

Lecture #17

Proof: Refer to Ch. 8 of the Book I sent via e-mail.

Why bother about contraction maps?

Proposition: Let F be the class of all functions

$z : \{1, \dots, n\} \mapsto \mathbb{R}$. Define $\|z\|_\infty := \max_{1 \leq i \leq n} |z(i)|$ and the nonlinear map $T : F \mapsto F$ as

$$(Tz)(i) := \inf_{u \in U} \left\{ c(i, u) + \beta \sum_{j=1}^n z(j) P_{ij}(u) \right\}$$

Then $T(\cdot)$ is a contraction.

Proof. P. 148
in Ch. 8 of the book I sent.

Recap: Bellman's eq^{**} / DP eq^{**} from
Infinite horizon discounted cost

OCP: $\min_{\underline{x} \in \Gamma} \mathbb{E}_w \left[\sum_{k=0}^{+\infty} \beta^k c(\underline{x}_k, \underline{u}_k, \underline{w}_k) \right]$

assume \underline{w}_k iid.

s.t.

$$\underline{x}_{k+1} = f_k(\underline{x}_k, \underline{u}_k, \underline{w}_k).$$

DP Eq $\hat{=}$: (nonlinear algebraic eq $\hat{=}$)

$$W_\infty(\underline{x}) = \min_{u \in U} \mathbb{E} \left[c(\underline{x}, \underline{u}, \underline{w}) + \beta W_\infty(f(\underline{x}, \underline{u}, \underline{w})) \right]$$

Special case of Bellman's eq $\hat{=}$

Instead of $\underline{x}_{k+1} = f_k(\underline{x}_k, \underline{u}_k, \underline{w}_k),$

we have $\underline{x}_k \sim \text{Markov}(P(u)),$

$$P(u) = [P_{ij}^{(u)}]_{i,j=1,\dots,n}$$

Then Bellman's eq^{**} : / DP eq^{**}:

$$W_\infty(i) = \min_{u \in U} \left[C(i, u) + \beta \sum_{j=1}^n p_{ij}(u) W_\infty(j) \right]$$

$$i = 1, \dots, n$$



The proposition says the RDS, thought of as a nonlinear operator $T \underline{W}_\infty$, is a contraction.

Theorem :

(1) There exists unique sol^{un} $(W_\alpha(1), W_\alpha(2), \dots, W_\alpha(n))^\top$ to the set of nonlinear eq^{un}s:

$$W_\alpha(i) = \inf_{u \in U} \left\{ C(i, u) + \beta \sum_{j=1}^n P_{ij}(u) W_\alpha(j) \right\}, \quad 1 \leq i \leq n.$$

(2) The answer/unique sol^{un} from (1) satisfies:

$$W_\alpha(i) = \inf_{\gamma \in \Gamma} \mathbb{E}^{\gamma} \left\{ \sum_{k=0}^{+\infty} \beta^k C(x_k, u_k) \mid x_0 = i \right\}$$


cost - f₀ - g₀

(3) The map T, defined as

$$T_k(i) = \inf_{u \in U} \left\{ C(i, u) + \beta \sum_{j=1}^n P_{ij}(u) z(j) \right\}, \quad i=1, \dots, n$$

is a contraction w.r.t. norm $\|z\| = \max_{1 \leq i \leq n} |z_i|$

$$(4) \lim_{n \rightarrow \infty} (T^n z)(i) = W_0(i) \quad \forall i = 1, \dots, n$$

for any z .

$$(5) \text{ If } z(i) = \underbrace{0}_{\pi} + \varepsilon, \text{ then } (T^n z)(i) = W_n(i),$$

$$\underline{z} = \text{zeros}(n, 1)$$

$$\text{where } W_n(i) := \inf_{\gamma \in \Gamma} \mathbb{E}^{\gamma} \left\{ \sum_{k=0}^{n-1} \beta^k C(x_k, u_k) \middle| x_0 = i \right\}.$$

- Everything we said so far holds even if \mathcal{X} (state space) is countable & \mathcal{U} is compact.

Next Main result:

We say a Markov policy is stationary time-invariant if

$$\{\underline{\delta}_0, \underline{\delta}_1, \dots\} = \{\underline{\delta}, \underline{\delta}, \underline{\delta}, \dots\}$$

Infinite sequence
of functions

Result: Stationary policy is optimal for Infinite horizon discounted.

Algorithms.

Value Iteration

- Good: • Guaranteed convergence
 • Rate of convergence is geometric

Bad(?): Only asymptotic convergence

(Simply implement the contractive iterates)

e.g. choose arbitrary
 $W_0(i) = \text{rand}(n, 1)$
 then iterate $W_{\infty}^{k+1} = TW_{\infty}^k$
 $k = 0, 1, \dots$

Policy Iteration

Finite step convergence.

Idea: $T(\cdot)$ is not only contractive, but also a monotone (nonlinear) operator

(i.e.) If $\underline{\pi} \leq \tilde{\pi}$,

then $T\underline{\pi} \leq T\tilde{\pi}$

This suggests we can generate approx. optimal policy $\underline{\pi}$ from a finite set.

Policy Iteration Algorithm:

(Instead of iterating value f^{π} , let us iterate policy)

→ Start with any policy $\{\underline{\gamma}_0, \underline{\gamma}_0, \underline{\gamma}_0, \dots\}$
with $\underline{\gamma}_0$ arbitrary

→ Calculate the cost $W_{\underline{\gamma}_0}$, i.e.

$$W_{\underline{\gamma}_0} = (I - \beta P_{\underline{\gamma}_0})^{-1} c_{\underline{\gamma}_0}$$

→ Is $W_{\underline{\gamma}_0} = \overbrace{TW_{\underline{\gamma}_0}}^{\text{?}}$?

$$\hookrightarrow := \inf_{u \in U} [c(i, u) + \beta \sum_{j=1}^n p_{ij}(u) W_{\underline{\gamma}_0}(j)]$$

If yes, then ~~stop~~.

If not, let

$$\gamma_1(i) = \underset{u(\cdot)}{\operatorname{arg\,min}} \text{ (above RHS)}$$

Again calculate

$$w_{\gamma_1} = (I - \beta P_{\gamma_1})^{-1} c_{\gamma_1} \dots \text{ continue.}$$

Another alternative : (Linear Programming solution)

$$\max \sum_{i=1}^n z(i)$$

subject to

$$z(i) \leq c(i, u) + \sum_{j=1}^m \beta p_{ij}(u) z(j) \quad \forall i \quad \forall u$$

decision variable

$$\# \text{ of constraints} = (\# \text{ of states}) \times (\# \text{ of controls})$$

The LP constraint \Leftrightarrow

$$z \leq Tz \leq T^2 z \leq T^3 z \dots \rightarrow w_\infty.$$

(i.e.) Sol^m of LP is w_∞ .

So far, assumed that \underline{z} is available for feedback, i.e., (completely observed) Markov decision process (MDP).

More realistic case, \underline{x} is NOT directly observable.

What if we can only observe $y \neq x$.

POMDP

(Partially Observed Markov Decision Process)

System:

$$p_{ij}(u), i \in \mathcal{X}, u \in \mathcal{U}$$

Noisy observation of state:

\mathcal{Y} = set of observations

$$p(y|x) = P(y(t) = y | x(t) = x)$$

History dependent policies:

$$u(t) = \delta_t(y(0), y(1), \dots, y(t), u(0), u(1), \dots, u(t-1))$$

$$\underline{\delta} = (\delta_0, \delta_1, \dots, \delta_{T-1}) \leftarrow \text{policy}$$

The main result (in this Markov Chain setting):
we will show that we can solve this POMDP
as a 2-step procedure by separating
estimation & control.

(Whenever these 2 problems can be
decoupled, we say "Separation Principle" holds)

State estimation for a Markov Chain:

p_{ij} : Markov chain (no control)

$p(y|x)$: observation probabilities.

Question: What is $x(t)$?

We want

$P(x(t)=i | y(0), y(1), \dots, y(t))$ for each $i=1, \dots, h$

Define $y^t := \underbrace{(y(0), \dots, y(t))}_{\text{History upto time } t}$

Denote:

$$p_{t|t}(i | y^t) := P(x(t) = i | y^t)$$

Define vector:

$$\begin{aligned} p_{t|t}(y^t) &= [p_{t|t}(1 | y^t), p_{t|t}(2 | y^t), \dots, p_{t|t}(N | y^t)] \\ &= \text{conditional probability distribution} \\ &\quad \text{of } x(t) \text{ given } y^t \end{aligned}$$

We will show : (Recursive sol[†])

