goal is
so that an agent interacts with the
environment appropriately

$$A = \{ A^{(1)}, \dots, A^{(A)} \}$$

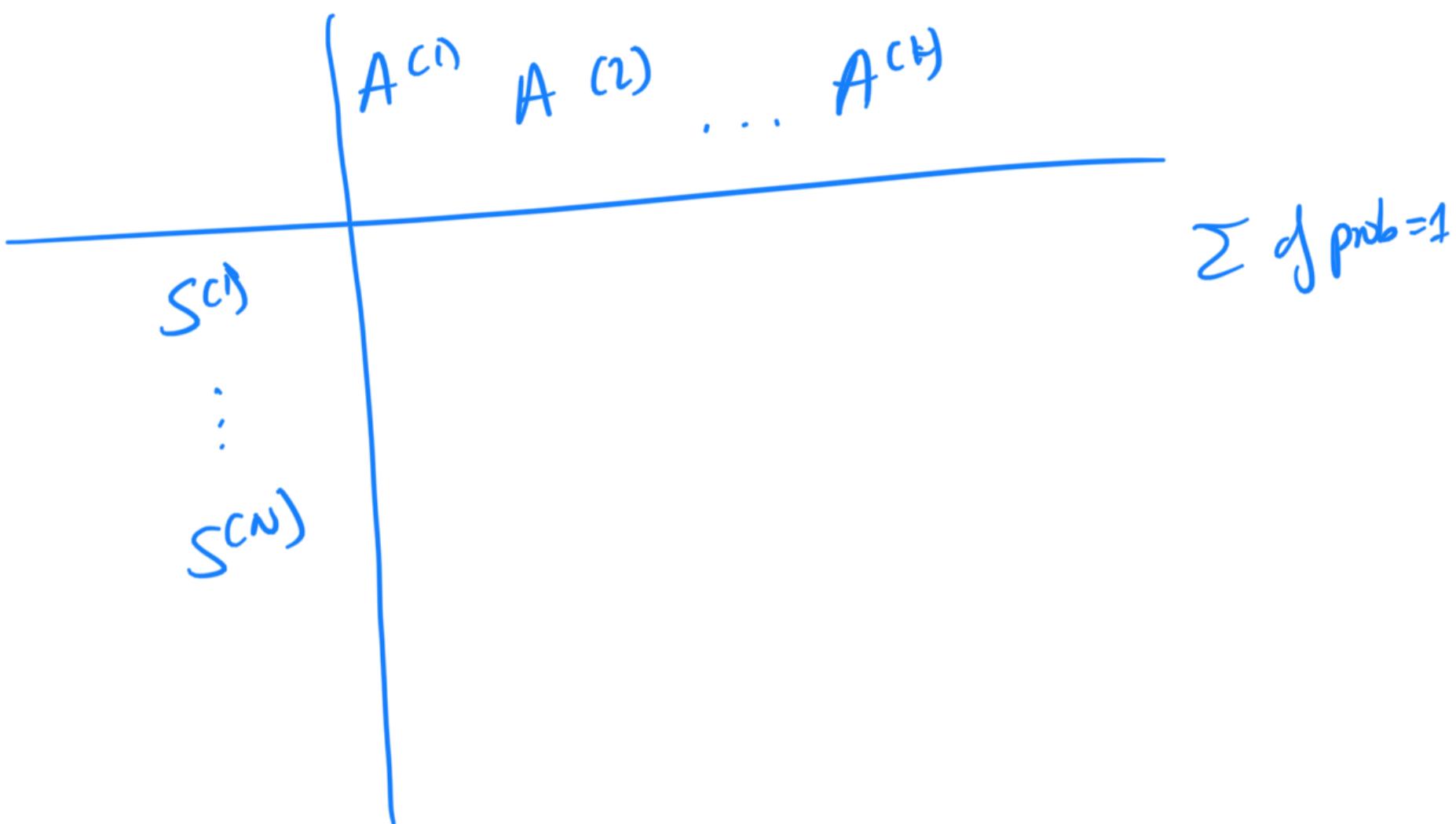
↑
index of action



This sequence is called an episode or a trace

what action do I select
given a state?

11



The sum of the probs must equal 1.

We may initialize the policy to random values.
What matters is not the initial condition but rather
how those values evolve.

How is the agent selecting the best action?

Subsequent sequence of rewards is maximized.

From time step i to the end of the episode
at time t .

We want to maximize the sum of future rewards.

Definition

Return

11

$$G_T = \sum_{k=t+1}^T \gamma^{k-(t+1)} R_k$$

Cumulative sum of all future rewards

$\gamma \rightarrow$ discount factor

\hookrightarrow A number between 0 and 1 if $T \rightarrow \infty$

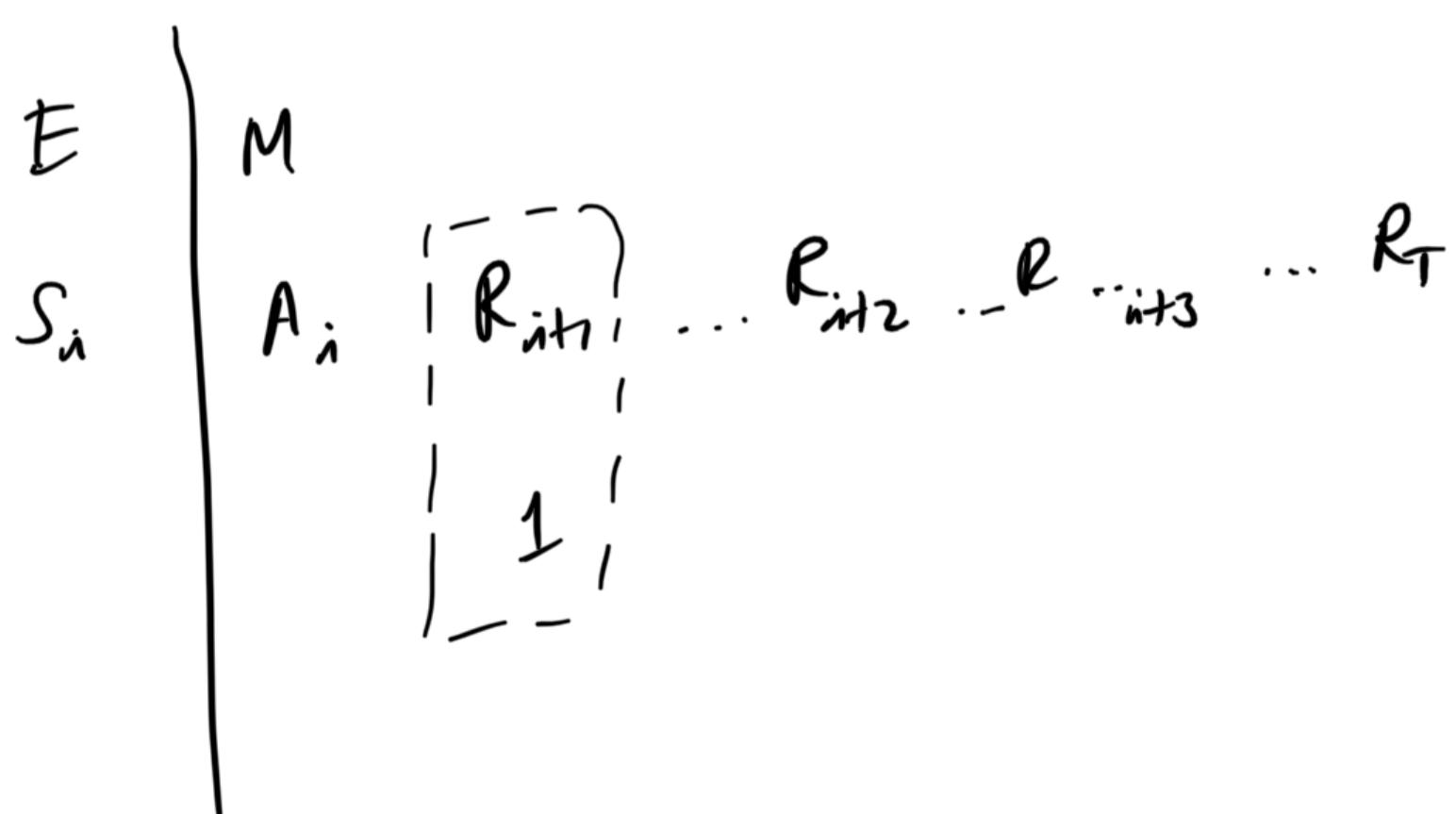
$$\gamma \in [0, 1) \quad T \rightarrow \infty$$

$$\gamma \in [0, 1] \quad T \rightarrow \text{finite}$$



As $\gamma \rightarrow 0 \rightarrow$ distant rewards get attenuated

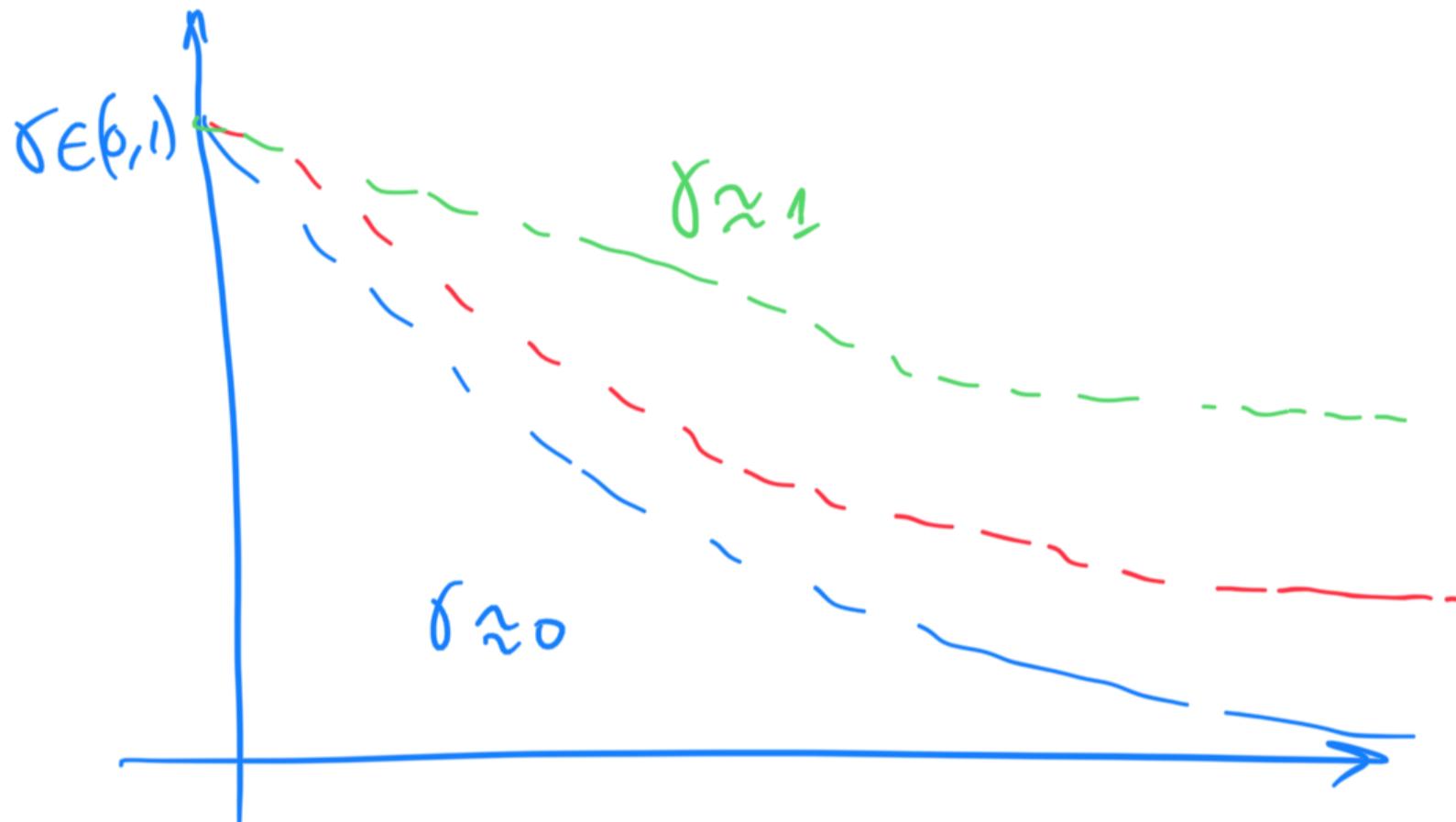
$\gamma \rightarrow 1 \rightarrow$ short term rewards get attenuated



Maximize feature reward

we want to maximize all the things related to my work

We always account for next reward with full strength



long term view $\rightarrow \gamma > 0$

short term view $\rightarrow \gamma = 0$.

The ACTION SELECTION PROBLEM

Given S_i , choose A_i such that we maximize the cumulative future rewards (G_{ti}) which is called the return.

Our goal is to create a policy that maximizes the cumulative return G_i .

→ How will the rewards look like for every combination of actions

in different episodes

→ Suppose we can learn by historical experience how being in state s_t relates to the return G_t for all states.

→ How being in s_t

Combinations of states and actions and how they relate to the returns that they generate

VALUE FUNCTIONS

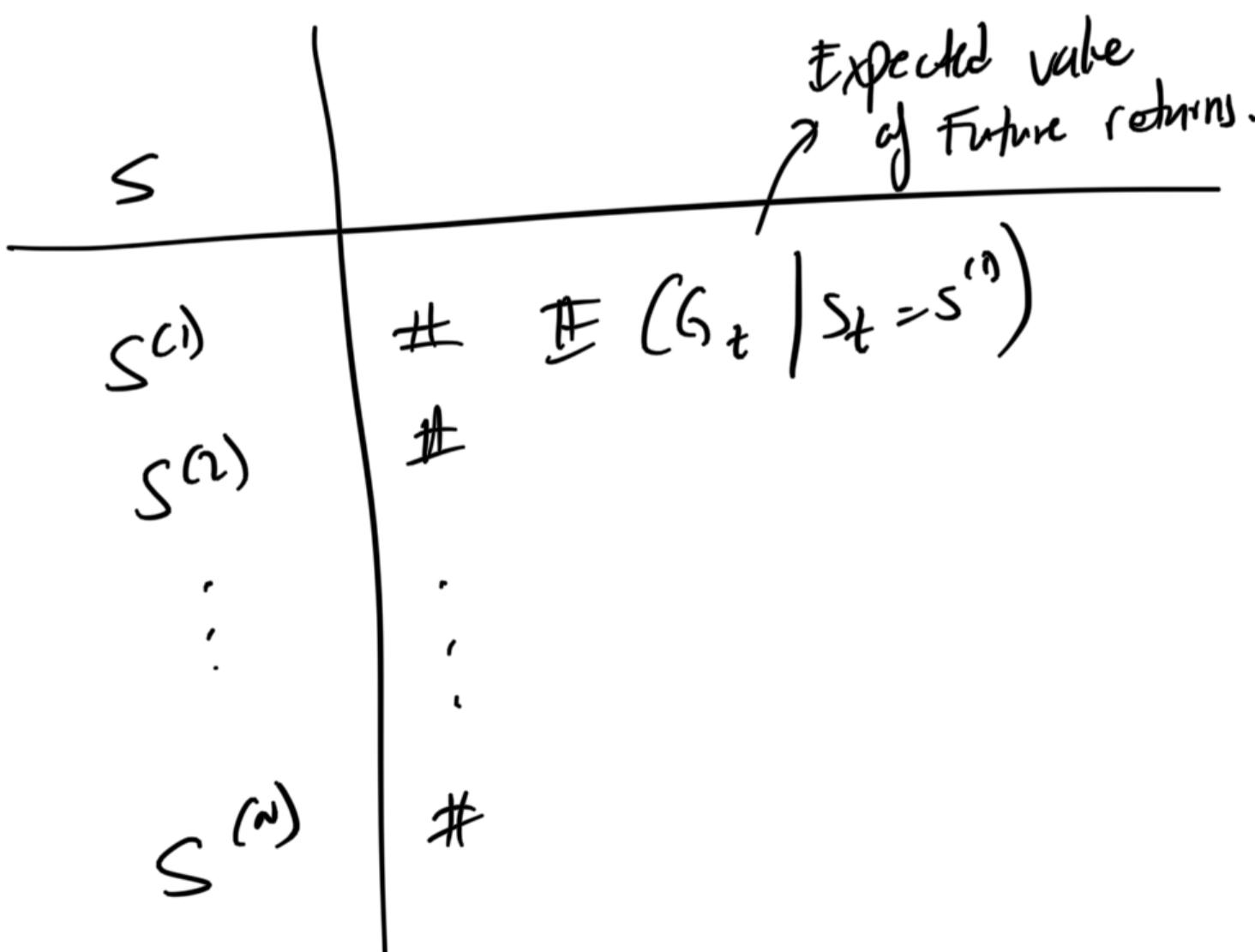
keep track of what results in the best result historically

state value function

$$V_{\pi}(s) = \mathbb{E} [G_t \mid S_t = s] \quad \forall s$$

↓
Return

for
all states



We populate this table by processing the traces or sequences that we pick up by having the agent interacting with the environment.

ACTION VALUE FUNCTIONS ($q_{\pi}(s,a)$)

Expected future return given a current state and a current action.

$$q_{\pi}(s,a) = \mathbb{E}[G_t | S_t = s, A_t = a]$$

$S \times A$		
$S^{(1)}$	$A^{(1)}$	#
	$A^{(2)}$	#
	\vdots	
	$A^{(k)}$	
$S^{(2)}$	$A^{(1)}$	
	$A^{(2)}$	
	\vdots	
	$A^{(k)}$	
$S^{(n)}$	$A^{(1)}$	
	$A^{(2)}$	
	\vdots	
	$A^{(k)}$	

Given a state, choose the action that maximizes the reward.

$Q_{\pi}(s, a) \rightarrow$ helps to solve the Action SELECTION PROBLEM
as the maximum value can be computed directly.
value function

If you have access to the V value function

$V_{\pi}(s) \rightarrow$ What is the expected return being in each state?
which state is the best one to be in?

$s^{(1)}$	#
$s^{(2)}$	#
:	:
$s^{(k)}$	#

Some states are more desirable than others,
 $V_{\pi}(s)$ tells you which are the most desirable States to be in.

Know what the tables look like for the TEST

You can select a value with q only
→ with V , you don't have a mean to select an

action.

With V alone you need a model of the environment to make it useful.

$$p(s', r, s, a) = \Pr \left\{ S_{t+1} = s', R_{t+1} = r \middle| S_t = s, A_t = a \right\}$$

$S \times A$	Next states (up to N)		next rewards up to M.	
	$S^{(1)}$	$S^{(2)}$	$\dots S^{(N)}$	$R^{(1)}$ $R^{(2)}$ $\dots R^{(M)}$
$S^{(1)}$	$A^{(1)}$	\dots	\dots	\dots
	$A^{(2)}$			
	\vdots			
	$A^{(k)}$			
	\vdots			
$S^{(n)}$	$A^{(1)}$			
	$A^{(2)}$			
	\vdots			
	$A^{(N)}$			

$\sum = 1$ for
each state action

$\sum = 1$ for
each state action

We want to choose an action that has the highest likelihood of getting me to the state that gives me the highest value. (S^*)

Given a stochastic model of the environment and a value function, how can you solve the Action selection Problem

With all you need

If you have $q \rightarrow$ then π or r

If you have $V \rightarrow$ You need V_T and P

$$V_T(s) = \sum_{\pi \in \Pi} \pi(G|s) \sum_{s' \in S, r \in R} p(s', r, s, a) [r + \gamma V_T(s')]$$

Next Lecture:

How to obtain V and Q given P

Next Next lecture

How to obtain V and Q if we do not have P .

S : State Space \rightarrow Surface or abstract state.

A : Action Space \rightarrow Click, Impression, Repin...

g : Transition rule function \rightarrow likelihood of going from one state to the other.

R : Reward function

$\left\{ \begin{array}{l} \text{- Assign points for things you know} \\ \text{- Initialize to a random value for things you don't know} \end{array} \right.$

π : Policy.

γ : Discount factor.

\hookrightarrow How much you care about the future.

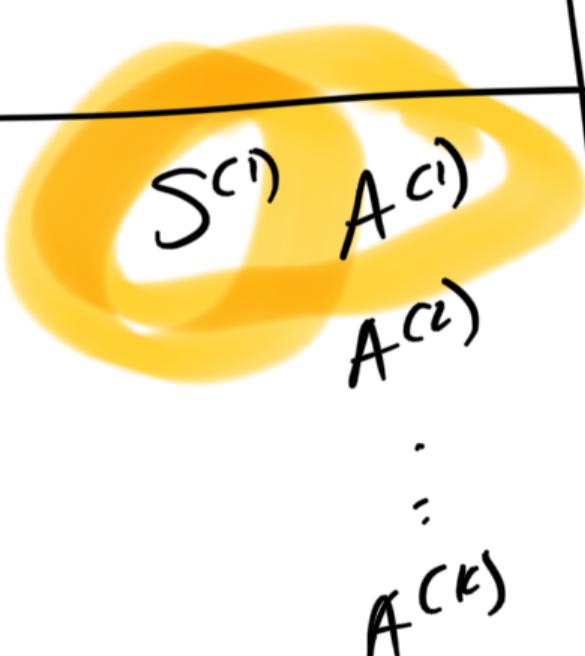
ACTION VALUE FUNCTION (one hop)

Next' states

Next state
Reward

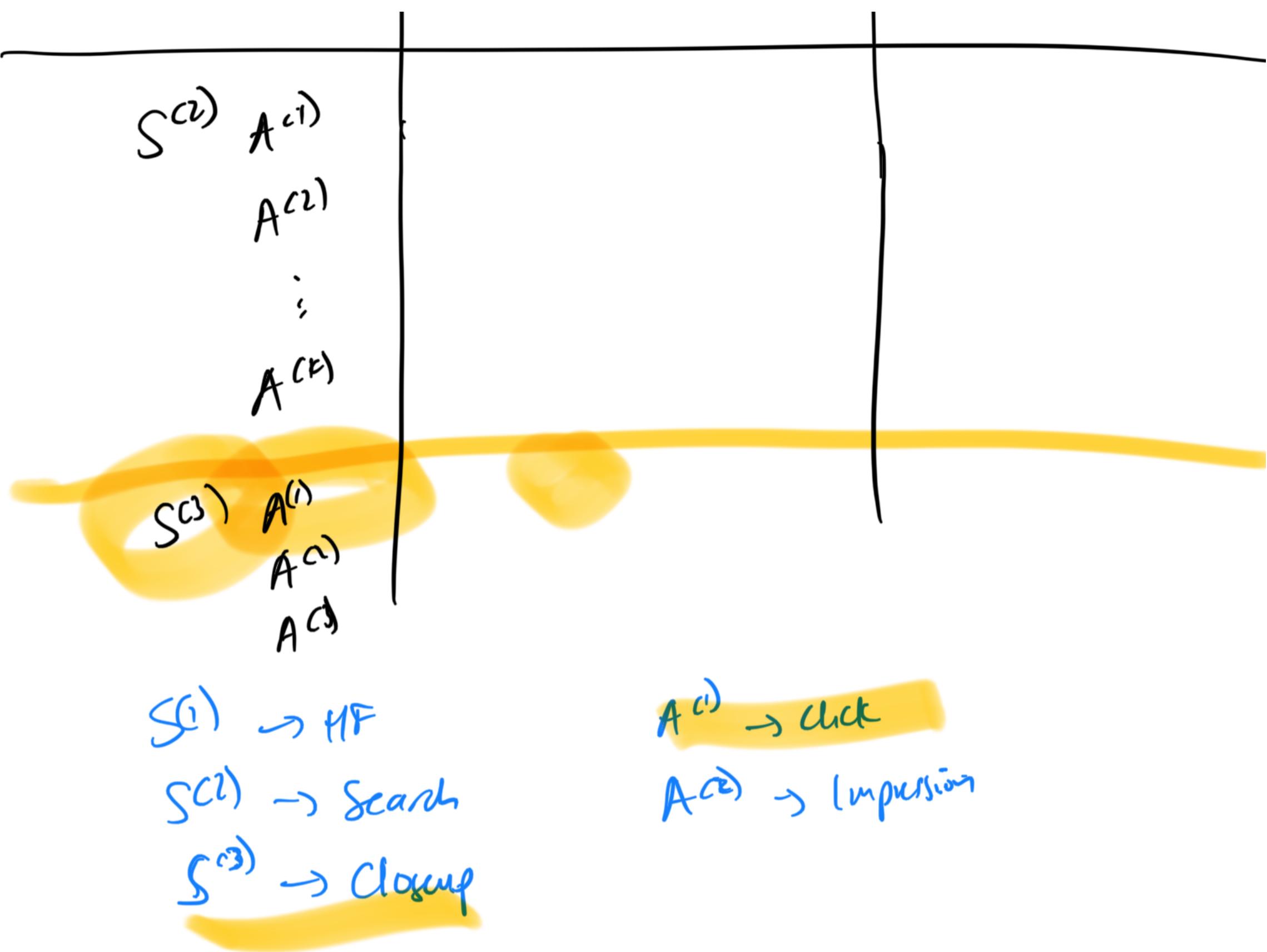
$s^{(1)} s^{(2)} \dots s^{(n)}$

$r^{(1)} r^{(2)} \dots r^{(n)}$



$$\sum = 1$$

$$\Sigma = 1$$



RETURNS AND EPISODES

The Agent's goal is to maximize the cumulative reward it receives over time.

Downstream Rewards

@ Nutella -

Instant gratification \Rightarrow Long term objective optimization



Replies.



optimize for
best sequence
(long term)

Markov Decision Processes.

$A: \{A_1, A_2, \dots, A_n\} \rightarrow \text{Action space}$

click
Repin
Impression ...

$S: \{S_1, S_2, \dots, S_k\} \rightarrow \text{State space}$

Abstract state \rightarrow Exploration, consideration,
Fulfillment.

Surfaces.

q : Transition rate function.

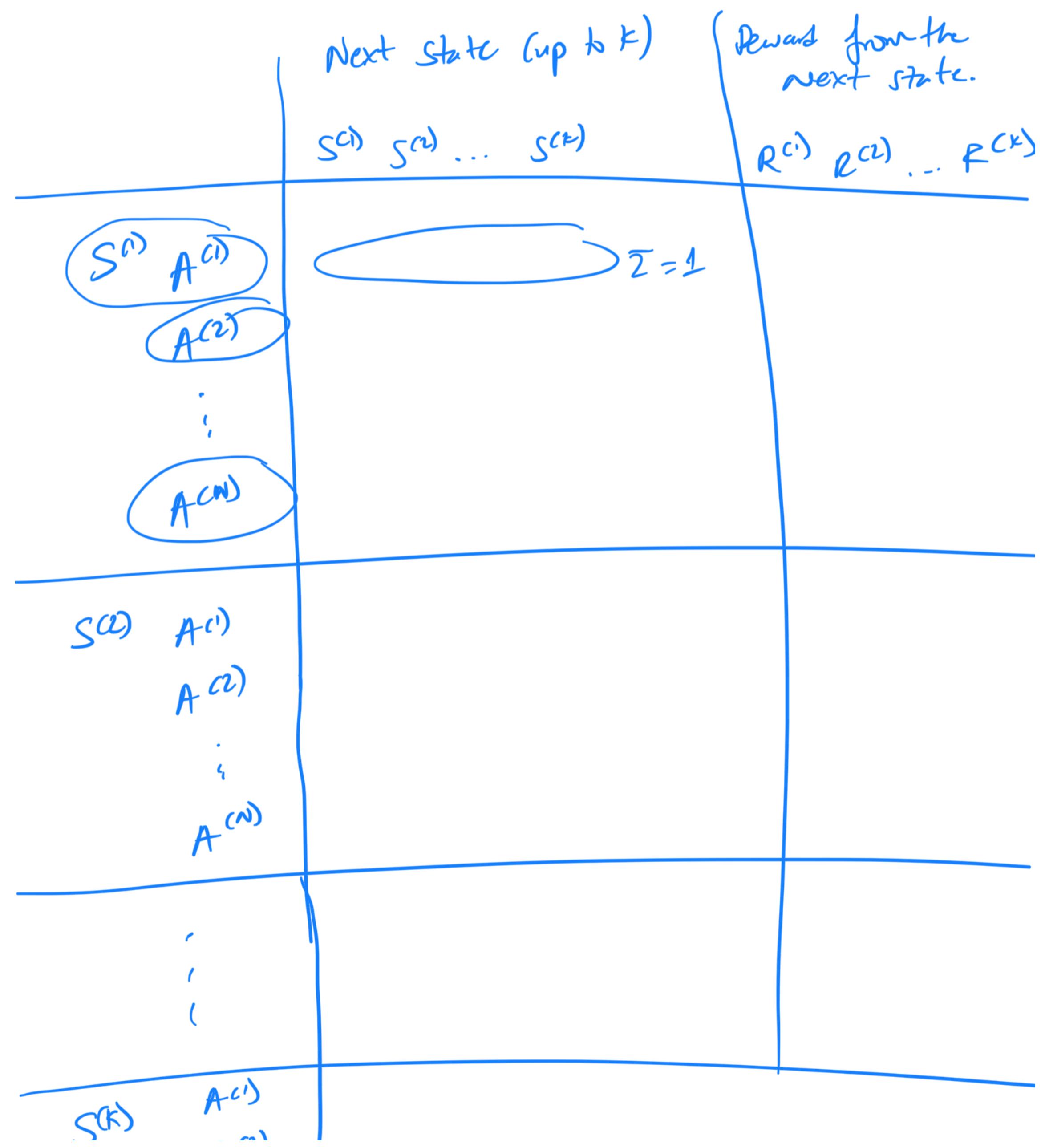
R : Reward function. \rightarrow Relative value of different Actions and states.

\rightarrow $\Pi : \text{Bilz}$

→ Hide → -100
→ Share → 100

γ : Discount factor.

↳ How much you care about the future.



$A^{(1)}$
:
 $A^{(N)}$

Choosing upto 3-5 Actions

Starting with one surface, i.e. Home feed.

$S^{(1)}$: Home feed
 $S^{(2)}$: Search
 $S^{(3)}$: Related Pins

{
- User ID
- time stamp begin
- time stamp end
- action (e.g., repin, click, etc)
- Surface origin
- Surface Destination
- Pin ID.

{
- web / mobile
- Country
- Date joined
- Gender
- Age
- L1 interest
- iOS / Android
- Version
} TD normalize population

$S \times A$

$S^{(1)}$	$A^{(1)}$	#	$\mathbb{E} [G_i \mid S_i = S^{(1)}, A_i = A^{(1)}]$
$S^{(2)}$	$A^{(1)}$	#	
$S^{(2)}$	$A^{(2)}$		
$S^{(3)}$	$A^{(1)}$		
$S^{(3)}$	$A^{(2)}$		
$S^{(4)}$	$A^{(1)}$		
$S^{(4)}$	$A^{(2)}$		
$S^{(5)}$	$A^{(1)}$		
$S^{(5)}$	$A^{(2)}$		
$S^{(6)}$	$A^{(1)}$		
$S^{(6)}$	$A^{(2)}$		
$S^{(7)}$	$A^{(1)}$		
$S^{(7)}$	$A^{(2)}$		
$S^{(8)}$	$A^{(1)}$		
$S^{(8)}$	$A^{(2)}$		
$S^{(9)}$	$A^{(1)}$		
$S^{(9)}$	$A^{(2)}$		
$S^{(10)}$	$A^{(1)}$		
$S^{(10)}$	$A^{(2)}$		