

# Dynamic Programming And Optimal Control

GianAndrea Müller

October 24, 2018

3.4	Converting non-standard Problems to the Standard Form . . . . .	4
3.4.1	Time Lags . . . . .	4
3.4.2	Correlated Disturbances . . . . .	4
3.4.3	Forecasts . . . . .	5
<b>4</b>	<b>Infinite Horizon Problem</b>	<b>5</b>
4.1	The Stochastic Shortest Path (SSP) Problem . . . . .	6
4.2	Theorem 4.1 . . . . .	6

## CONTENTS

<b>1</b>	<b>Random Variables</b>	<b>2</b>
1.1	Discrete Random Variables (DRV) . . . . .	2
1.1.1	Generalization for multiple variables . . . . .	2
1.2	Continuous Random Variables (CRV) . . . . .	2
1.3	Expectation . . . . .	2
1.3.1	Multi-variable generalizations . . . . .	3
<b>2</b>	<b>Basics</b>	<b>3</b>
2.1	Cost Function . . . . .	3
2.1.1	Expected Cost . . . . .	3
2.2	Open Loop Control . . . . .	3
2.3	Closed Loop Control . . . . .	4
2.4	Discrete State and Finite State Problems . . . . .	4
<b>3</b>	<b>The Dynamic Programming Algorithm</b>	<b>4</b>
3.1	The standard problem formulation . . . . .	4
3.2	Principle of Optimality . . . . .	4
3.3	DPA . . . . .	4

# 1 RANDOM VARIABLES

## 1.1 DISCRETE RANDOM VARIABLES (DRV)

$\mathcal{X}$  set of all possible outcomes  
 $p_x(\cdot)$  probability density function (PDF)

1.  $p_x(\bar{x}) \geq 0 \forall \bar{x} \in \mathcal{X}$

2.  $\sum_{\bar{x} \in \mathcal{X}} p_x(\bar{x}) = 1$

**Definition 1.**  $p_x(\cdot)$  and  $\mathcal{X}$  define a **discrete random variables (DRV)**  $x$ .

The probability that a random variable  $x$  is equal to some value  $\bar{x} \in \mathcal{X}$  is  $p_x(\bar{x})$ . This is written as  $Pr(x = \bar{x}) = p_x(\bar{x})$ .

**Definition 2.** The **joint PDF**  $p_{xy}(\cdot, \cdot)$  is a real valued function that satisfies:

1.  $p_{xy}(\bar{x}, \bar{y}) \geq 0 \forall \bar{x} \in \mathcal{X}, \forall \bar{y} \in \mathcal{Y}$ ,

2.  $\sum_{\bar{x} \in \mathcal{X}} \sum_{\bar{y} \in \mathcal{Y}} p_{xy}(\bar{x}, \bar{y}) = 1$ .

**Definition 3. Marginalization or Sum Rule axiom:**

Given  $p_{xy}(\cdot, \cdot)$  define  $p_x(\bar{x}) := \sum_{\bar{y} \in \mathcal{Y}} p_{xy}(\bar{x}, \bar{y})$ .

**Definition 4. Conditioning or Product Rule axiom:**

Given  $p_{xy}(\cdot, \cdot)$  the PDF of  $y$  is  $p_{x|y}(\bar{x}|\bar{y}) := \frac{p_{xy}(\bar{x}, \bar{y})}{p_y(\bar{y})}$  when  $p_y(\bar{y}) \neq 0$ .

- Sum rule applied to a conditional PDF:

$$\text{Given } p_{xy|z}(\bar{x}, \bar{y}|\bar{z}), p_{x|z}(\bar{x}|\bar{z}) := \sum_{\bar{y} \in \mathcal{Y}} p_{xy|z}(\bar{x}, \bar{y}|\bar{z}).$$

- Short form:  $p(x|y)$

- Product rule usually written as:  $p(x, y) = p(x|y)p(y) = p(y|x)p(x)$

**Definition 5. Total Probability Theorem:**

$$p_x(\bar{x}) = \sum_{\bar{y} \in \mathcal{Y}} p_{x|y}(\bar{x}|\bar{y})p_y(\bar{y})$$

## 1.1.1 GENERALIZATION FOR MULTIPLE VARIABLES

**Definition 6. Marginalization**

$$p_x(\bar{x}) = \sum_{\bar{y} \in \mathcal{Y}} p_{xy}(\bar{x}, \bar{y})$$

as a short form of:

$$p_{x_1, \dots, x_N}(\bar{x}_1, \dots, \bar{x}_N) = \sum_{(\bar{y}_1, \dots, \bar{y}_L) \in \mathcal{Y}} p_{x_1, \dots, x_N}(\bar{x}_1, \dots, \bar{x}_N, \bar{y}_1, \dots, \bar{y}_L).$$

**Definition 7. Conditioning**

$$p(x, y) = p(x|y)p(y)$$

as a short form of:

$$p(x_1, \dots, x_N, y_1, \dots, y_L) = p(x_1, \dots, x_N|y_1, \dots, y_L)p(y_1, \dots, y_L)$$

**Definition 8.** Random variables  $x$  and  $y$  are said to be **independent** if  $p(x|y) = p(x)$ . Equivalently:  $p(x, y) = p(x)p(y)$ .

**Definition 9.** Random variables are said to be **conditionally independent** if:  $p(x|y, z) = p(x|z)$ . Knowledge of  $z$  makes  $x$  and  $y$  independent.

## 1.2 CONTINUOUS RANDOM VARIABLES (CRV)

$\mathcal{X}$  subset of the real line  
 $p(\cdot)$  PDF

1.  $p_x(\bar{x}) \geq 0 \forall \bar{x} \in \mathcal{X}$

2.  $\int_{\mathcal{X}} p_x(\bar{x})d\bar{x} = 1$

**Definition 10.** The **probability of being in an interval** is:

$$Pr(x \in [a, b]) := \int_a^b p_x(\bar{x})d\bar{x}$$

## 1.3 EXPECTATION

**Definition 11.** The **expected value** of a random variable is defined as:

$$E[x] := \sum_{\bar{x} \in \mathcal{X}} \bar{x}p_x(\bar{x})$$

- $E[ax + b] = aE[x] + b$  where  $a, b$  constant

- $E[g(x)] = \sum_{\bar{x} \in \mathcal{X}} \bar{x}p_{x|y}(\bar{x}|\bar{y})$ .

- For conditional PDF's:

$$E[x|y = \bar{y}] := \sum_{\bar{x} \in \mathcal{X}} \bar{x}p_{x|y}(\bar{x}|\bar{y})$$

### 1.3.1 MULTI-VARIABLE GENERALIZATIONS

If  $x$  is a vector:

$$E[x] = \sum_{\bar{x} \in \mathcal{X}} \bar{x} p_x(\bar{x}) = \sum_{\bar{x}_1 \in \mathcal{X}} \cdots \sum_{\bar{x}_N \in \mathcal{X}} [\bar{x}_1, \dots, \bar{x}_N]^T p_{(x_1, \dots, x_N)}(\bar{x}_1, \dots, \bar{x}_N)$$

Given  $g(x) : \mathbb{R}^N \rightarrow \mathbb{R}$  and DRV  $x$

$$E[g(x)] = \sum_{\bar{x} \in \mathcal{X}} g(\bar{x}) p_x(\bar{x}) = \sum_{\bar{x}_1 \in \mathcal{X}} \cdots \sum_{\bar{x}_N \in \mathcal{X}} g(\bar{x}_1, \dots, \bar{x}_N) p_{(x_1, \dots, x_N)}(\bar{x}_1, \dots, \bar{x}_N)$$

If the two random variables are **independent**, then:

$$\begin{aligned} E[g(x, y)] &= \sum_{\bar{y} \in \mathcal{Y}} \sum_{\bar{x} \in \mathcal{X}} g(\bar{x}, \bar{y}) p_{xy}(\bar{x}, \bar{y}) = \sum_{\bar{y} \in \mathcal{Y}} \sum_{\bar{x} \in \mathcal{X}} g(\bar{x}, \bar{y}) p_{xy}(\bar{x}, \bar{y}) = \\ &= \sum_{\bar{y} \in \mathcal{Y}} \sum_{\bar{x} \in \mathcal{X}} g(\bar{x}, \bar{y}) p_x(\bar{x}) p_y(\bar{y}) = E_x[E_y[g(x, y)]] \end{aligned}$$

**Mean and Variance:**

**Definition 12.**  $E[x]$  is called the **mean**, generally a vector.

**Definition 13.**  $Var[x] := E_x \left[ \left( x - E_x[x] \right) \left( x - E_x[x] \right)^T \right]$  is called the **variance**, generally a matrix.

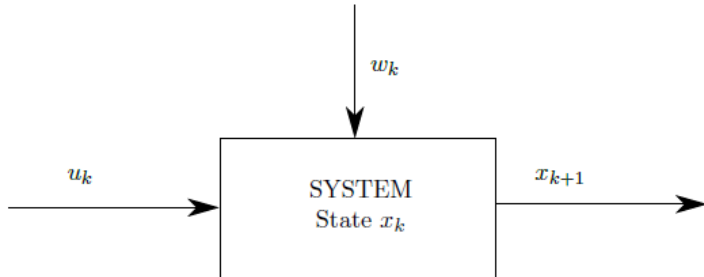
**Linerarity:**

$$E_{xy}[x + y] = \sum_{\bar{y} \in \mathcal{Y}} \sum_{\bar{x} \in \mathcal{X}} (\bar{x} + \bar{y}) p_{xy}(\bar{x}, \bar{y}) = \sum_{\bar{x} \in \mathcal{X}} \bar{x} \sum_{\bar{y} \in \mathcal{Y}} p_{xy}(\bar{x}, \bar{y}) + \sum_{\bar{y} \in \mathcal{Y}} \bar{y} \sum_{\bar{x} \in \mathcal{X}} p_{xy}(\bar{x}, \bar{y}) = E_x[x] + E_y[y]$$

**Law of Total Expectation:**

$$E_y[E_{x|y}[x]] = \sum_{\bar{y} \in \mathcal{Y}} p_y(\bar{y}) \left( \sum_{\bar{x} \in \mathcal{X}} \bar{x} p_{x|y}(\bar{x}|\bar{y}) \right) = \sum_{\bar{x} \in \mathcal{X}} \bar{x} \sum_{\bar{y} \in \mathcal{Y}} p_{x|y}(\bar{x}|\bar{y}) p_y(\bar{y}) = \sum_{\bar{x} \in \mathcal{X}} \bar{x} p_x(\bar{x}) = E_x[x]$$

## 2 BASICS



$$x_{k+1} = f_k(x_k, u_k, w_k), \quad k = 0, 1, \dots, N-1$$

$k$	discrete time index
$N$	given time horizon
$x_k \in \mathcal{S}_k$	system state vector at time $k$
$u_k \in \mathcal{U}_k(x_k)$	control input vector at time $k$
$w_k$	disturbance vector at time $k$
$f_k(\cdot, \cdot, \cdot)$	function capturing system evolution at time $k$

- It is assumed that the conditional probability of  $w_k$  depending on  $u_k$  and  $x_k$  is known and  $w_k$  is independent of any other variables.

### 2.1 COST FUNCTION

$$\underbrace{g_N(x_N)}_{\text{terminal cost}} + \underbrace{\sum_{k=0}^{N-1} g_k(x_k, u_k, w_k)}_{\text{accumulated cost}} \quad \text{stage cost}$$

- Cost is a random variable.

#### 2.1.1 EXPECTED COST

$X_1 := (x_1, \dots, x_N)$	set of all states
$U_0 := (u_0, \dots, u_{N-1})$	set of all inputs
$W_0 := (w_0, \dots, w_{N-1})$	set of disturbances

$$E_{(X_1, U_0, W_0 | x_0)} \left[ g_N(x_N) + \sum_{k=0}^{N-1} g_k(x_k, u_k, w_k) \right]$$

- All variables  $x_k, u_k, w_k$  are in general all random variables, since at least the disturbance is random, and thus coupled by the dynamics and possibly a control law depending on the state, all variables are coupled in general.

### 2.2 OPEN LOOP CONTROL

$$\bar{U}_0 := (\bar{u}_0, \dots, \bar{u}_{N-1}) \quad \text{fixed set of control inputs}$$

$$g_N(x_N) + \sum_{k=0}^{N-1} g_k(x_k, \bar{u}_k, w_k) \quad \text{Open Loop Cost}$$

$$E_{(X_1, W_0 | x_0)} \left[ g_N(x_N) + \sum_{k=0}^{N-1} g_k(x_k, \bar{u}_k, w_k) \right] \quad \text{Expected Open Loop Cost}$$

- The bar emphasizes that the control inputs are fixed.
- Open loop control means that the control law is defined and fixed at time zero. The measurements of the state are not used to adapt the control law.

## 2.3 CLOSED LOOP CONTROL

$$u_k = \mu_k(x), \quad u_k \in \mathcal{U}_k(x), \quad \forall x \in \mathcal{S}_k, \quad k = 0, \dots, N-1 \quad \text{control law}$$

$$\boxed{\pi := (\mu_0(\cdot), \mu_1(\cdot), \dots, \mu_{N-1}(\cdot))} \quad \text{Admissible Policy}$$

$$g_N(x_N) + \sum_{k=0}^{N-1} g_k(x_k, \mu_k(x_k), w_k) \quad \text{Closed Loop Cost}$$

$$\boxed{J_\pi(x) := \underset{(X_1, W_0 | x_0=x)}{E} \left[ g_N(x_N) + \sum_{k=0}^{N-1} g_k(x_k, \mu_k(x_k), w_k) \right]} \quad \text{Expected Closed Loop Cost}$$

$\Pi$  set of all admissible policies

$\pi^*$  is called the optimal policy if

$$J_{\pi^*}(x) \leq J_\pi(x) \quad \forall \pi \in \Pi, \quad \forall x \in \mathcal{S}_0$$

## 2.4 DISCRETE STATE AND FINITE STATE PROBLEMS

$$P_{ij}(u, k) := p_{x_{k+1}|x_k, u_k}(j|i, u) = \Pr(x_{k+1} = j | x_k = i, u_k = u) \quad \text{transition probability}$$

- $p_{x_{k+1}|x_k, u_k}(\cdot|\cdot, \cdot)$  denotes the PDF of  $x_{k+1}$  given  $x_k$  and  $u_k$ .

# 3 THE DYNAMIC PROGRAMMING ALGORITHM

## 3.1 THE STANDARD PROBLEM FORMULATION

$$x_{k+1} = f_k(x_k, u_k, w_k), \quad k = 0, 1, 2, \dots, N-1$$

where  $x_k \in \mathcal{S}_k$ ,  $u_k \in \mathcal{U}_k$  and  $w_k \sim p_{w_k|x_k, u_k}$

The control inputs are generated by the admissible policy  $\pi \in \Pi$

$$\pi = (\mu_0(\cdot), \mu_1(\cdot), \dots, \mu_{N-1}(\cdot))$$

The expected closed loop cost, given  $x \in \mathcal{S}_0$ , associated with policy  $\pi$  is:

$$J_\pi(x) = \underset{(X_1, W_0 | x_0=x)}{E} \left[ g_N(x_N) + \sum_{k=0}^{N-1} g_k(x_k, \mu_k(x_k), w_k) \right]$$

Now the objective is to construct  $\pi^*$ , an optimal policy such that

$$\pi^* = \arg \min_{\pi \in \Pi} J_\pi(x).$$

$\pi^*$  is not a function of the state, thus  $\pi^*$  has to work for all possible  $x$ .

## 3.2 PRINCIPLE OF OPTIMALITY

Let  $\pi^* = (\mu_0^*(\cdot), \mu_1^*(\cdot), \dots, \mu_{N-1}^*(\cdot))$  be an optimal policy.

$$\underset{(X_{i+1}, W_i | x_i=x)}{E} \left[ g_N(x_N) + \sum_{k=i}^{N-1} g_k(x_k, \mu_k^*(x_k), w_k) \right]$$

where  $X_{i+1} := (x_{i+1}, \dots, x_N)$  and  $W_i := (w_i, \dots, w_{N-1})$ . Then the **truncated** policy  $(\mu_i^*(\cdot), \mu_{i+1}^*(\cdot), \dots, \mu_{N-1}^*(\cdot))$  is optimal for this problem.

## 3.3 DPA

**Theorem 1.** For any initial state  $x \in \mathcal{S}_0$ , the optimal cost  $J^*(x)$  is equal to  $J_0(x)$  given the following recursive algorithm:

**Initialization**

$$J_N(x) = g_N(x), \quad \forall x \in \mathcal{S}_N$$

**Recursion**

$$J_k(x) := \min_{u \in \mathcal{U}_k(x)} \underset{(w_k | x_k=x, u_k=u)}{E} [g_k(x_k, u_k, w_k) + J_{k+1}(f_k(x_k, u_k, w_k))]$$

furthermore, if for each  $k$  and  $x \in \mathcal{S}_k$ ,  $u^* =: \mu_k^*(x)$  minimizes the recursion equation, the policy  $\pi^* = (\mu_0^*(\cdot), \mu_1^*(\cdot), \dots, \mu_{N-1}^*(\cdot))$  is optimal.

## 3.4 CONVERTING NON-STANDARD PROBLEMS TO THE STANDARD FORM

### 3.4.1 TIME LAGS

Assume the dynamics have the following form:

$$x_{k+1} = f_k(x_k, x_{k-1}, u_k, u_{k-1}, w_k)$$

- Let  $y_k := x_{k-1}$ ,  $s_k := u_{k-1}$  and the augmented state vector  $\tilde{x}_k := (x_k, y_k, s_k)$ .
- The dynamics of the augmented state then become

$$\tilde{x}_{k+1} = \begin{bmatrix} x_{k+1} \\ y_{k+1} \\ s_{k+1} \end{bmatrix} = \begin{bmatrix} f_k(x_k, y_k, u_k, s_k, w_k) \\ x_k \\ u_k \end{bmatrix} =: \tilde{f}_k(\tilde{x}_k, u_k, \omega_k)$$

### 3.4.2 CORRELATED DISTURBANCES

If the disturbance is correlated across time (colored noise) it can be modelled as follows:

$$\begin{aligned} w_k &= C_k y_{k+1} \\ y_{k+1} &= A_k y_k + \xi_k \end{aligned}$$

where  $A_k, C_k$  are given and  $\xi_k$ ,  $k = 0, \dots, N-1$  are independent random variables.

- Let the augmented state vector  $\tilde{x}_k := (x_k, y_k)$ . Note that now  $y_k$  must be observed at time  $k$ , which can be done using a state estimator.
- $y_k$  is an internal state of the filter that generates the noise realization  $w_k$
- The dynamics of the augmented state vector then become:

$$\tilde{x}_{k+1} = \begin{bmatrix} x_{k+1} \\ y_{k+1} \end{bmatrix} = \begin{bmatrix} f_k(x_k, u_k, C_k(A_k y_k + \xi_k)) \\ A_k y_k + \xi_k \end{bmatrix} =: \tilde{f}_k(\tilde{x}_k, u_k, \xi_k)$$

- When augmenting the state, the cost function becomes increasingly complex  $\rightarrow$  **curse of dimensionality!**

### 3.4.3 FORECASTS

Each time period has its own forecast that reveals the probability distribution of  $w_k$  and possibly of future disturbances.

At the beginning of each period  $k$ , we receive a prediction  $y_k$  that  $w_k$  will attain a probability distribution out of a given finite collection of distributions  $\{p_{w_k|y_k}(\cdot|1), p_{w_k|y_k}(\cdot|2), \dots, p_{w_k|y_k}(\cdot|m)\}$ . In particular, we receive a forecast that  $y_k = i$  and thus  $p_{w_k|x_k}(\cdot, |i)$  is used to generate  $w_k$ . Furthermore the forecast itself has a given a priori probability distribution, namely

$$y_{k+1} = \xi_k$$

where  $\xi_k$  are independent random variables taking value  $i \in \{1, 2, \dots, m\}$  with probability  $p_{\xi_k}(i)$ .

- Let the augmented state vector  $\tilde{x}_k = (x_k, y_k)$ . Since the forecast  $y_k$  is known at time  $k$  we still have perfect information.
- We defined our new disturbance as  $\tilde{w}_k := (w_k, \xi_k)$  with the distribution

$$\begin{aligned} p(\tilde{w}_k|\tilde{x}_k, u_k) &= p(w_k, \xi_k|x_k, y_k, u_k) \\ &= p(w_k|x_k, y_k, u_k, \xi_k)p(\xi_k|x_k, y_k, u_k) \\ &= p(w_k|x_k)p(\xi_k) \end{aligned}$$

- The dynamics become

$$\tilde{x}_{k+1} = \begin{bmatrix} x_{k+1} \\ y_{k+1} \end{bmatrix} = \begin{bmatrix} f_k(x_k, u_k, w_k) \\ \xi_k \end{bmatrix} =: \tilde{f}_k(\tilde{x}_k, u_k, \tilde{w}_k)$$

- The DPA becomes:

#### Initialization

$$J_N(\tilde{x}) = J_N(x, y) = g_N(x), \quad x \in S_N, y \in \{1, \dots, m\}$$

#### Recursion

$$J_k(\tilde{x}) = J_k(x, y) = \min_{u \in \mathcal{U}(x)(w_k|y_k=y)} E \left[ g_k(x, u, w_k) + \sum_{i=1}^m p_{\xi_k}(i) J_{k+1}(f_k(x, u, w_k), i) \right] \\ \forall x \in S_k, \forall y \in \{1, \dots, m\}, \forall k = N-1, \dots, 0$$

## 4 INFINITE HORIZON PROBLEM

We are dealing with a linear time-invariant system:

$$x_{k+1} = f(x_k, u_k, w_k), \quad x_k \in \mathcal{S}, \quad u_k \in \mathcal{U}(x_k), \quad w_k \sim p_{w|x,u}, \quad k = 0, \dots, N-1$$

The control inputs are generated by an admissible policy  $\pi \in \Pi$ :

$$\pi = (\mu_0(\cdot), \mu_1(\cdot), \dots, \mu_{N-1}(\cdot))$$

such that

$$\mu_k = \mu_k(x_k), \quad u_k \in \mathcal{U}(x_k), \quad \forall x_k \in \mathcal{S}, \quad \forall k$$

The cost is a function of time-invariant stage costs:

$$\sum_{k=0}^{N-1} g(x_k, u_k, w_k)$$

The expected closed loop cost of starting at  $x$  associated with policy  $\pi \in \Pi$ :

$$J_\pi(x) = E_{(X_1, W_0|x_0=x)} \left[ \sum_{k=0}^{N-1} g(x_k, \mu_k(x_k), w_k) \right]$$

**Objective:**

$$\pi^* = \arg \min_{\pi \in \Pi} J_\pi(x)$$

Analyse behaviour when  $N \rightarrow \infty$ :

$$J_N(x) = 0 \quad \forall x \in \mathcal{S}$$

$$J_k(x) = \min_{u \in \mathcal{U}(x)(w|x=x, u=u)} E [g(x, u, w) + J_{k+1}(f(x, u, w))], \quad \forall x \in \mathcal{S}, \quad \forall k = N-1, \dots, 0$$

By index substitution:  $l := N - k$  and  $V_l(\cdot) := J_{N-l}(\cdot)$ :

$$V_0(x) = 0 \quad \forall x \in \mathcal{S}$$

$$V_l(x) = \min_{u \in \mathcal{U}(x)(w|x=x, u=u)} E [g(x, u, w) + V_{l-1}(f(x, u, w))], \quad \forall x \in \mathcal{S}, \quad \forall l = 1, \dots, N$$

Now assume that for each  $x \in \mathcal{S}$  the sequence  $V_l(x)$  approaches a certain value as  $N \rightarrow \infty$ :

$$J(x) = \min_{u \in \mathcal{U}(x)(w|x=x, u=u)} E [g(x, u, w) + V_{l-1}(f(x, u, w))], \quad \forall x \in \mathcal{S} \quad \text{Bellman Equation (BE)}$$

#### 4.1 THE STOCHASTIC SHORTEST PATH (SSP) PROBLEM

Consider a finite state, time-invariant system:

$$\begin{aligned} x_{k+1} &= w_k, \quad x_k \in \mathcal{S} \\ \Pr(w_k = j | x_k = i, u_k = u) &= P_{ij}(u), \quad u \in \mathcal{U}(i) \end{aligned}$$

The expected closed loop cost of starting at  $i$  associated with policy  $\pi$  becomes:

$$J_\pi(i) = \mathbb{E}_{X_1, W_0 | x_0 = i} \left[ \sum_{k=0}^{N-1} g(x_k, \mu_k(x_k), w_k) \right]$$

We **assume** that there exists a cost-free termination state, which we designate as state 0. In particular, there are  $n + 1$  states with  $\mathcal{S} = \{0, 1, \dots, n\}$  where

$$P_{00}(u) = 1 \text{ and } g(0, u, 0) = 0, \quad \forall u \in \mathcal{U}(0)$$

The objective is then:

$$\pi^* = \arg \min_{\pi \in \Pi} J_\pi(i)$$

**Definition 14.** A policy is **stationary** if it is the same for all times such that  $\pi = (\mu(\cdot), \mu(\cdot), \dots, \mu(\cdot))$  which is written just as  $\mu$ .

**Definition 15.** A stationary policy is said to be **proper** if, when using this policy there exists an integer  $m$  such that

$$\Pr(x_m = 0 | x_0 = i) > 0$$

If a policy is not proper it is said to be **improper**.

Further it is **assumed** that there exists at least one proper policy  $\mu \in \Pi$ . Furthermore, for every improper policy  $\mu'$ , the corresponding cost function  $J_{\mu'}$  is infinity for at least one state  $i \in \mathcal{S}$ .

Based on that it can be proven that the final state will be reached with probability 1:

$$\begin{aligned} \Pr(x_m = 0 | x_0 = i) &= \alpha > 0 \\ \Pr(x_m \neq 0 | x_0 = i) &= 1 - \alpha, \quad i \in \mathcal{S} \setminus \{0\} \\ \Pr(x_{2m} \neq 0 | x_0 = i) &= \Pr(x_{2m} \neq 0, x_m \neq 0 | x_0 = i) \\ &= \underbrace{\Pr(x_{2m} \neq 0 | x_m \neq 0, x_0 = i)}_{1-\alpha} \underbrace{\Pr(x_m \neq 0 | x_0 = i)}_{1-\alpha} \\ \Rightarrow \Pr(x_{2m} = 0 | x_0 = i) &= 1 - (1 - \alpha)^2 \\ &\Rightarrow \lim_{N \rightarrow \infty} \Pr(x_N = 0 | x_0 = 1) = 1 \end{aligned}$$

#### 4.2 THEOREM 4.1

1. Given any initial conditions  $V_0(1), \dots, V_0(n)$ , the sequence  $V_l(i)$  generated by the iteration

$$V_{l+1}(i) = \min_{u \in \mathcal{U}} \left( q(i, u) + \sum_{j=1}^n P_{ij}(u) V_l(j) \right), \quad \forall i \in \mathcal{S} \setminus \{0\}$$

where

$$q(i, u) := \mathbb{E}_{(x | x=i, u=u)} [g(x, u, w)]$$

2. The optimal costs satisfy the Bellman Equation:

$$J^*(i) = \min_{u \in \mathcal{U}} \left( q(i, u) + \sum_{j=1}^n P_{ij}(u) J^*(j) \right) \quad \forall i \in \mathcal{S} \setminus \{0\}$$

3. The solution to the BE is unique
4. The minimizing  $u$  for each  $i \in \mathcal{S} \setminus \{0\}$  of the BE gives an optimal policy, which is proper.

See lecture 4 for an intuition why this is true.