

A high level idea about the RL setting:

Task: Policy evaluation, given a policy π
in your favorite MDP (SSP, discounted, etc)

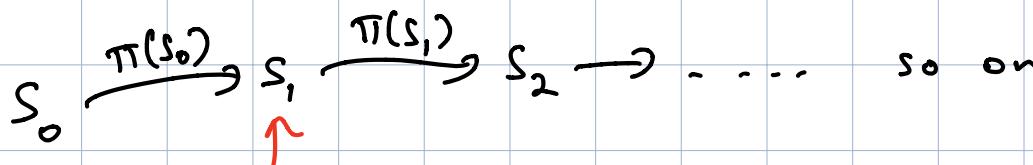
If we know the transition probabilities, then
we know T_π & hence, can solve

$$J_\pi = T_\pi J_\pi \text{ directly}$$

or start with some J & keep applying T_π .

In RL setting: $P_{ij}(\alpha)$ is not known.

Instead, we get to see a sample path.



Sampled from

$$P_{s_0 s_1}(\pi(s_0))$$

$$T^\pi J = \sum_j P_{ij}(\pi(i)) (g(i, \pi(i), j) + J(j))$$

This ain't known, but we get to see
a sample path that follows $P_{ij}(\alpha)$.

Assume cost function is known & is a deterministic function.
For simplicity, we may assume $g(\cdot, \cdot)$ is a function
of the current state & current action.

The goal is to solve MDP as before,
but the information available is only through
the sample path.

So, an element of "learning" involved.

Policy evaluation! Temporal difference (TD) learning
Control: Q-learning

Background on Stochastic iterative algorithms

Ref: Chapter 4 of TD book.

Aim: Solve the system of equations

$H\tau = \tau$, where $H: \mathbb{R}^n \rightarrow \mathbb{R}^n, \tau \in \mathbb{R}^n$

in RL applications, think of
 H as the Bellman operator T or
 T^π & τ as the value function

A direct approach: Start with τ_0 & do $\tau_{n+1} = H\tau_n$

a variation to the above:

$$r_{n+1} = (1 - \beta) r_n + \beta H r_n$$

\uparrow \nearrow
Iterative scheme Step size $0 < \beta < 1$

An interesting special case:

$$Hr = r - \nabla f(r) \quad \text{for some } f: \mathbb{R}^n \rightarrow \mathbb{R}$$

Then $Hr = r \quad (\Rightarrow) \quad \nabla f(r) = 0$

problem of finding the minima of f .

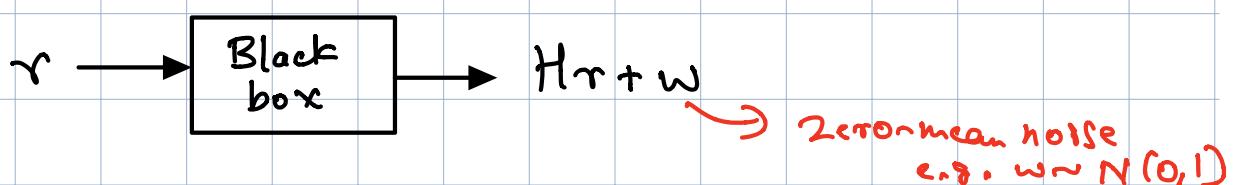
$$r_{n+1} = (1 - \beta) r_n + \beta (r_n - \nabla f(r_n))$$

(or) $r_{n+1} = r_n - \beta \nabla f(r_n)$ ← Gradient descent algorithm

Remark: If $r_n \rightarrow r^*$ and H is continuous at r^* ,
then $Hr^* = r^*$.

So far, we assumed Hr is perfectly observable $\forall r$

Now, we consider a setting where H is not
precisely known.



$$r_{m+1} = (1 - \beta) r_m + \beta \underbrace{(H r_m + w_m)}_{\text{noisy observation}} \quad (*)$$

Simplicity: assume $E(w_i) = 0, Ew_i^2 < \infty, w_i$ i.i.d.
independent & identically distributed.

(*) is a stochastic iterative scheme.

Want: $r_m \rightarrow r^*$ with probability (w.p.) 1
as $m \rightarrow \infty$

same as
almost surely
(a.s.)

We shall identify conditions on

(i) H

e.g. contraction

(ii) step sizes

e.g. constant β or diminishing
 $\beta_m \approx k_m$

(iii) noise w_m

e.g. $Ew_i = 0, Ew_i^2 < \infty$
& $\{w_i\}$ is iid.

so that we can infer $r_m \rightarrow r^*$ w.p. 1 as $m \rightarrow \infty$

An application: Mean estimation

Consider a random variable (r.v.) X with mean μ
& finite variance, say σ^2 .

Suppose we are given iid samples X_1, \dots, X_m

Let r_m be the estimate of μ .

$$r_m = \frac{1}{m} \sum_{k=1}^m x_k$$

$$r_{m+1} = \frac{1}{m+1} \sum_{k=1}^{m+1} x_k$$

$$= \frac{m}{m+1} \left(\frac{1}{m} \sum_{k=1}^m x_k \right) + \frac{1}{m+1} x_{m+1}$$

$$r_{m+1} = \frac{m}{m+1} r_m + \frac{1}{m+1} x_{m+1}$$

← iterative scheme
for updating sample
mean

$$r_{m+1} = r_m + \frac{1}{m+1} (x_{m+1} - r_m)$$

$$r_{m+1} = r_m + \beta_m (x_{m+1} - r_m) \quad - (\text{xx})$$

where

$$\beta_m = \frac{1}{m+1}$$

(xx) resembles the sto. iter.algo (x)

What does Strong Law of Large numbers say about r_m ?

$r_m \rightarrow \mu$ a.s. as $m \rightarrow \infty$

with step size $\beta_m = \frac{1}{m+1}$

$$r_{m+1} = r_m + \beta_m (x_{m+1} - r_m)$$

$$= r_m + \beta_m ((\mu - r_m) + (x_{m+1} - \mu))$$

$$r_{m+1} = (1 - \beta_m) r_m + \beta_m \left(\underbrace{\mu}_{Hr_m} + \underbrace{(x_{m+1} - \mu)}_{w_m} \right)$$

$$\mathbb{E} w_m = 0, \quad \mathbb{E} w_m^2 < \infty, \quad \text{if } r^* = \bar{r} (\Leftrightarrow) \bar{r} = \mu.$$

So, (**) is really a sto. itkrative algorithm.

And, we get $r_m \rightarrow r^* (= \mu)$ a.s. as $m \rightarrow \infty$.

The theory of sto. itkr. algo that we shall develop ensures $r_m \rightarrow \mu$ a.s. as $m \rightarrow \infty$ for more general step sizes that satisfy

$$\sum_m \beta_m = \infty, \quad \sum_m \beta_m^2 < \infty$$

Question: Why do we need these conditions?

On step-size requirements:

$$r_{m+1} = (1 - \beta_m) r_m + \beta_m (H r_m + w_m)$$

$$r_{m+1} = r_m + \beta_m (H r_m + w_m - r_m)$$

increment

Suppose $\{w_i\}$ is independent of r_m, H_m

& has variance σ^2

Variance of r_{m+1}

$$= \text{Var}(r_{m+1})$$

$$= \text{Var}\left[(1 - \beta_m) r_m + \beta_m H r_m\right] + \beta_m^2 \text{Var}(w_m)$$

$$= \text{Var}\left[(1 - \beta_m) r_m + \beta_m H r_m\right] + \beta_m^2 \sigma^2$$

$$\geq \beta_m^2 \sigma^2$$

Now, if $\beta_m = \beta$, then

$$\text{Var}(r_{m+1}) \geq \beta^2 \sigma^2$$

So, fixing $\beta > 0$, $r_m \rightarrow r^*$ for any r^* .

So, the step-size has to vanish asymptotically.

But, β_m cannot go down too fast.

$$r_{m+1} = r_m + \beta_m (Hr_m + w_m - r_m)$$

$$|r_m - r_0| \leq \sum_{\tau=0}^{m-1} \beta_\tau |Hr_\tau - r_\tau + w_\tau|$$

So, if $|Hr_\tau - r_\tau + w_\tau| \leq C_1$, and

if $\sum_{\tau=0}^{\infty} \beta_\tau \leq C_2 < \infty$, then

$|r_m - r_0|$ is bounded above

(\Rightarrow) r_m is within a certain radius of r_0

problematic if r^* lies outside the radius.

So, we need $\sum_{\tau} \beta_\tau = \infty$

Lecture-19*

Notions of convergence of r.v.s:

I Almost sure or w.p. 1 Convergence!

$\{X_m\}$ r.v.s.

$X_m \rightarrow X$ a.s. or $X_m \rightarrow X$ w.p. 1 as $n \rightarrow \infty$

if $P \left(\omega \mid \lim_{m \rightarrow \infty} X_m(\omega) = X(\omega) \right) = 1$

e.g. SLLN

II Convergence in probability

$X_m \xrightarrow{P} X$ if $\lim_{m \rightarrow \infty} P(|X_m(\omega) - X(\omega)| > \epsilon) = 0, \forall \epsilon > 0.$

usually written as

$$P(|X_m - X| > \epsilon) \rightarrow 0$$

e.g. WLLN

III L^2 -Convergence

$X_m \xrightarrow{L^2} X$ if $E |X_m - X|^2 \rightarrow 0$ as $m \rightarrow \infty$.

IV Convergence in distribution

$X_m \xrightarrow{d} X$ if $E f(X_m) \rightarrow E f(X)$

f bdd continuous f.

(Check textbook for equivalent definitions). e.g. CLT.

Note: a.s. $\xrightarrow{\text{conv. in}} \text{prob} \xrightarrow{\text{conv. in}} \text{dist.}$

In this lecture, we provide a.s. convergence guarantees for the popular RL algorithms, e.g., TD-learning & Q-learning.

A crash course in selected topics in prob

Def: A collection \mathcal{F} of subsets of Ω is a **σ -field** if

- (a) $\emptyset \in \mathcal{F}$
- (b) $A_1, A_2, \dots \in \mathcal{F} \Rightarrow \bigcup_{i=1}^{\infty} A_i \in \mathcal{F}$
- (c) $A \in \mathcal{F} \Rightarrow A^c \in \mathcal{F}$

Def: A probability measure P on

(Ω, \mathcal{F}) is a function

$P: \mathcal{F} \rightarrow [0, 1]$ satisfying

(a) $P(\Omega) = 1$

(b) A_1, A_2, \dots disjoint, i.e., $A_i \cap A_j = \emptyset$ $\forall i \neq j$

then, $P\left(\bigcup_{i=1}^{\infty} A_i\right) = \sum_{i=1}^{\infty} P(A_i)$

Conditional expectation! Check CS601S course notes, or
any standard probability textbook.

(e.g. Grimmett - Stirzaker,
Prob & random processes)

σ -field generated by a r.v. X , denoted by $\sigma(X)$:

$\sigma(X)$ is the smallest σ -field containing all sets

of the form $\{X \leq \alpha\} = \{\omega \mid X(\omega) \leq \alpha\}$

\downarrow
Take all such sets, add complements
& countable unions to get
 $\sigma(X)$.

Given $\sigma(X)$, $E(Y \mid \sigma(X))$ is denoted by $E(Y|X)$

Extend $\sigma(X)$ to the σ -field generated
by a collection of r.v.s $\{X_1, \dots, X_m\}$ &
denote it by $\sigma(X_1, \dots, X_m)$

Also, $E(Y \mid \sigma(X_1, \dots, X_m)) = E(Y \mid X_1, \dots, X_m)$

Filtration: increasing sequence of σ -fields.

e.g. $\mathcal{F}_1 = \sigma(X_1)$, $\mathcal{F}_2 = \sigma(X_1, X_2)$, and so on

$\mathcal{F}_1 \subseteq \mathcal{F}_2 \subseteq \mathcal{F}_3 \dots$ $\{\mathcal{F}_k\}$ is a filtration

$\mathcal{F}_k = \sigma(X_1, \dots, X_k) \leftarrow$ filtration
 (X_k, \mathcal{F}_k) is a martingale if

$$E(X_{k+1} | \mathcal{F}_k) = X_k, \forall k$$

Note! If $E(X_{k+1} | \mathcal{F}_k) = 0$ then (X_k, \mathcal{F}_k) is a martingale difference sequence.

(*) $\{X_k\}$ are not necessarily i.i.d.

To see martingale difference, contraction mapping, etc in a stoc. iter. algo in RL context,
Let's look at TD(0) learning.

Fix policy π . Consider a discounted MDP.

Also, assume single stage cost $g(i, \pi(i))$ & is indep. of next state

$$(H\pi)(i) = g(i, \pi(i)) + \gamma \sum_j p_{ij}(\pi(i)) \pi(j)$$

Want to solve $\pi(i) = (H\pi)(i) \quad \forall i - (\infty)$

Sto. Ikr. algo for solving (*)

$$J_{t+1}(i) = J_t(i) + \beta_t \left(g(i, \pi(i)) + \alpha J_t(\bar{i}) - J_t(i) \right)$$

\bar{i} is a state observed along the sample path, i.e., in state i , take action $\pi(i)$, & we transition to \bar{i} & \bar{i} follows $P_{i|\bar{i}}(\pi(i))$

$$= J_t(i) + \beta_t \left(g(i, \pi(i)) + \alpha \sum_j P_{i|\bar{i}}(\pi(i)) J_t(j) - J_t(i) \right)$$

$$+ \left\{ g(i, \pi(i)) + \alpha J_t(\bar{i}) - g(i, \pi(i)) - \alpha \sum_j P_{i|\bar{i}}(\pi(i)) J_t(j) \right\}$$

↑ added & subtracted this

$$J_{t+1}(i) = J_t(i) + \beta_t \left(((H J_t)(i) - J_t(i)) + \omega_t(i) \right)$$

where $\omega_t(i)$ is the quantity in flower braces.

In the absence of noise term $\omega_t(i)$, it is easy to infer

$$J_t \rightarrow J_\pi, \text{ where } J_\pi = H J_\pi$$

What conditions on $w_t(i)$ would ensure the above convergence?

Let $\mathcal{F}_t = \sigma(\mathcal{T}_0, \dots, \mathcal{T}_t, w_0, \dots, w_{t-1})$

$$E(w_t(i) | \mathcal{F}_t) = 0$$

We shall get back to TD(0) later.

Coming next:

Sto. iter. algo. $r_{m+1}(i) = r_m(i) + \beta_m (Hr_m(i) + w_m(i) - r_m(i))$

Then,

if ① H is a contraction mapping wrt max-norm
 $Hr^* = r^*$

$$\textcircled{2} \quad E(w_m(i) | \mathcal{F}_m) = 0$$

$$E(w_m^2(i) | \mathcal{F}_m) \leq \text{"(bounded)"}$$

$$\textcircled{3} \quad \sum \beta_m = \infty, \quad \sum \beta_m^2 < \infty$$

Then, $r_m \rightarrow r^*$ a.s. as $m \rightarrow \infty$.

Lecture-20

Convergence of stochastic iterative algorithms

Ref: Sec 4.3 of NOP book

$$r_{m+1}(i) = (1 - \beta_m) r_m(i) + \beta_m (H r_m(i) + \omega_m(i)), \quad (1)$$

$i = 1, \dots, n$

If helps to look at per-component update because RL applications would estimate the value function, say $\mathcal{J}_\pi(i)$, with start state i using an iterate of the form $r_{m+1}(i)$

The underlying σ -field is

$$\mathcal{F}_m = \sigma(r_0(i), \dots, r_m(i), \omega_0(i), \dots, \omega_{m-1}(i), i=1, \dots, n)$$

all the random variables upto time m

Assumptions:

<conditions on noise factors $\{\omega_m\}$ >

$$(A1) \quad H_i \text{ and } H_m, (i) \quad E(\omega_m(i) | \mathcal{F}_m) = 0$$

conditionally
zero mean

$$(ii) \quad E(\omega_m^2(i) | \mathcal{F}_m) \leq A + B \|r_m\|^2$$

bound on the conditional variance

Note: For mean-estimation example,

$w_m \rightarrow$ iid, indep. of $\{\gamma_m\}$, bounded variance

So, $E(w_m | \mathcal{F}_m) = E(w_m) = 0$

$$E(w_m^2 | \mathcal{F}_m) = E(w_m^2) \leq (\text{const})^2$$

↓
since $\text{Var}(w_m)$ is bounded.

i.e., $B=0$ in (A1) part(i)

(A2) Condition on H

Define weighted norm

$$\|\boldsymbol{x}\|_{\xi} = \max_i \frac{|x(i)|}{\xi(i)}$$

→ we saw an example of such a norm in SSP

If $\xi(i) \equiv 1 \forall i$, then we get the max-norm. → see discounted MDP

Assume H is a "weighted max-norm pseudo-contraction",
i.e., \exists a positive $\xi = (\xi(1), \dots, \xi(n))$ and
a constant $\beta \in [0, 1)$ s.t.

$$(*) \rightarrow \|Hr - r^*\|_{\xi} \leq \beta \|\boldsymbol{x} - \boldsymbol{x}^*\|_{\xi}, \quad \forall r \in \mathbb{R}^n.$$

$$(\Rightarrow) \quad \frac{|Hr(i) - r^*(i)|}{\xi(i)} \leq \beta \max_j \frac{|x(j) - x^*(j)|}{\xi(j)}, \quad \forall i, r.$$

[Recall \rightarrow full contraction if
 $\|Hr - Hr'\|_S \leq \beta \|r - r'\|_S$]

(claim w/o proof: $(*) \Rightarrow H\hat{r}^* = \hat{r}^*$ & \hat{r}^* is the unique fixed point.)

A full-contraction is trivially a pseudo-contraction.

Connection to MDPs:

- ① In a SSP with all proper policies,
the Bellman operator is a contraction.
 - ② In a discounted MDP with bounded single stage cost,
the Bellman operator is a contraction.
-

(A3) "Step size conditions"

$$\sum_m \beta_m = \infty, \quad \sum_m \beta_m^2 < \infty$$

Theorem 1 "Convergence of Sto. iter. algo"

Assume (A1), (A2), (A3). For r_m updated
using Eq.(1), we have $r_m \rightarrow \hat{r}^*$ a.s. as $m \rightarrow \infty$

A rough intuitive argument:

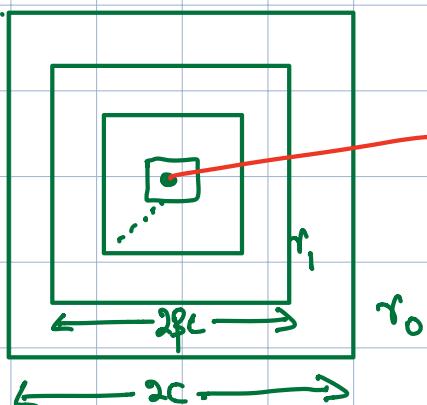
Take a very special case where

$H_0 = 0$ i.e., $r^* = 0$, H is - contraction with modulus β

$$w_m = 0 \quad \forall m$$

$$r_{m+1}(i) = H r_m(i)$$

Suppose the initial point $|r_0(i)| \leq c$



In the sto. iter-algo (Eq ①),

$$r_{m+1}(i) = (1 - \beta_m) r_m(i) + \beta_m (H r_m(i) + w_m(i))$$

noise

Because of noise, the shrinkage is not immediate, but the iterate r_m gets into smaller & smaller box as the algorithm proceeds.

A convergence result under monotonicity

Motivation: SSP where \exists at least one proper policy & improper policies have infinite cost. In such SSPs, the Bellman operator is not "contractive", but we still have "monotonicity".

Instead of (A2), assume the following:

(A2') (i) H is monotone, i.e., $\gamma \leq \gamma' \Rightarrow H\gamma \leq H\gamma'$, $\forall \gamma, \gamma'$.

(ii) $\exists \gamma^*$ s.t. $H\gamma^* = \gamma^*$ & γ^* is unique.

(iii) e = vector of all ones, $e \in \mathbb{R}^n$.
 $\delta > 0$

$$H\gamma - \delta e \leq H(\gamma - \delta e) \leq H(\gamma + \delta e) \leq H\gamma + \delta e$$

Theorem 2: Assume (A1), (A2') & (A3)
 Consider the gto. ikr. algo that updates τ_m using Eq(1).

Suppose τ_m is bounded w.p. 1 ie., $\sup_{m,i} |\tau_m(i)| < \infty$

Then, $\gamma_m \rightarrow \gamma^*$ a.s. as $m \rightarrow \infty$

Note!: Unlike Thm 1, for the monotone case,
 we require the iterate r_m to satisfy a
 "boundedness" requirement.

A sufficient condition to ensure boundedness of the iterates:

The conditions are

(i) $\{\omega_m(i)\}$ satisfies (A1)

(ii) $\{\beta_m\}$ satisfies (A3)

(iii) \exists a positive vector ζ , a $\beta \in [0,1]$ and
 a scalar $D > 0$ s.t.

$$\|Hr_m\|_\zeta \leq \beta \|r_m\|_\zeta + D, \quad \forall m$$

Under (i), (ii), (iii), r_m is bounded w.p. 1.

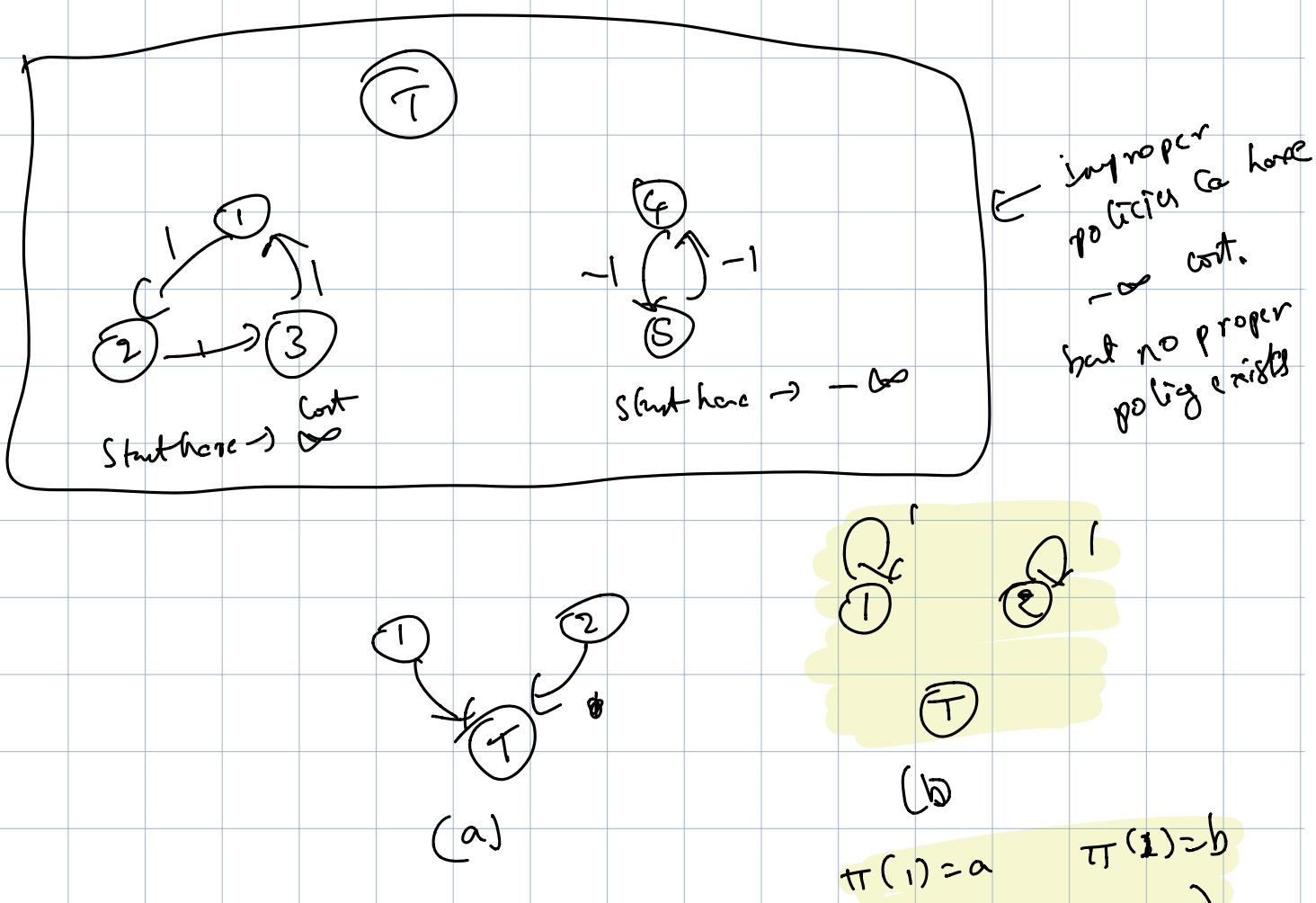
Remark!: If H is a pseudo-contraction, then

$$\begin{aligned} \|Hr_t\|_\zeta &= \|Hr_t - r^* + r^*\|_\zeta \leq \|Hr_t - r^*\|_\zeta + \|r^*\|_\zeta \\ &\leq \beta \|r_t - r^*\|_\zeta + \|r^*\|_\zeta \\ &\leq \beta \|r_t\|_\zeta + (\text{if } \beta) \|r^*\|_\zeta \end{aligned}$$

So, a pseudo-contractive Π sat's fics condit.(iii)
 & hence, the iterate is bounded in the
 pseudo-contractive case.

Summary

$\Sigma \text{ of } \Pi$ converges
 Π is Contractive
 Π is monotone + iterate bounded.



$$J(\pi) = (-\infty, \infty)$$

π

improper policy

but no state without