

Reinforcement Learning Notes

Varad Vaidya

March 2, 2024

Lecture notes from the various YouTube playlist related to Reinforcement Learning, combined with the course E1277 Reinforcement Learning at IISc Bangalore by Prof. Gagan Thoppe. The plan is to merge the notes from david silver's course and the course at IISc Bangalore, at appropriate places, to make a single set of notes.

Since the notes are being merged from two different sources, proper crediting of the two (and many more) is hard. In general, the notes will follow, from the standard books of Reinforcement Learning, and are my interpretation of the same.

Disclaimer: This document will inevitably contain some mistakes— both simple typos and legitimate errors. Keep in mind that these are the notes of a graduate student in the process of learning the material, so take what you read with a grain of salt. If you find mistakes and feel like telling me, I will be grateful and happy to hear from you, even for the most trivial of errors. You can reach me by email at vaidyavarad2001@gmail.com.

For more notes like this, visit [varadVaidya](https://varadVaidya.com).

Varad Vaidya,
Fall Term: 2023,
Last Update: March 2, 2024,

Contents

I	E1277 — Reinforcement Learning	1
1	Markov Chains	2
1.1	Markov Processes or Markov Chains	2
1.2	Communicating Classes	7
1.3	Strong Markov Property	8
1.4	Recurrence and Transience	8
1.5	Invariant Distributions	12
2	Markov Decision Processes and Dynamic Programming	18
2.1	Controllrd Mrkov Chain	18
2.2	Markov Decision Process	18
2.3	Finite Horizon Problems	19
2.4	Stochastic Shortest Path	21

Part I

E1277 — Reinforcement Learning

1 Markov Chains

1.1 Markov Processes or Markov Chains

A Markov chain is simply a Markov Decision Process without decision. It is one the most simplest stochastic process, and has no “memory” of the past. So, just the present state determines its future dynamics. In this context, we will be considering Discrete Time Markov Chains (DTMCs).

DTMC involves two concepts:

- Discrete Time Stochastic Process (DTSP)
- Row Stochastic Matrix.

The two are defined as follows:

Definition 1.1 (Discrete Time Stochastic Process (DTSP)). A DTSP is a sequence $(X_n)_{n \geq 0}$ of random variables defined on the same probability space $(\Omega, \mathcal{F}, \mathbb{P})$, taking values in the same set \mathcal{S} , i.e.

$$X_n : \Omega \rightarrow \mathcal{S} \quad \forall n \geq 0$$

where \mathcal{S} is called a state space, and is finite or countably infinite. We call the a variable as a “state” when the state is an element of the state space \mathcal{S} . The cardinality of the state space is denoted by $|\mathcal{S}|$.

Definition 1.2 (Row Stochastic Matrix). A matrix $\mathcal{P} \in \mathbb{R}^{|\mathcal{S}| \times |\mathcal{S}|}$ is called a row stochastic matrix if it satisfies the following conditions:

$$\mathcal{P}_{ij} \in [0, 1] \quad \forall i, j \in \mathcal{S}$$

$$\sum_{j=1}^{|\mathcal{S}|} \mathcal{P}_{ij} = 1 \quad \forall i \in \mathcal{S}$$

Thus, with these two definitions we can define a Markov Chain (or DTMC) as follows:

Definition 1.3 (Markov Chain). Let \mathcal{S} be a finite state space, and ν be a distribution over \mathcal{S} .

$$\nu = (\nu_1, \nu_2, \dots, \nu_{|\mathcal{S}|}) \quad \text{s.t.} \quad \nu_i \in [0, 1] \quad \forall i \in \mathcal{S}, \quad \sum_{i \in \mathcal{S}} \nu_i = 1$$

Further, let \mathcal{P} be a row stochastic matrix over \mathcal{S} .

Then a DTSP $(X_n)_{n \geq 0}$ is called a Markov Chain with initial distribution ν and transition matrix \mathcal{P} if it satisfies the following conditions:

- Initial state is distributed according to ν , i.e.

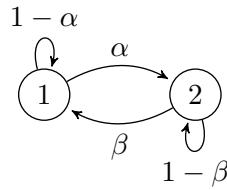
$$\mathbb{P}(X_0 = i) = \nu_i \quad \forall i \in \mathcal{S}$$

- Present is independent of the past given the present.

$$\begin{aligned}\mathbb{P}(X_{n+1} = i_{n+1} | X_n = i_n, X_{n-1} = i_{n-1}, \dots, X_0 = i_0) &= \mathbb{P}(X_{n+1} = i_{n+1} | X_n = i_n) \\ &= \mathcal{P}_{i_n i_{n+1}} \quad \forall n \geq 0\end{aligned}$$

Example (Markov Chains). Let a Markov chain be defined as follows:

$$\mathcal{S} = \{1, 2\} \quad \nu = (p, 1-p) \quad \mathcal{P} = \begin{bmatrix} 1-\alpha & \alpha \\ \beta & 1-\beta \end{bmatrix}$$

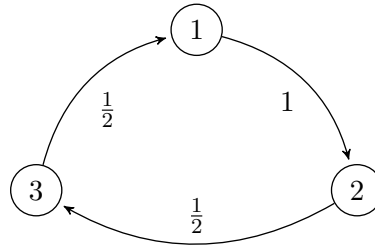


Thus, we can have say:

$$\mathbb{P}(X_3 = 1 | X_0 = 1, X_1 = 1, X_2 = 2) = \beta$$

Another example that we can have is:

$$\mathcal{S} = \{1, 2, 3\} \quad \nu = (\nu_1, \nu_2, \nu_3) \quad \mathcal{P} = \begin{bmatrix} 0 & 1 & 0 \\ 0 & \frac{1}{2} & \frac{1}{2} \\ \frac{1}{2} & 0 & \frac{1}{2} \end{bmatrix}$$



given, $\nu = (1, 0, 0)$, the outcomes of the Markov chain if we sample it are:

$$\begin{aligned}X_n &= 1, 2, 3, 3, 1, 2, \dots \\ &= 1, 2, 3, 1, 2, \dots\end{aligned}$$

Note that, to define the DTMC, we need the initial distribution, but the results won't change on the choice of the distribution ν , assuming that the Markov chain is ergodic. Hence, while defining the Markov chain, the initial distribution is not mentioned, and can be assumed arbitrarily.

Theorem 1.1 (Necessary and Sufficient Conditions for DTSP to be Markov Chains). A DTSP $(X_n)_{n \geq 0}$ on \mathcal{S} is a Markov chain $\langle \mathcal{S}, \mathcal{P}, \nu \rangle$ if and only if:

$$\mathbb{P}(X_n = i_n, X_{n-1} = i_{n-1}, \dots, X_0 = i_0) = \nu_{i_0} \mathcal{P}_{i_0 i_1} \mathcal{P}_{i_1 i_2} \dots \mathcal{P}_{i_{n-1} i_n} \quad \forall n \geq 0$$

Proof. To simplify the proof, assume that $\mathcal{P}_{ij} > 0 \quad \forall i, j \in \mathcal{S}$. The theorem remains true in the general case, but requires more book-keeping.

Suppose, $(X_n)_{n \geq 0}$ is a Markov chain $\langle \mathcal{S}, \mathcal{P}, \nu \rangle$.

Using the fact,

$$\begin{aligned} \mathbb{P}(A \cap B) &= \mathbb{P}(A) \mathbb{P}(B|A) \\ \mathbb{P}(A \cap B \cap C) &= \mathbb{P}(A) \mathbb{P}(B|A) \mathbb{P}(C|A \cap B) \end{aligned}$$

We get,

$$\begin{aligned} \mathbb{P}(X_0 = i_0, \dots, X_n = i_n) &= \mathbb{P}(\{X_0 = i_0\} \cap \{X_1 = i_1\} \cap \dots \cap \{X_n = i_n\}) \\ &= \mathbb{P}(X_0 = i_0) \mathbb{P}(X_1 = i_1 | X_0 = i_0) \dots \mathbb{P}(X_n = i_n | X_{n-1} = i_{n-1}) \\ &= \mathbb{P}(X_0 = i_0) \mathbb{P}(X_1 = i_1 | X_0 = i_0) \dots \mathbb{P}(X_n = i_n | X_{n-1} = i_{n-1}) \\ &\dots \text{Using the memoeryless property of Markov chains} \\ &= \nu_{i_0} \mathcal{P}_{i_0 i_1} \dots \mathcal{P}_{i_{n-1} i_n} \end{aligned}$$

This proves the forward claim. To show the reverse claim:

Put $n = 0$, in the claim, to trivially get the initial distribution back, showing one part of the definition. For the other part:

$$\begin{aligned} \mathbb{P}(X_n = i_n | X_{n-1} = i_{n-1}, \dots, X_0 = i_0) &= \frac{\mathbb{P}(X_n = i_n, \dots, X_0 = i_0)}{\mathbb{P}(X_{n-1} = i_{n-1}, \dots, X_0 = i_0)} \\ &= \frac{\nu_{i_0} \mathcal{P}_{i_0 i_1} \dots \mathcal{P}_{i_{n-1} i_n}}{\nu_{i_0} \mathcal{P}_{i_0 i_1} \dots \mathcal{P}_{i_{n-1} i_{n-1}}} \\ &= \mathcal{P}_{i_{n-1} i_n} \end{aligned}$$

Now we need to show:

$$\begin{aligned} &\mathbb{P}(X_n = i_n | X_{n-1} = i_{n-1}) \\ &= \frac{\mathbb{P}(X_{n-1} = i_{n-1}, X_n = i_n)}{\mathbb{P}(X_{n-1} = i_{n-1})} \\ &= \frac{\sum_{i_0 \in \mathcal{S}} \mathbb{P}(X_0 = i_0 \dots X_{n-1} = i_{n-1}, X_n = i_n)}{\sum_{i_0 \in \mathcal{S}} \mathbb{P}(X_{n-1} = i_{n-1}, X_0 = i_0)} \\ &= \mathcal{P}_{i_{n-1} i_n} \end{aligned}$$

☺

In the above proof we have used the following series of fact:

$$\begin{aligned}\mathbb{P}(X_2 = j) &= \mathbb{P}(\Omega \cap \{X_2 = j\}) \\ &= \mathbb{P}\left[\bigcup_{i=1}^{|S|} \{X_i = i\} \cap \{X_2 = j\}\right]\end{aligned}$$

Using the fact:

$$(A \cup B) \cap C = (A \cap C) \cup (B \cap C)$$

we get:

$$\mathbb{P}(X_2 = j) = \sum_{i=1}^{|S|} \mathbb{P}(X_i = i, X_2 = j)$$

Theorem 1.2 (Markov Property). Let $(X_n)_{n \geq 0}$ be the markov chain denoted by $\langle \mathcal{S}, \mathcal{P}, \mu \rangle$, then conditional on $\{X_m = i\}$, $\{X_{m+n}\}$ is the markov chain $\langle \mathcal{S}, \mathcal{P}, \delta_i \rangle$ and is independent of X_1, X_2, \dots, X_m , where δ_i is the distribution that is 1 at i and 0 everywhere else, i.e

$$\delta_i(j) \equiv \delta_{ij} = \begin{cases} 1 & \text{if } j = i \\ 0 & \text{otherwise} \end{cases}$$

Proof. It suffices to show the following, for any $n \geq 0$ and on an event A , related to X_0, \dots, X_m :

$$\begin{aligned}\mathbb{P}[(X_m = i_m, X_{m+1} = i_{m+1}, \dots, X_{m+n} = i_{m+n}) \cap A \mid X_m = i_m] \\ = \mathbb{P}(A \mid X_m = i_m) \delta_{ij} \mathcal{P}_{i_m i_{m+1}} \dots \mathcal{P}_{i_{m+n-1} i_{m+n}}\end{aligned}$$

The above equation is just conditional independence:

$$\mathbb{P}(E_1 \cap E_2 \mid X_m = i) = \mathbb{P}(E_1 \mid X_m = i) \mathbb{P}(E_2 \mid X_m = i)$$

where $E_2 = (X_m = i_m, X_{m+1} = i_{m+1}, \dots, X_{m+n} = i_{m+n})$ and $E_1 = A$. The proof of the above statement goes as follows:

Let A be an elementary event, i.e. $A = \{X_0 = j_0, X_1 = j_1, \dots, X_m = j_m\}$. Consider the LHS of the above equation: Then, we have two cases:

- **Case 1:** $j_m \neq i_m$. Then, the LHS is trivially 0, since A and E_2 are disjoint.
- **Case 2:** $j_m = i_m$. Then, we have further cases. When, $i \neq i_m = j_m$, then the LHS and RHS are again 0. But when $i = i_m = j_m$, then the LHS is:

$$\begin{aligned}\mathbb{P}(X_0 = j_0, X_1 = j_1, \dots, X_m = j_m = i, X_m = i_m = i, X_{m+1} = i_{m+1}, \\ \dots, X_{m+n} = i_{m+n} \mid X_m = i_m = i)\end{aligned}$$

$$\begin{aligned}
&= \frac{\nu(j_0) \mathcal{P}_{j_0 j_1} \cdots \mathcal{P}_{j_{m-1} j_m} \mathcal{P}_{j_m i_{m+1}} \cdots \mathcal{P}_{i_{m+n-1} i_{m+n}}}{\mathbb{P}(X_m = i_m = i)} \\
&= \frac{\mathbb{P}(A) \mathcal{P}_{i_m i_{m+1}} \cdots \mathcal{P}_{i_{m+n-1} i_{m+n}}}{\mathbb{P}(X_m = i_m = i)} \\
&= \mathbb{P}(A \mid X_m = i_m = i) \mathcal{P}_{i_m i_{m+1}} \cdots \mathcal{P}_{i_{m+n-1} i_{m+n}}
\end{aligned}$$

Thus, completing the proof. This can be extended to non elementary event A by using the fact that any event A can be written as a union of elementary events, i.e

$$A = \bigcup_{i=1}^{|S|} \{X_0 = i_0, X_1 = i_1, \dots, X_m = i_m\}$$

⊙

Theorem 1.3 (Linear Algebra and Markov Chains). Let $(X_n)_{n \geq 0}$ be the markov chain $\langle \mathcal{S}, \mathcal{P}, \nu \rangle$. Then:

1. $\mathbb{P}(X_n = j) = (\nu \mathcal{P}^n)_j$
2. $\mathbb{P}_i(X_n = j) = \mathbb{P}(X_{m+n} = j \mid X_m = i) = \mathcal{P}_{ij}^{(n)} \equiv \text{ij entry of } \mathcal{P}^n$

Proof. The second statement is obvious from the first statement. The proof of the first statement goes as follows:

Proof by induction: For $n = 0$, we have:

$$\mathbb{P}(X_0 = j) = \nu_j \equiv \nu(j)$$

Now, assume that the statement is true for n , then we have:

$$\begin{aligned}
\mathbb{P}(X_{n+1} = j) &= \mathbb{P} \left(\bigcup_{i=1}^{|S|} \{X_n = i, X_{n+1} = j\} \right) \\
&= \sum_{i=1}^{|S|} \mathbb{P}(X_n = i, X_{n+1} = j) \\
&= \sum_{i=1}^{|S|} \mathbb{P}(X_n = i) \mathbb{P}(X_{n+1} = j \mid X_n = i) \\
&= \sum_{i=1}^{|S|} \mathbb{P}(X_n = i) \mathcal{P}_{ij} \\
&= \sum_{i=1}^{|S|} (\nu \mathcal{P}^n)_i \mathcal{P}_{ij} \\
&= (\nu \mathcal{P}^{n+1})_j
\end{aligned}$$

Thus, the proof is complete. \ominus

1.2 Communicating Classes

Definition 1.4. Let $\langle \mathcal{S}, \mathcal{P}, \nu \rangle$ be a markov chain. Then, we say i leads to j denoted as $i \rightarrow j$ if:

$$\mathbb{P}_i(X_n = j \text{ for some } n \geq 0) > 0$$

which is equivalent to:

$$\mathbb{P}(X_{m+n} = j \text{ for some } n \geq 0 \mid X_m = i) > 0$$

Definition 1.5. Let $\langle \mathcal{S}, \mathcal{P}, \nu \rangle$ be a markov chain. Then, we say i and j communicate denoted as $i \leftrightarrow j$ if:

$$i \rightarrow j \text{ and } j \rightarrow i$$

Theorem 1.4. The following statements are equivalent:

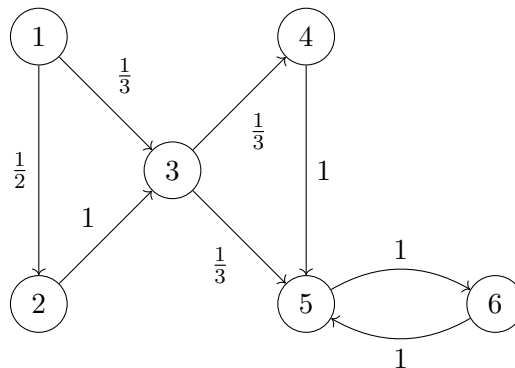
- $i \leftrightarrow j$
- $\mathcal{P}_{i_0 i_1} \mathcal{P}_{i_1 i_2} \dots \mathcal{P}_{i_{n-1} i_n} > 0$ for some $n \geq 0$ and some sequence of states i_0, i_1, \dots, i_n such that $i_0 = i$ and $i_n = j$.
- $\mathcal{P}_{ij}^{(n)} > 0$ for some $n \geq 0$

Proof. i did not understood this :(Revisit this later. \ominus

Definition 1.6 (Communicating Class). Each equivalence class of the relation \leftrightarrow is called a communicating class.

A communicating class is called closed if $\forall i \in C, i \rightarrow j \Rightarrow j \in C$.

Example. The properties that hold in this example are:



$$1 \leftrightarrow 2$$

$$2 \leftrightarrow 3 \Rightarrow 1 \leftrightarrow 3$$

We also have the following:

$$3 \rightarrow 4, 5, 6$$

$$4 \rightarrow 5, 6$$

$$5 \leftrightarrow 6$$

Thus, the following can be concluded:

$$\text{Communicating Classes} = \{1, 2, 3\}, \{4\}, \{5, 6\}$$

$$\text{Closed Groups} = \{5, 6\}$$

1.3 Strong Markov Property

Definition 1.7 (Stopping Time). A random variable $T : \Omega \rightarrow \{0, 1, \dots\}$ is called a stopping time if for any $n \geq 0$, the occurrence or non occurrence of the event $\{T = n\}$ can be determined based only on the values of X_0, X_1, \dots, X_n .

Example. The following are examples of stopping times:

First passage time:

$$T_j = \inf\{n \geq 1 : X_n = j\}$$

Last passage time:

$$T_j = \sup\{n \geq 0 : X_n = j\}$$

Definition 1.8 (Strong Markov Property). Let $(X_n)_{n \geq 0}$ be a markov chain $(\mathcal{S}, \mathcal{P}, \mu)$ and let T , be a stopping time. Then, conditional on $T < \infty$ and $X_T = i$, the process $(X_{T+n})_{n \geq 0}$ is a markov chain with initial distribution δ_i and is independent of X_0, X_1, \dots, X_T .

1.4 Recurrence and Transience

Let $(X_n)_{n \geq 0}$ be the markov chain $\langle \mathcal{S}, \mathcal{P}, \nu \rangle$.

Definition 1.9 (Recurrence). A state i is called recurrent if:

$$\mathbb{P}_i(X_i = 1 \text{ for infinitely many } n) = 1$$

Definition 1.10 (Transience). A state i is called transient if:

$$\mathbb{P}_i(X_i = 1 \text{ for infinitely many } n) < 1$$

Loosely, a recurrent state is one that the markov chain keeps visiting, while a transient state is one that the markov chain never visits from some time on.

Definition 1.11 (Passage Time). First Passage Time to state j :

$$T_j := \inf\{n \geq 1 : X_n = j\}$$

or,

$$T_j : \Omega \rightarrow \{0, 1, 2, \dots\} \cup \{\infty\}$$

r th passage time to state j : We can define this inductively as follows:

$$\begin{aligned} T_j^{(0)} &= 0 \Rightarrow T_j^{(1)} = T_j \\ \Rightarrow T_j^{(r+1)} &= \inf\{n > T_j^{(r)} : X_n = j\} \end{aligned}$$

r th excursion or sojourn time:

$$S_j^{(r)} = \begin{cases} T_j^{(r)} - T_j^{(r-1)} & \text{if } T_j^{(r-1)} < \infty \\ 0 & \text{otherwise} \end{cases}$$

Lemma 1.1. For $r = 2, 3, \dots$ conditional on $T_i^{(r-1)} < \infty$, S_i^r is independent of X_m for $0 < m < T_i^{(r-1)}$. Further,

$$\mathbb{P}_i(S_i^r = n \mid T_i^{(r-1)} < \infty) = \mathbb{P}_i(T_i = n)$$

Proof.

$$\begin{aligned} \{T_i^{(r-1)}\} &\subseteq \{X_{T_i^{(r-1)}}\} \\ \therefore \{T_i^{(r-1)}\} &= \{T_i^{(r-1)}\} \cap \{X_{T_i^{(r-1)}} = i\} \end{aligned}$$

This implies conditioning on the event $\{T_i^{(r-1)} < \infty\}$ is equivalent to conditioning on $\{T_i^{(r-1)}\} < \infty$ and $\{X_{T_i^{(r-1)}} = i\}$. Further $T_i^{(r-1)}$ is a stopping time. Thus, conditional on $\{T_i^{(r-1)} < \infty\}$ and $\{X_{T_i^{(r-1)}} = i\}$

$$\left(X_{T_j^{(r-1)}}\right)_{(n \geq 0)} \text{ is a markov chain}$$

Now,

$$\begin{aligned} &\mathbb{P}(S_i^{(r)} = n \mid T_i^{(r-1)} < \infty) \\ &= \mathbb{P}(S_i^{(r)} = n \mid T_i^{(r-1)} < \infty, X_{T_i^{(r-1)}} = i) \\ &= \mathbb{P}(X_{T_i^{(r-1)}+1} \neq i, \dots, X_{T_i^{(r-1)}+n} = i \mid T_i^{(r-1)} < \infty, X_{T_i^{(r-1)}} = i) \\ &= \mathbb{P}(X_1 \neq i, \dots, X_n = i \mid X_0 = i) \\ &= \mathbb{P}_i(T_i = n) \end{aligned}$$

Completing the proof. ⊙

Definition 1.12 (Number of Visits).

$$V_i = \sum_{n=1}^{\infty} \mathbb{I}\{X_n = i\}$$

Thus, we have the following:

$$\mathbb{E}_i(V_i) = \mathbb{E}_i \sum_{n=0}^{\infty} \mathbb{I}\{X_n = i\} = \sum_{n=0}^{\infty} \mathbb{P}_i(X_n = i) = \sum_{n=0}^{\infty} \mathcal{P}_{ii}^{(n)}$$

we can replace the sum with the expectation, using the monotone convergence theorem.

Lemma 1.2. Let $f_i := \mathbb{P}_i(T_i < \infty)$. Now for $r = 0, 1, \dots$ we have:

$$\mathbb{P}_i(V_i > r) = f_i^r$$

Proof. We can prove this using induction.

Base case: $r = 0$:

$$r = 0 \Rightarrow \mathbb{P}_i(V_i > 0) = 1 \text{ and } f_i^0 = 1$$

Thus, the base case holds true trivially.

Suppose, the given claim holds true for some r . Then,

$$\begin{aligned} \mathbb{P}_i(V_i > r + 1) &= \mathbb{P}_i(T_i^{(r+1)} < \infty) \\ &= \mathbb{P}_i(T_i^{(r)} < \infty, S_i^{(r+1)} < \infty) \\ &= \mathbb{P}_i(T_i^{(r)} < \infty) \mathbb{P}_i(S_i^{(r+1)} < \infty, T_i^{(r)} < \infty) \\ &= \mathbb{P}_i(V_i > r) \mathbb{P}_i(T_i < \infty) \\ &= f_i^r f_i = f_i^{r+1} \end{aligned}$$

☺

Thus, using the above lemmas we can state the following theorem:

Theorem 1.5. Let $(X_n)_{n \geq 0}$ be the markov chain $\langle \mathcal{S}, \mathcal{P}, \nu \rangle$. Then,

- for $i \in \mathcal{S}$, $\mathbb{P}_i(T_i < \infty) = 1 \Rightarrow i$ is recurrent, and

$$\sum_{n=0}^{\infty} \mathcal{P}_{ii}^{(n)} = \infty$$

- for $i \in \mathcal{S}$, $\mathbb{P}_i(T_i < \infty) < 1 \Rightarrow i$ is transient, and

$$\sum_{n=0}^{\infty} \mathcal{P}_{ii}^{(n)} < \infty$$

Proof. Consider the case of recurrence. We have:

$$\mathbb{P}_i(V_i = \infty) = \mathbb{P}_i\left(\bigcap_{r=1}^{\infty} \{V_i > r\}\right)$$

Since, $\{V_i > r\} \supseteq \{V_i > r+1\}$, using the continuity of probability measure, we get:

$$\mathbb{P}_i(V_i = \infty) = \lim_{r \rightarrow \infty} \mathbb{P}_i(V_i > r) = \lim_{r \rightarrow \infty} f_i^r = 1$$

$$\Rightarrow \mathbb{E}_i(V_i) = \infty \Rightarrow \sum_{n=0}^{\infty} \mathcal{P}_{ii}^{(n)} = \infty$$

Now, consider the case of transience. Suppose, $f_i = \mathbb{P}_i(T_i < \infty) < 1$. Then, we have:

$$\sum_{n=0}^{\infty} \mathcal{P}_{ii}^{(n)} = \sum_{n=0}^{\infty} \mathbb{P}_i(T_i = n) = \sum_{n=0}^{\infty} f_i^n = \frac{1}{1-f_i} < \infty$$

Hence, $\mathbb{P}_i(V_i = \infty) = 0$, which implies that i is transient. \ominus

Theorem 1.6. Let C be a communicating class in a Markov Chain. Then either all states in C are recurrent or all states in C are transient.

Proof. Let C be a communicating class. Suppose $i \in C$ is transient.

$$\Rightarrow \sum_{n=0}^{\infty} \mathcal{P}_{ii}^{(n)} < \infty$$

Let $j \in C$ and $\exists n, m \geq 0$ such that $\mathcal{P}_{ij}^{(n)} > 0$ and $\mathcal{P}_{ji}^{(m)} > 0$. Then, for $r \geq 0$:

$$\begin{aligned} \mathcal{P}_{jj}^{(n+m+r)} &\geq \mathcal{P}_{ji}^{(m)} \mathcal{P}_{ii}^{(r)} \mathcal{P}_{ij}^{(n)} \\ \Rightarrow \mathcal{P}_{jj}^{(r)} &\leq \frac{\mathcal{P}_{ii}^{(n+r+m)}}{\mathcal{P}_{ii}^{(n)} \mathcal{P}_{ii}^{(m)}} \\ \Rightarrow \sum_{r=0}^{\infty} \mathcal{P}_{jj}^{(r)} &\leq \frac{\sum_{r=0}^{\infty} \mathcal{P}_{ii}^{(n+r+m)}}{\mathcal{P}_{ij}^{(n)} \mathcal{P}_{ji}^{(m)}} < \infty \end{aligned}$$

\ominus

Note that recurrence or transience is a class property. And thus the following theorem can be stated:

Theorem 1.7. Every recurrent class must be closed. And, Every finite closed class is recurrent.

Note that, infinite closed classes need not be necessarily recurrent. The classic random walk on a line of integers is a classic example.

1.5 Invariant Distributions

Definition 1.13 (Invariant Distribution). A distribution Π on \mathcal{S} is called invariant for the markov chain $\langle \mathcal{S}, \mathcal{P}, \nu \rangle$ if:

$$\Pi = \Pi \mathcal{P}$$

Thus, a stationary distribution is a left eigenvector of the transition matrix \mathcal{P} with eigenvalue 1. Note that in contrast to linear algebra, vectors are represented as row vectors.

Theorem 1.8. Let $(X_n)_{n \geq 0}$ be the markov chain $\langle \mathcal{S}, \mathcal{P}, \nu \rangle$. Then $X_{m+n} \geq 0$ is also a markov chain with the distribution $\langle \mathcal{S}, \mathcal{P}, \Pi \rangle$, where Π is the invariant distribution of \mathcal{P} .

Proof. We know that,

$$\mathbb{P}(X_n = j) = (\nu \mathcal{P}^n)_j \Rightarrow \mathbb{P}(X_n = j) = (\Pi \mathcal{P}^n)_j = \Pi_j$$

Now,

$$\begin{aligned} & \mathbb{P}(X_m = i_0, X_{m+1} = i_1, \dots, X_{m+n} = i_n) \\ &= \mathbb{P}(X_m = i_0) \mathbb{P}(X_{m+1} = i_1 \mid X_m = i_0) \dots \mathbb{P}(X_{m+n} = i_n \mid X_{m+n-1} = i_{n-1}) \\ &= \Pi_{i_0} \mathcal{P}_{i_0 i_1} \dots \mathcal{P}_{i_{n-1} i_n} \\ &\Rightarrow (X_{m+n})_{n \geq 0} \text{ is a markov chain with the distribution } \Pi \quad \ominus \end{aligned}$$

Theorem 1.9. Every markov chain on a finite set space has at least one invariant distribution.

Proof. Since \mathcal{P} is a row stochastic matrix:

$$\mathcal{P} \mathbf{1} = \mathbf{1}$$

Since \mathcal{P} and \mathcal{P}^\top have the same eigenvalues, we have:

$$\Rightarrow \exists \nu \neq 0 \text{ such that } \nu \mathcal{P} = \nu$$

Since, \mathcal{P} is a real valued matrix, taking the complex conjugate of the above equation, we get:

$$\bar{\nu} \mathcal{P} = \bar{\nu}$$

Thus, adding and subtracting the above two equations, we get:

$$\Re(\nu) \mathcal{P} = \Re(\nu) \quad \text{and} \quad \Im(\nu) \mathcal{P} = \Im(\nu)$$

Since, $\nu \neq 0$, atleast one of $\Re(\nu)$ or $\Im(\nu)$ is non zero. Thus, if \mathcal{P} has a complex left eigenvector, then it has a real left eigenvector.

Now, without loss of generality, let u be a real valued vector such that:

$$u \mathcal{P} = u$$

Defining,

$$u_+(i) := \max\{u(i), 0\} \quad \text{and} \quad u_-(i) := \max\{-u(i), 0\}$$

we have:

$$\Rightarrow u = u_+ - u_-$$

letting,

$$u_+\mathcal{P} =: u_+ \quad \text{and} \quad u_-\mathcal{P} =: u_-$$

we get,

$$u\mathcal{P} = u_+\mathcal{P} - u_-\mathcal{P} = y_+ - y_-$$

Suppose: $u_+(i) > 0 \Rightarrow u_-(i) = 0$.

$$\Rightarrow y_+(i) - y_-(i) = u_+(i) \Rightarrow y_+(i) = u_+(i) \text{ and } y_-(i) = 0$$

Similiarly

$$u_-(i) > 0 \Rightarrow u_+(i) = 0 \text{ and } y_-(i) = u_-(i)$$

Thus, we can conclude:

$$u_+\mathcal{P} = u_+ \quad \text{and} \quad u_-\mathcal{P} = u_-$$

Since $u \neq 0$, either of u_+ or u_- is non zero. Thus, we have found a non zero left real valued eigenvector of \mathcal{P} .

Let z be the non zero vector among u_+ or u_- . Then, we have:

$$z\mathcal{P} = z \quad \text{and} \quad z \neq 0, z_i > 0 \Rightarrow \sum_i z_i > 0$$

Rescaling the vector we get,

$$\frac{z}{\sum_i z_i} \mathcal{P} = \frac{z}{\sum_i z_i}$$

Defining $\Pi(i) = \frac{z}{\sum_i z_i}$, we get:

$$\Pi\mathcal{P} = \Pi$$

⊕

Definition 1.14. A markov chain is said to be irreducible if the whole state space is one communicating class. i.e

$$i \leftrightarrow j \quad \forall i, j \in \mathcal{S}$$

Theorem 1.10. Let $\langle \mathcal{S}, \mathcal{P}, \nu \rangle$ be an irreducible and recurrent markov chain. Further let

$$\gamma_i^k = \mathbb{E}_k \sum_{n=0}^{T_k-1} \mathbb{I}\{X_n = i\}$$

and $\gamma^k = (\gamma_1^k, \gamma_2^k, \dots, \gamma_{|\mathcal{S}|}^k)$. Then, the following holds:

1. $\gamma_k^k = 1$
2. $\gamma^k P = \gamma^k$
3. $0 < \gamma^k < \infty$

Proof. (1): This is true by the definition of γ_i^k .

(2): We have:

$$\begin{aligned}
 \gamma_j^k &= \mathbb{E}_k \sum_{n=0}^{T_k-1} \mathbb{I}\{X_n = j\} \\
 &= \mathbb{E}_k \sum_{n=1}^{T_k} \mathbb{I}\{X_n = j\} = \mathbb{E}_k \sum_{n=0}^{\infty} \mathbb{I}\{X_n = j, n \leq T_k\} \\
 &= \sum_{n=1}^{\infty} \mathbb{P}_k(X_n = j, n \leq T_k) \quad // \text{by monotone convergence theorem} \\
 &= \sum_{n=1}^{\infty} \sum_{i \in \mathcal{S}} \mathbb{P}_k(X_n = j, X_{n-1} = i, n \leq T_k) \\
 &= \sum_{n=1}^{\infty} \sum_{i \in \mathcal{S}} \mathbb{P}_k(X_n = j \mid X_{n-1} = i, n \leq T_k) \mathbb{P}_k(X_{n-1} = i, n \leq T_k) \\
 &= \sum_{n=1}^{\infty} \sum_{i \in \mathcal{S}} \mathcal{P}_{ij} \mathbb{P}_k(X_{n-1} = i, n \leq T_k)
 \end{aligned}$$

Since we are summing over non negative sum, we can interchange the summation.

Thus, we get:

$$\begin{aligned}
 \gamma_j^k &= \sum_{i \in \mathcal{S}} \mathcal{P}_{ij} \sum_{n=1}^{\infty} \mathbb{P}_k(X_{n-1} = i, n \leq T_k) \\
 &= \sum_{i \in \mathcal{S}} \mathcal{P}_{ij} \underbrace{\sum_{m=0}^{\infty} \mathbb{P}_k(X_m = i, m \leq T_k - 1)}_{\mathbb{E}_k \sum_{m=0}^{T_k-1} \mathbb{I}\{X_m = i\}} \\
 &= \sum_{i \in \mathcal{S}} \mathcal{P}_{ij} \gamma_i^k \\
 &\Rightarrow \gamma^k P = \gamma^k
 \end{aligned}$$

(3): We will use the fact that the given markov chain is irreducible.

$$\Rightarrow \mathcal{P}_{ij}^{(n)} > 0 \quad \forall i, j \in \mathcal{S}$$

Since we know that γ^k is invariant $\Rightarrow \gamma^k = \gamma^k \mathcal{P} = \gamma^k \mathcal{P}^n$. Thus, we have:

$$\gamma_k^k = \sum_j \gamma_j^k \mathcal{P}_{jk}^{(n)} \geq \gamma_i^k \mathcal{P}_{ik}^{(n)}$$

The above statement holds, since sum of non negative quantities is greater than or equal to any of the individual quantities. Now,

$$\gamma_k^k = 1 \quad \because \mathcal{P}_{ik}^{(n)} > 0 \Rightarrow \gamma_i^k < \infty$$

Now, to show $\gamma_i^k > 0$. Let $\mathcal{P}_{ki}^{(m)} > 0$

$$\Rightarrow \gamma_i^k = \sum_j \gamma_j^k \mathcal{P}_{ji}^{(m)} \geq \gamma_k^k \mathcal{P}_{ki}^{(m)} = \mathcal{P}_{ki}^{(m)} > 0$$

$$\Rightarrow 0 < \gamma_i^k < \infty$$

Completing the proof. ☺

Theorem 1.11. Let $\langle \mathcal{S}, \mathcal{P}, \nu \rangle$ be an irreducible markov chain. Further, let λ be an invariant measure such that $\lambda_k = 1$, then $\lambda \geq \gamma^k$. If $\langle \mathcal{S}, \mathcal{P}, \nu \rangle$ is additionally recurrent, then $\gamma^k = \lambda$.

Proof. NOTE: This theorem also applies with infinite markov chain.

Since λ , is a measure, thus, $\lambda(i)$ is non negative. Now, since λ is invariant, we have:

$$\lambda = \lambda \mathcal{P} \Rightarrow \lambda(j) = \sum_{i \in \mathcal{S}} \lambda(i) \mathcal{P}_{ij} = \sum_{i_1 \neq k} \lambda(i) \mathcal{P}_{ij} + \mathcal{P}_{kj}$$

Massaging \mathcal{P}_{kj} , we get:

$$\mathcal{P}_{kj} = \mathbb{P}(X_1 = j \mid X_0 = k) = \mathbb{P}_k(X_1 = j) = \mathbb{P}_k(X_1 = j, T_k > 1)$$

Thus,

$$\begin{aligned} \lambda_j &= \sum_{i_1 \neq k} \sum_{i_2 \in \mathcal{S}} \lambda_{i_2} \mathcal{P}_{i_2 i_1} \mathcal{P}_{i_1 j} + \mathbb{P}_k(X_1 = j, T_k > 1) \\ &= \sum_{i_1 \neq k} \sum_{i_2 \neq k} \lambda_{i_2} \mathcal{P}_{i_2 i_1} \mathcal{P}_{i_1 j} + \mathbb{P}_k(X_1 = j, T_k > 1) + \mathbb{P}_k(X_2 = j, T_k > 2) \end{aligned}$$

From induction, we get:

$$\begin{aligned} \lambda_j &= \underbrace{\text{some expression}}_{\text{non negative}} + \sum_{n=1}^m \mathbb{P}_k(X_n = j, T_k > n) \\ \Rightarrow \lambda_j &\geq \sum_{n=1}^{\infty} \mathbb{P}_k(X_n = j, T_k > n) = \gamma_j^k \quad // \text{by Theorem 1.10} \end{aligned}$$

Thus, we have:

$$\lambda \geq \gamma^k$$

Now, if the markov chain is recurrent, then we have to show that $\lambda = \gamma^k$. Let $\mu = \lambda - \gamma^k$. Thus, $\mu\mathcal{P} = (\lambda - \gamma^k)\mathcal{P} = \lambda - \gamma^k = \mu$. Now,

$$\mu_k = \lambda_k - \gamma_k^k = 1 - 1 = 0 \quad \mu = \mu\mathcal{P}$$

Choose $i, k \in \mathcal{S}$ such that $\mathcal{P}_{ik}^{(n)} > 0$. Then,

$$\begin{aligned} 0 = \mu_k &= \sum_{j \in \mathcal{S}} \mu_j \mathcal{P}_{jk}^{(n)} \\ &\geq \mu_i \mathcal{P}_{ik}^{(n)} \Rightarrow \mu_i = 0 \quad \forall i \in \mathcal{S} \\ &\Rightarrow \mu = 0 \Rightarrow \lambda = \gamma^k \end{aligned}$$

Completing the proof. \ominus

Theorem 1.12. Let $\langle \mathcal{S}, \mathcal{P}, \nu \rangle$ be a finite irreducible markov chain. Then, Π is unique and $\Pi_i = \frac{1}{m_i}$, where $m_i = \mathbb{E}_i(T_i)$.

Proof. The intuitive meaning of m_i is the expected time to return to state i , starting from state i .

Since the markov chain is finite, there exists an invariant distribution Π such that $\Pi = \Pi\mathcal{P}$.

Now,

$$\pi_k = \sum_{j \in \mathcal{S}} \pi_j \mathcal{P}_{jk}^{(n)} \geq \pi_i \mathcal{P}_{ik}^{(n)} > 0$$

Thus, $\pi_i > 0$. Now, let $\lambda := \frac{\Pi}{\pi_k}$. Thus, $\lambda_k = 1$

$$\Rightarrow \lambda\mathcal{P} = \frac{1}{\pi_k} \Pi\mathcal{P} = \frac{1}{\pi_k} \Pi = \lambda$$

$\Rightarrow \lambda$ is an invariant measure

\because chain is irreducible and from [Theorem 1.10](#), $\lambda = \gamma^k$

$$\Rightarrow \frac{\Pi}{\pi_k} = \gamma^k \Rightarrow \sum_i \frac{\Pi_i}{\pi_k} = \sum_i \gamma_i^k \Rightarrow \frac{1}{\pi_k} = \sum_i \gamma_i^k$$

$$\because \gamma_i^k > 0 \Rightarrow \pi_k = \frac{1}{\sum_i \gamma_i^k}$$

Now,

$$\Rightarrow \sum_i \gamma_i^k = \sum_{i \in \mathcal{S}} \sum_{n=1}^{\infty} \mathbb{P}_k(X_n = i, n \leq T_k)$$

Now, sum is over all positive quantities, summation can be interchanged. Thus, we get:

$$\begin{aligned}\sum_i \gamma_i^k &= \sum_{n=1}^{\infty} \sum_{i \in \mathcal{S}} \mathbb{P}_k(X_n = i, n \leq T_k) \\ &= \sum_{n=1}^{\infty} \mathbb{P}_k(T_k \geq n) = \mathbb{E}_k(T_k) = m_k \Rightarrow \pi_k = \frac{1}{\mathbb{E}_k(T_k)}\end{aligned}$$

Thus, completing the proof. \ominus

Definition 1.15. A state i , is said to be aperiodic if $\mathcal{P}_{ii}^n > 0$ for sufficient large n .

Lemma 1.3. $\langle \mathcal{S}, \mathcal{P}, \sigma \rangle$ is an irreducible markov chain further, suppose $k \in \mathcal{S}$ is aperiodic. Thus, every $i \in \mathcal{S}$, is aperiodic. Infact, $\mathcal{P}_{ii}^{(n)} > 0$, $\forall i, j \in \mathcal{S}$ for sufficient large n .

Theorem 1.13. Let $\langle \mathcal{S}, \mathcal{P}, \nu \rangle$ be a finite irreducible and aperiodic markov chain. Then, irrespective of the initial distribution, the markov chain converges to the unique invariant distribution Π . i.e.

$$\lim_{n \rightarrow \infty} \mathbb{P}(X_n = i) = \Pi_i \quad \forall i \in \mathcal{S}$$

Further,

$$\lim_{n \rightarrow \infty} \mathcal{P}_{ij}^{(n)} = \Pi_j$$

2 Markov Decision Processes and Dynamic Programming

2.1 Controllrd Mrkov Chain

let \mathcal{S} be some finite state space and for all $i \in \mathcal{S}$ let $\mathcal{A}(i)$ denote the set of feasible actions at state i . Then a discrete time stochastic process $\{X_n\}_{n \geq 0}$ is said to be controlled markov chia, with initial distribution ν and transition matrix $\mathcal{P} \in \mathbb{R}^{n \times n}$ or a markov chain controlled by \mathcal{Z}_n if

$$\begin{aligned} \mathbb{P}(X_{n+1} = i_{n+1} \mid X_0 = i + 0, \mathcal{Z}_0 = a_0, \dots, X_n = i_n, \mathcal{Z}_n = a_n) \\ = \mathbb{P}(X_{n+1} = i_{n+1} \mid X_n = i_n, \mathcal{Z}_n = a_n) \equiv \mathcal{P}(i_{n+1} \mid i_n, a_n) \end{aligned}$$

where $\mathcal{S} = |\mathcal{S}|, A = \left| \bigcup_i \mathcal{A}(i) \right|$, if $\mathbb{P}(X_0 = i) = \nu(i)$

2.2 Markov Decision Process

An Markov Decision Process is a acontrolled markov chain with an additional cost strtucture.

$$g : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{R}$$

That is $g(i, a, j)$ is the cost incurred if $X_n = i, \mathcal{Z}_n = a, X_{n+1} = j$.

We denote a MDP with $\langle \mathcal{S}, \mathcal{A}, \mathcal{P}, g, \nu \rangle$

The loose goal with a MDP is to find the sequence of action that minimizes some long term cost. We can formulate this in 3 different yet similar settings:

Finite Horoxom MDP: here $N < \infty$ and the end of the horizon is deterministic. Thus we can define:

$$J_z(i) = \mathbb{E} \left[\sum_{n=0}^{N-1} \alpha^n g(X_n, \mathcal{Z}_n, X_{n+1}) + \alpha^N G(X_N) \mid X_0 = i \right]$$

The term $J_z(i)$ is cost to go from the state i to the end of the horizon. and $\alpha = [0, 1]$ is called as the discount factor.

Stochastic Shortest Path: Here $N < \infty$ where N itself is random. We can define the cost to go in this setting as:

$$J_z(i) = \mathbb{E} \left[\sum_{n=0}^{N-1} \alpha^n g(X_n, \mathcal{Z}_n, X_{n+1}) + \alpha^N G(X_N) \mid X_0 = i \right]$$

Note that here the discount factor α can be 1.

Inifinite Horizon MDP: Here $N = \infty$ and we can define the cost to go as:

$$J_z(i) = \mathbb{E} \left[\sum_{n=0}^{\infty} \alpha^n g(X_n, \mathcal{Z}_n, X_{n+1}) \mid X_0 = i \right]$$

In the case of

Definition 2.1 (Policy). Let $M \equiv \langle \mathcal{S}, \mathcal{A}, \mathcal{P}, g, \rangle$ be a MDP. Then $\pi \equiv \mu_k$, $k = 0$ is said to be a policy or a strategy w.r.t. M if:

$$\mu_k : \mathcal{S} \rightarrow \mathcal{A} = \bigcup_{i \in \mathcal{S}} \mathcal{A}(i)$$

that is each μ_k is a function that maps a state to an action.

We say that a policy π is admissible if $\mu_k \in \mathcal{A}(i) \forall i, k$, where $\mathcal{A}(i)$ is the set of feasible or allowable actions at that state i

Definition 2.2 (Expected Cost to Go). For finite horizon MDP or SSP, the expected cost to go under policy π is given by:

$$J^\pi(i) = \mathbb{E} \left[\sum_{n=0}^{N-1} \alpha^n g(X_n, \mathcal{Z}_n, X_{n+1}) + \alpha^N G(X_N) \mid X_0 = i \right]$$

For Infinite Horizon MDP, the expected cost to go under policy π is given by:

$$J^\pi(i) = \mathbb{E} \left[\sum_{n=0}^{\infty} \alpha^n g(X_n, \mathcal{Z}_n, X_{n+1}) \mid X_0 = i \right]$$

Definition 2.3 (Optimal Cost to Go and Optimal Policy).

$$J^*(i) = \inf_{\pi} J^\pi(i)$$

And the policy π^* is said to be optimal if:

$$J^{\pi^*}(i) = J^*(i) \quad \forall i \in \mathcal{S}$$

2.3 Finite Horizon Problems

To find the optimal policy for finite horizon problems, we will make use of the dynamic programming approach. Consider the case where there is only one stage i.e. $N = 1$. Then the cost to go following a policy π is given by:

$$\begin{aligned} J^\pi(i) &= \mathbb{E} [g(i, \mu_0(i), j) + \alpha G(j) \mid X_0 = i] \\ &= \sum_{j=1}^S \mathbb{P}(j \mid i, \mu_0(i)) [g(i, \mu_0(i), j) + \alpha G(j)] \end{aligned}$$

Then, optimal cost to go is by definition:

$$J^*(i) = \min_{\pi} J^\pi(i) = \min_{\mu_0} J^\pi(i)$$

For any fixed state i , the minimisation over π is same as minimisation over μ_0 by definition. Thus,

$$J^\pi(i) = \min_{a \in \mathcal{A}_i} \sum_{j=1}^S \mathbb{P}(j \mid i, a) [g(i, a, j) + \alpha G(j)]$$

where μ_0 is such that,

$$\mu_0(i) = \arg \min_{a \in \mathcal{A}_i} \sum_{j=1}^S \mathbb{P}(j \mid i, a) [g(i, a, j) + \alpha G(j)] \quad \forall i \in \mathcal{S}$$

The interpretation is that optimal control choice with one stage to go must minimise the sum of expected cost to go. The DP algoirthm extends this idea to a N stage problem. It states that the optimal control choice with N stages to go must minimise the sum of expected present stage cost and expected optimal cost J_{k-1}^* with $k-1$ stages to go, discounted by α . Thus optimal N stage cost to go $J_N^*(i)$ can be computed by:

$$J_N^*(i) = \min_{a \in \mathcal{A}_i} \sum_{j=1}^S \mathbb{P}(j \mid i, a) [g(i, a, j) + \alpha J_{N-1}^*(j)]$$

starting with,

$$J_0^* = G(i) \quad \forall i \in \mathcal{S}$$

To prove this, we can note that any policy can be broken down as: $\pi^k \equiv (\mu_{N-k}, \pi^{k-1})$

$$\begin{aligned} \Rightarrow J_k^*(i) &= \min_{(\mu_{N-k}, \pi^{k-1})} \mathbb{E} \left[g(X_{N-k}, \mu_{N-k}(X_{N-k}), X_{N-k+1}) + \alpha J_{k-1}^{\pi^{k-1}}(X_{N-k+1}) \mid X_{N-k} = i \right] \\ &= \min_{\mu_{N-k}} \sum_{j=1}^S \mathbb{P}(j \mid i, \mu_{N-k}(i)) \left[g(i, \mu_{N-k}(i), j) + \alpha J_{k-1}^{\pi^{k-1}}(j) \right] \\ &= \min_{a \in \mathcal{A}_i} \sum_{j=1}^S \mathbb{P}(j \mid i, a) \left[g(i, a, j) + \alpha \min_{\pi^{k-1}} J_{k-1}^{\pi^{k-1}}(j) \right] \\ &= \min_{a \in \mathcal{A}_i} \sum_{j=1}^S \mathbb{P}(j \mid i, a) [g(i, a, j) + \alpha J_{k-1}^*(j)] \end{aligned}$$

Example (Chess Match). Aim is to find an optimal strategy for the player for the game i . The player can opt for

$$\underbrace{\text{timid play}}_{\text{never wins}} = \begin{cases} p_d, & \rightarrow \text{probability of draw} \\ 1 - p_d, & \rightarrow \text{probability of loss} \end{cases}$$

$$\underbrace{\text{bold play}}_{\text{never draws}} = \begin{cases} p_w, & \rightarrow \text{probability of win} \\ 1 - p_w, & \rightarrow \text{probability of loss} \end{cases}$$

There are N games to be played. If the scores are tied after N games, then the game goes into sudden death mode. i.e. the player who wins the first wins the game.

We have to find the policy that maximised the prob for win, when the single game reward is:

$$\begin{cases} 1, & \text{if win} \\ 0.5, & \text{if draw} \\ 0, & \text{if loss.} \end{cases}$$

when the state is given by the net score i.e. points of the player - points of the opponent, with the ternimal reward being:

$$G(i) = \begin{cases} 1, & \text{if } i > 0; \\ p_w, & \text{if } i = 0; \\ 0, & \text{otherwise } i < 0. \end{cases}$$

Thus, we let,

$$J_N(i) = G(i) \quad \forall i \in \mathcal{S}$$

Now, we can clearly see:

$$J_k^*(i) = \max\{p_d J_{k+1}^*(i) + (1 - p_d) J_{k+1}^*(i - 1), p_w J_{k+1}^*(i + 1) + (1 - p_w) J_{k+1}^*(i - 1)\}$$

Thus, it is optimal to play bold, starting from i if:

$$p_w J_{k+1}^* + (1 - p_w) J_{k+1}^*(i - 1) > p_d J_{k+1}^* + (1 - p_d) J_{k+1}^*(i - 1)$$

That is if:

$$p_w (J_{k+1}^*(i + 1)) - J_{k+1}^*(i) > p_d (J_{k+1}^*(i)) - J_{k+1}^*(i - 1)$$

which holds when,

$$\frac{p_w}{p_d} > \frac{J_{k+1}^*(i) - J_{k+1}^*(i - 1)}{J_{k+1}^*(i + 1) - J_{k+1}^*(i)}$$

Assume that $p_d > p_w$, then from dynamic programming:

$$J_{N-1}^*(i) = \max\{p_d J_N^*(i) + (1 - p_d) J_N^*(i - 1), p_w J_N^*(i + 1) + (1 - p_w) J_N^*(i - 1)\}$$

Thus, we have two cases:

Case 1: $i > 1$, Then, $J_{N-1}^*(i) = 1$. Thus, any action is optimal action.

Case 2: $i = 1$, Then $J_{N-1}^*(i) = \max\{p_d + (1 - p_d)p_w, p_w + (1 - p_w)p_w\}$

Since we assumed $p_d > p_w$:

$$\Rightarrow \frac{p_w}{p_d} > \frac{1 - p_w}{1 - p_w}$$

Thus, optimal action is to play timid

2.4 Stochastic Shortest Path

Since the end of the horizon is random, we cannot directly apply the principles of finite horizon problems. Thus, we create a slightly generalised approach. Here we assume that there is no discounting. To make the cost to go reasonable, we assume that there exists a cost-free termination state 0, with terminal cost of 0. Thus the goal of SSP problems is to reach the terminal state with least expected cost.

Definition 2.4 (Admissible Policy). A policy π is said to be admissible if $\mu_k \in \mathcal{A}(i) \forall i, k$

Definition 2.5 (Stationary Policy). An admissible policy is said to be stationary if $\mu_0 = \mu_1 = \mu_2 = \dots = \mu$. We will often denote the stationary policy by μ .

Definition 2.6 (Proper Policy). A stationary policy is said to be proper if there exists a positive probability that the termination state will be reached in at most n stages, regardless of the initial state.

$$\Rightarrow \rho_\mu = \max_{i=1, \dots, n} \mathbb{P}(X_n \neq 0 \mid X_0 = i, \mu) < 1$$

where $n = |\mathcal{S}| - 1$, thus n represents the number of non terminal states.

A improper policy is the one which is not proper.

Lemma 2.1. Suppose μ is proper, then:

$$\mathbb{P}(X_k \neq 0 \mid X_0 = i, \mu) \leq \rho_\mu^{\lfloor k/n \rfloor}$$

Proof. For $k < n$, the claim is trivially true. Suppose $n \leq k \leq 2n$, for some n . Then we have:

$$\begin{aligned} & \mathbb{P}(X_k \neq 0 \mid X_0 = i, \mu) \\ &= \sum_{j \in \mathcal{S}} \mathbb{P}(X_k \neq 0, X_n = j \mid X_0 = i, \mu) \\ &= \sum_{j=1}^n \mathbb{P}(X_k \neq 0, X_n = j \mid X_0 = i, \mu) && \because 0 \text{ is terminal state} \\ &= \sum_{j=1}^n \mathbb{P}(X_k \neq 0 \mid X_n = j, X_0 = i, \mu) \mathbb{P}(X_n = j \mid X_0 = i, \mu) \\ &\leq \sum_{j=1}^n \mathbb{P}(X_n = j \mid X_0 = i, \mu) && \because \mu \text{ is proper} \\ &= \mathbb{P}(X_n \neq 0 \mid X_0 = i, \mu) \leq \rho_\mu \end{aligned}$$

using induction, we can show:

$$mn \leq k \leq (m+1)n$$

⊕

Assumption 1. There exists at least one proper policy

Assumption 2. For every improper policy μ , the corresponding cost to go $J^\mu(i)$ is infinite

for some i

$$\Rightarrow J^\mu(i) = \mathbb{E} \left[\sum_{n=0}^{\infty} g(X_n, \mu_n(X_n), X_{n+1}) \mid X_0 = i \right] = \infty$$

Thus, we will search through proper policies to find the optimal policy.

Lemma 2.2. Suppose μ is proper and $|g(i, a, j)| \leq k, \forall i, a, j$, then:

$$|J_\mu(i)| < \infty \quad \forall i$$

Proof.

$$\begin{aligned} J_\mu(i) &= \mathbb{E} \left[\sum_{m=0}^{N-1} g(X_m, \mu(X_m), X_{m+1}) \mid X_0 = i, \mu \right] \\ \Rightarrow |J_\mu(i)| &= \left| \mathbb{E} \left[\sum_{m=0}^{\infty} g(X_m, \mu(X_m), X_{m+1}) \mid X_0 = i, \mu \right] \right| \quad \because \text{terminal cost} = 0 \\ &\leq \mathbb{E} \left[\sum_{m=0}^{N-1} |g(X_m, \mu(X_m), X_{m+1})| \mid X_0 = i, \mu \right] \quad \because |\mathbb{E}X| \leq \mathbb{E}|X| \\ &= \sum_{m=0}^{\infty} \mathbb{E} [|g(X_m, \mu(X_m), X_{m+1})| \mid X_0 = i, \mu] \quad // \text{monotone convergence} \\ &= \sum_{m=0}^{\infty} \sum_{j, k \in \mathcal{S}} \mathbb{P}(X_m = j, \mu(j), X_{m+1} = k \mid X_0 = i, \mu) |g(j, \mu(j), k)| \\ &= \sum_{m=0}^{\infty} \sum_{j=1}^n \mathbb{P}(X_m = j \mid X_0 = i, \mu) \underbrace{\mathbb{P}(\mu(j) \mid X_0 = i, \mu, X_m = j)}_1 \\ &\quad \underbrace{\sum_{k \in \mathcal{S}} \mathbb{P}(X_{m+1} = k \mid X_m = j, \mu(j))}_{1} \underbrace{|g(j, \mu(j), k)|}_{\leq k} \\ &\leq k \sum_{m=0}^{\infty} \sum_{j=1}^n \mathbb{P}(X_m = j \mid X_0 = i, \mu) \\ &= k \sum_{m=0}^{\infty} \mathbb{P}(X_m \neq 0 \mid X_0 = i, \mu) \leq k \sum_{m=0}^{\infty} \rho_\mu^{\lfloor m/n \rfloor} = k \sum_{l=0}^{\infty} \sum_{m=l-n}^{(l+1)n-1} \rho_\mu^l \\ &= kn \sum_{l=0}^{\infty} \rho_\mu^l = \frac{kn}{1 - \rho_\mu} < \infty \end{aligned}$$

Thus, completing the proof. \ominus

Definition 2.7 (T operator).

$$T : \mathbb{R}^n \rightarrow \mathbb{R}^n \quad n = |\mathcal{S}| - 1$$

Thus, for any vector $J = (J(1), J(2), \dots, J(n))$ we consider the vector TJ obtained by one iteration of the DP algorithm to J . Thus,

$$(TJ)(i) = \min_{a \in \mathcal{A}_i} \sum_{j \in \mathcal{S}} \mathbb{P}(j | i, a) (g(i, a, j) + J(j))$$

TJ is the optimal cost-to-go vector for the one stage problem that has one stage cost g and terminal cost J . Similarly, for any vector J and any stationary policy μ we consider the vector $T_\mu J$ with components,

$$T_\mu J(i) = \sum_{j \in \mathcal{S}} \mathbb{P}(j | i, \mu(i)) (g(i, \mu(i), j) + J(j))$$

Given a stationary policy μ , we can define the $n \times n$ matrix \mathcal{P}_μ whose ij th entry is $\mathcal{P}_\mu(i, j) = \mathbb{P}(j | i, \mu(i))$. Then, we can write:

$$T_\mu J = \bar{g}_\mu + \mathcal{P}_\mu J$$

where $\bar{g}_\mu \in \mathbb{R}^n$ is the vector with i th component,

$$\bar{g}_\mu(i, \mu(i)) = \sum_{j \in \mathcal{S}} \mathbb{P}(j | i, \mu(i)) g(i, \mu(i), j)$$

Definition 2.8 (T^k operator). We define the T^k operator as the composition of T with itself k times. That is:

$$(T^k J)(i) = (T \circ T^{k-1})J = T(T^{k-1}J)$$

Similarly,

$$T_\mu^k J = (T_\mu \circ T_\mu^{k-1})J = T_\mu(T_\mu^{k-1}J)$$

Thus, it can be seen, $(T^k J)(i)$ is the optimal cost to go for the k stage problem with one stage cost g and terminal cost J . Similarly, $(T_\mu^k J)(i)$ is the cost to go for the stationary policy μ for the same problem.

Lemma 2.3 (Monotonicity Lemma). For any n dimensional vector J, \bar{J} , if $J \leq \bar{J}$, then

- $TJ \leq T\bar{J}$
- $T_\mu J \leq T_\mu \bar{J}$

Proof.

$$\begin{aligned} TJ(i) &= \min_{a \in \mathcal{A}_i} \sum_{j \in \mathcal{S}} \mathbb{P}(j \mid i, a) (g(i, a, j) + J(j)) \\ &\leq \min_{a \in \mathcal{A}_i} \sum_{j \in \mathcal{S}} \mathbb{P}(j \mid i, a) (g(i, a, j) + \bar{J}(j)) = T\bar{J}(i) \end{aligned}$$

Similarly,

$$T_\mu J \leq T_\mu \bar{J}$$

⊖

Corollary 2.1.

$$T^k J \leq T^k \bar{J} \quad \text{and} \quad T_\mu^k J \leq T_\mu^k \bar{J}$$

The above corollary can be proved using induction.

Lemma 2.4. For any $J \in \mathbb{R}^n$ and stationary policy μ and $r \geq 0$:

- $T(J + re)(i) \leq TJ(i) + r$
- $T_\mu(J + re)(i) \leq T_\mu J(i) + r$

where $e = (1, 1, \dots, 1)^\top$

Proof.

$$\begin{aligned} T(j + re)(i) &= \min_{a \in \mathcal{A}_i} \sum_{j \in \mathcal{S}} \mathbb{P}(j \mid i, a) (g(i, a, j) + (J + re)(j)) \\ &= \min_{a \in \mathcal{A}_i} \sum_{j \in \mathcal{S}} \mathbb{P}(j \mid i, a) (g(i, a, j) + J(j) + r) \\ &= \min_{a \in \mathcal{A}_i} \sum_{j \in \mathcal{S}} \mathbb{P}(j \mid i, a) (g(i, a, j) + J(j)) + r \underbrace{\sum_{j \in \mathcal{S}} \mathbb{P}(j \mid i, a)}_{\leq 1} \\ &\leq \min_{a \in \mathcal{A}_i} \sum_{j \in \mathcal{S}} \mathbb{P}(j \mid i, a) (g(i, a, j) + J(j)) + r \\ &= TJ(i) + r \end{aligned}$$

⊖