



# Stochastic Decision Theory III

APP MATH 3020 & 7090



Giang Nguyen  
(bringing together notes of  
Joshua Ross, David Green, and many others)



Discipline of Applied Mathematics  
School of Mathematical Sciences  
University of Adelaide



Last updated September 27, 2018

© 2017 Giang Nguyen  
All rights reserved

Publishing Company, Adelaide, SA.

Printed in the World

The paper used in this publication may meet the minimum requirements of the American National Standard for Information Sciences — Permanence of Paper for Printed Library Materials, ANSI Z39.48–1984.

10 09 08 07 06 05 04 03 02 01 15 14 13 12 11 10 9

First edition: July, 2017

# Contents

<b>Contents</b>	<b>i</b>
<b>Preface</b>	<b>iv</b>
<b>1 Introduction and Revision</b>	<b>1</b>
1.1 Three exemplar problems . . . . .	1
1.2 Probability . . . . .	4
Basic Probability Elements . . . . .	5
Probability Distributions . . . . .	8
Discrete distributions . . . . .	9
Continuous distributions . . . . .	11
Discrete Joint Distributions . . . . .	11
Continuous Joint Distributions . . . . .	11
Expectation . . . . .	12
1.3 Convex and Quasi-Convex Functions, and Nonlinear Optimi- sation . . . . .	15
Convex Combinations and Sets . . . . .	15
Concave and Convex Functions . . . . .	16
Quasi-Concave and Quasi-Convex Functions . . . . .	19
1.4 Discrete-Time Markov Chains . . . . .	23
Classification of States . . . . .	26
Equilibrium Behaviour of DTMCs . . . . .	27

1.5	Linear Programming . . . . .	30
	MATLAB linprog.m . . . . .	36
	Duality Theory . . . . .	40
<b>2</b>	<b>Stochastic Linear Programming</b>	<b>43</b>
2.1	Motivation . . . . .	43
2.2	General Formulation . . . . .	53
2.3	Naïve DEP . . . . .	54
2.4	Recourse DEP . . . . .	54
	Recourse DEP: Examples . . . . .	54
	Recourse DEP: General Formulation . . . . .	62
	Recourse DEP: Multi-stage Recourse Problem . . . . .	64
	Stochastic Linear Programs with Recourse . . . . .	66
	Recourse DEP: Properties . . . . .	70
	Recourse DEP: Induced Constraints . . . . .	72
2.5	Chance Constrained DEP . . . . .	79
	Chance Constrained DEP: Properties . . . . .	82
2.6	Lagrange Multipliers, LPs, and Duality . . . . .	86
2.7	Dual Decomposition Method . . . . .	91
<b>3</b>	<b>Markov Decision Chains</b>	<b>107</b>
3.1	Motivation . . . . .	108
3.2	Elements of a Markov decision chain . . . . .	111
	Policies . . . . .	112
3.3	Finite Horizon Problems . . . . .	114
3.4	Indefinite Horizon Problems . . . . .	129
	Stochastic Shortest Path MDP (SPMDP) . . . . .	129
3.5	Infinite Horizon Problems . . . . .	134
	Infinite Horizon Average Value (Cost) MDPs . . . . .	134
	Discounted MDPs . . . . .	141
3.6	Positive Programming & The Value Iteration Algorithm . . . . .	147
	Positive Programming . . . . .	147
	Characterisation of the optimal policy . . . . .	148

	The Value Iteration Algorithm . . . . .	152
3.7	Negative Programming & Optimal Stopping . . . . .	154
	Characterisation of the Optimal Policy . . . . .	155
	Optimal Stopping Over a Finite Horizon . . . . .	156
	Optimal Stopping Over an Infinite Horizon . . . . .	159
	<b>Bibliography</b>	<b>161</b>

# Preface

People make decisions everyday:

- whether to take an umbrella to work;
- to take an available park for their car or continue to search for a better one;
- which of several possible methods to implement to attempt to save a species from extinction; and,
- which people in the population to give a vaccine to.

All of these decisions are being made under uncertainty:

- there exists a certain chance of rain today;
- a certain chance all the remaining car parks are occupied;
- uncertainty about how many individuals of the species exist and how they will respond to each of the possible interventions; and,
- the actual dynamics of the infection and the uptake of the vaccine by the population.

This course will focus on formulating problems of this type in a mathematical framework and provide methods for making the best decision possible taking into account the uncertainty. Three topics will be covered:

- I. Stochastic Linear Programming;
- II. Markov Decision Chains (MDCs); and,
- III. Hidden Markov Chains (HMCs).

This course is interesting as it brings together techniques from optimisation and applied probability, and uses these to address problems which have all the features of real-world decision problems.

# Chapter 1

## Introduction and Revision

### §1.1 THREE EXEMPLAR PROBLEMS

Let us consider three problems which highlight the topics of this course.

#### 1. Stochastic Linear Programming Example:

Consider a manufacturer which can make two products using two resources. The amount of each resource required to make each product, the demand for each product, the price of each product, and the supply and costs of each resource are provided in the Table 1.1.

	Resource 1	Resource 2	Demand	Price
Product 1	2	1	10	9
Product 2	3	2	10	12
Supply	50	40		
Cost	2	2		

Table 1.1: Example 1 data.

Given this information, in Optimisation and Operations Research II, we



would have formulated the *Linear Program*:

$$\begin{array}{ll}
 \text{maximise} & z = 9x_1 + 12x_2 - 6x_1 - 10x_2 \\
 \text{such that} & x_1 \geq 10 \quad \text{---fixed demand \#1} \\
 & x_2 \geq 10 \quad \text{---fixed demand \#2} \\
 & 2x_1 + 3x_2 \leq 50 \\
 & x_1 + 2x_2 \leq 40 \\
 & x_1, x_2 \geq 0
 \end{array}$$

where  $x_i, i = 1, 2$  is the number of product  $i$  produced.

However, if we consider the real-world scenario, it is typical for the manufacturer to be uncertain as to *demand* prior to producing the goods. For example, consider a shop and customers. Hence,  $D_i, i = 1, 2$  – the demand for product  $i$  – is really a random variable. Furthermore, other parts of the problem may also be similarly unknown, and stochastic.

Consequently, we might want to consider the “mathematical program”:

$$\begin{array}{ll}
 \text{'maximise'} & z = 9x_1 + 12x_2 - 6x_1 - 10x_2 \\
 \text{such that} & \text{"}x_1 \leq D_1\text{"} \\
 & \text{"}x_2 \geq D_2\text{"} \\
 & 2x_1 + 3x_2 \geq 50 \\
 & x_1 + 2x_2 \leq 40 \\
 & x_1, x_2 \geq 0
 \end{array}$$

but note that  $z$  is now a random variable, so what does “maximise” mean in this context? Furthermore, what does “ $x_i \geq D_i$ ” mean?

Given a realisation of  $(D_1, D_2)$ , then we could choose decision variables  $(x_1, x_2)$  to maximise the realisation of  $z$ , but what about the situation where we wish / need to make a decision now, before the realisation?

We will consider this problem; we’ll assume we know the distribution function of  $(D_1, D_2)$  (the random components), and that these random

variables are independent of the decision variables. We will focus on how to formulate so-called *deterministic equivalent problems*, which may (potentially) be solved using existing mathematical programming techniques.

## 2. Markov Decision Chains Example:

Eleanor is looking for a parking space on their way to a concert. Each parking space is unoccupied with probability (w.p.)  $p$  independently of whether other parking spaces are occupied or not. She cannot see if other parking spaces are occupied (due to a thick fog!), and hence can only observe the occupancy of the parking space immediately in front of her car.

If she parks  $s$  spaces from the entrance, she will miss the first  $s$  minutes of the concert,  $s = 0, 1, \dots$ . If she passes the entrance without having parked, she will miss  $D$  minutes.

**What is the optimal policy to adopt?** That is, if the driver finds a free park  $\ell \in \{s, s-1, \dots, 0\}$  spaces from the entrance, should she park there?

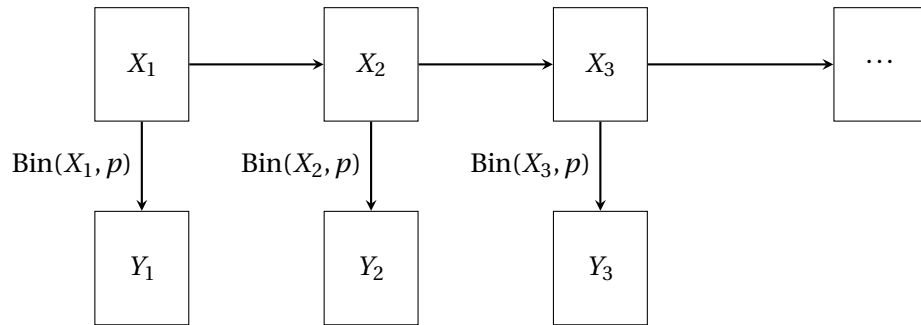
Some of the key features of these problems are:

- a Markov chain – here very simple, Discrete-time Markov chain (DTMC) with state space  $\mathcal{S} = \{0, 1\}$ ,  $t = s$  – number of parking spaces until entrance =  $0, 1, \dots, s+1$ .
- Actions – if full, no choice, but if empty, we may choose to park or not (state dependent).
- Costs – upon parking at  $t$ , cost is  $s - t$ , for  $t < s+1$ , and costs  $D$  otherwise.

## 3. Hidden Markov Chain Example:

Consider a DTMC  $\mathcal{X} = \{X_n, n = 0, 1, \dots\}$  evolving in time,  $n = 0, 1, \dots$ , and assume that at each time step we can only observe a partial observation,  $Y_n$ , of the process.

For example, at each time step we might have  $Y_n \sim \text{Bin}(X_n, p)$ .



This is a common situation in the real-world. For example, consider an epidemic – the spread of an infectious disease – where  $X_n$  is the actual number of people infected in a city on day  $n$ . Each of these individuals is only going to be recorded (known to government authorities) with some probability, say  $p$ ; this corresponds to a probability of displaying symptoms and seeking medical attention.

In such a situation, an obvious question is, how many people are really infected in the city on each day, given I only observe a realisation of  $(Y_n)$ ? We might additionally be interested in estimating the parameters of the underlying DTMC given these partial observations, or even wish to determine the best use of an antiviral. We will consider questions of this type.

## §1.2 PROBABILITY

This section follows **very closely** the reference [4]. Another excellent reference is [1].

As we are dealing with random variables (rvs) and discrete-time Markov chains (DTMCs) in this course, we will spend a little bit of time reviewing the necessary elements of basic probability and DTMCs.

### BASIC PROBABILITY ELEMENTS

Here we introduce the basic fundamental elements of probability theory.

**Defn 1.1.** The **sample space**  $\Omega$  of a random experiment is the set of all possible outcomes of the experiment.

**Example 1.** Tossing a coin:  $\Omega = \{H, T\}$ .  
Rolling a dice:  $\Omega = \{1, 2, \dots, 6\}$ .

**Defn 1.2.** An **event** is a subset of the sample space  $\Omega$  to which a probability can be assigned.

**Example 2.** Tossing a coin – getting a head:  $\{H\}$ .  
Rolling a dice – getting greater than 3:  $\{4, 5, 6\}$ .

**Defn 1.3.** A  $\sigma$ -algebra  $\mathcal{F}$  on  $\Omega$  is a collection of subsets of  $\Omega$  that satisfies:

1.  $\Omega \in \mathcal{F}$ ,
2. If  $A \in \mathcal{F}$  then also  $A^c \in \mathcal{F}$ ,
3. If  $A_1, A_2, \dots \in \mathcal{F}$ , then  $\cup_n A_n \in \mathcal{F}$ .

**Defn 1.4.** A probability (or measure)  $P$  is a function which assigns a number between 0 and 1 to each event, and which satisfies the following rules:

1.  $0 \leq P(A) \leq 1$ .
2.  $P(\Omega) = 1$ .
3. For any sequence  $A_1, A_2, \dots$  of disjoint events we have

$$P(\cup_i A_i) = \sum_i P(A_i).$$

For completeness, we note the following.

**Defn 1.5.** The triple  $(\Omega, \mathcal{F}, P)$  is called a **probability space**.

How do we evaluate probabilities of particular events of interest? In other words, what are the basic properties of a probability? We now consider such basic properties.

**Theorem 1.1.** Let  $A$  and  $B$  be events and  $P$  a probability. Then,

1.  $P(\emptyset) = 0$ ,
2. if  $A \subset B$ , then  $P(A) \leq P(B)$ ,
3.  $P(A^c) = 1 - P(A)$ ,
4.  $P(A \cup B) = P(A) + P(B) - P(A \cap B)$ .

**Defn 1.6.** The **conditional probability** of  $A$  given  $B$  (with  $P(B) \neq 0$ ) is defined as:

$$P(A|B) = \frac{P(A \cap B)}{P(B)}.$$

**Example 3.** Consider rolling 2 dice without looking. If someone informs you that the sum is 10, what is the probability that precisely one 6 was rolled?

Let  $B$  = event that sum is 10 =  $\{(4, 6), (5, 5), (6, 4)\}$  and let  $A$  = event that one 6 is rolled =  $\{(1, 6), (2, 6), \dots, (5, 6), (6, 1), \dots, (6, 5)\}$ . Then we have  $A \cap B = \{(4, 6), (6, 4)\}$ . Therefore,

$$P(A|B) = \frac{2/36}{3/36} = \frac{2}{3},$$

as all outcomes equally likely.

**Theorem 1.2. Product Rule.** Let  $A_1, \dots, A_n$  be a sequence of events with  $P(A_1, \dots, A_n) > 0$ . Then,

$$P(A_1, \dots, A_n) = P(A_1)P(A_2|A_1)P(A_3|A_1, A_2) \dots P(A_n|A_1, \dots, A_{n-1}).$$

Note,  $P(A_1, \dots, A_n) = P(A_1 \cap A_2 \cap \dots \cap A_n)$ .

**Example 4.** Do you think 2 of us in this class share the same birthday (i.e. day of the year)? What do you think the probability is?

Let  $A_i$  be the event that the first  $i$  people have different birthdays,  $i = 1, 2, \dots, n$ , where  $n$  is the class size. Note that  $A_n \subset A_{n-1} \subset \dots \subset A_2 \subset A_1$ . Therefore  $A_n = A_1 \cap A_2 \cap \dots \cap A_n$ , and thus by the product rule we have

$$P(A_n) = P(A_1)P(A_2|A_1)P(A_3|A_2) \dots P(A_n|A_{n-1}).$$

Now,  $P(A_k|A_{k-1}) = \frac{365-k+1}{365}$ . Thus

$$P(A_n) = \frac{364 \times 363 \times \dots \times (365 - n + 1)}{365 \times 365 \times \dots \times 365}, \quad n \geq 1.$$

Note, for  $n = 23$  (approximately our class size) we have  $P(A_{23}) < 1/2$  ( $\approx 1/2$ ).

**Theorem 1.3. Law of Total Probability.** *Let  $A$  be an event and let  $B_1, B_2, \dots, B_n$  be a partition of  $\Omega$ . Then*

$$P(A) = \sum_{i=1}^n P(A|B_i)P(B_i).$$

**Defn 1.7.** *The events  $A_1, A_2, \dots$ , are said to be **independent** if for any  $k$  and any choice of distinct indices  $i_1, \dots, i_k$ ,*

$$P(A_{i_1} \cap A_{i_2} \cap \dots \cap A_{i_k}) = P(A_{i_1})P(A_{i_2}) \dots P(A_{i_k}).$$

## PROBABILITY DISTRIBUTIONS

A formal definition of a rv is the following.

**Defn 1.8.** *A **random variable** is a function from the sample space  $\Omega$  to  $\mathbb{R}$ .*

More precisely, we have the following definition.

**Defn 1.9.** Let  $(\Omega, \mathcal{F}, P)$  be a probability space. A **random variable**  $X$  is a numeric function on  $\Omega$  that satisfies

$$\{X \leq x\} \in \mathcal{F},$$

for all  $x \in \mathbb{R}$ . Here  $\{X \leq x\}$  is an abbreviation of the set  $\{\omega \in \Omega : X(\omega) \leq x\}$ .

Let  $X$  be a rv. We will often desire the probabilities of events, such as  $\{X = x\}$  and/or  $\{a < X \leq b\}$ . If we can specify all such probabilities involving  $X$ , then we say we have determined the probability distribution of  $X$ .

**Defn 1.10.** The **cumulative distribution function** (cdf) of a random variable  $X$  is the function  $F : \mathbb{R} \rightarrow [0, 1]$  defined by

$$F(x) = P(X \leq x), \quad x \in \mathbb{R}.$$

## DISCRETE DISTRIBUTIONS

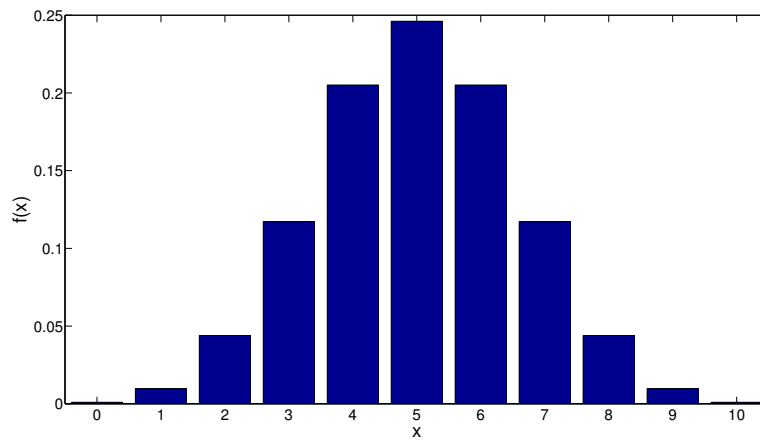
**Defn 1.11.** A rv  $X$  is said to have a **discrete distribution** if  $P(X = x_i) > 0$ ,  $i = 1, 2, \dots$ , such that  $\sum_i P(X = x_i) = 1$ . The **probability mass function** (pmf) of  $X$  is the function  $f$  defined by  $f(x) = P(X = x)$ .

**Example 5.** The *binomial distribution* models the number of successes from  $n$  trials, where success occurs wp  $p$  and all trials are independent. As a fixed example, think of the number of heads from  $n$  tosses of a fair coin, so  $p = 1/2$ . The binomial pmf is

$$f(x) = \binom{n}{x} p^x (1-p)^{n-x}, \quad x \in \{0, 1, \dots, n\}.$$



A plot of this pmf with  $n = 10$  and  $p = 1/2$  is in Figure 1.1.



**Figure 1.1:** Binomial pmf with  $n = 10$  and  $p = 1/2$ .

## CONTINUOUS DISTRIBUTIONS

**Defn 1.12.** A random variable  $X$  with cdf  $F(x)$  is said to have a **continuous distribution** if there exists a positive function  $f$ , with total integral 1, such that for  $a < b$ ,

$$P(a < X \leq b) = F(b) - F(a) = \int_a^b f(x) dx.$$

The **probability density function** (pdf) of  $X$  is the function  $f$ .

**Example 6.** An important continuous distribution is the exponential distribution. The holding time until a transition in a continuous-time Markov chain has an exponential distribution, where  $\lambda$  is the rate at which the chain leaves the current state. The pdf is

$$f(x) = \lambda e^{-\lambda x}, \quad x \geq 0$$

where  $\lambda$  is a fixed parameter. A plot of this pdf with  $\lambda = 1/2$  is in Figure 1.2.

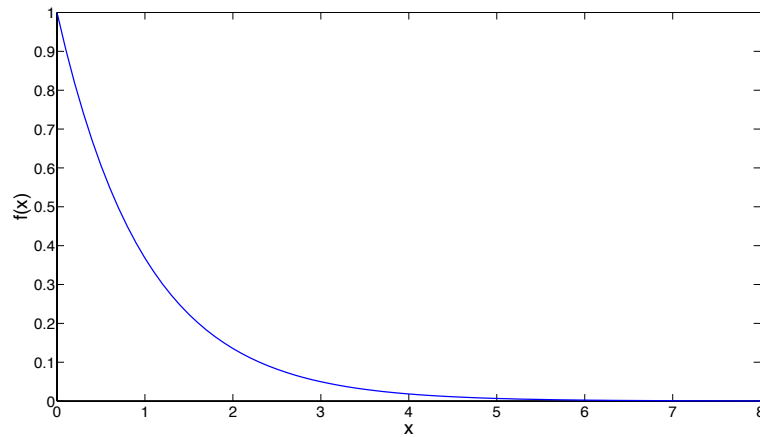
## DISCRETE JOINT DISTRIBUTIONS

**Defn 1.13** (Discrete Joint pmf). The joint pmf  $f$  of discrete rvs  $X_1, X_2, \dots, X_n$  is given by the function

$$f(x_1, x_2, \dots, x_n) = P(X_1 = x_1, X_2 = x_2, \dots, X_n = x_n).$$

Sometimes we may write  $f_{X_1, \dots, X_n}$  instead of  $f$ , to be emphatic about which random vector the pmf pertains to. Alternatively, if  $\mathbf{X} = (X_1, \dots, X_n)$  has been defined, we may write  $f_{\mathbf{X}}$ .

## CONTINUOUS JOINT DISTRIBUTIONS



**Figure 1.2:** Exponential PDF with  $\lambda = 1$ .

**Defn 1.14** (Continuous Joint pdf). *Continuous random variables  $X_1, X_2, \dots, X_n$  are said to have a joint pdf  $f$  if*

$$P(a_1 < X_1 \leq b_1, \dots, a_n < X_n \leq b_n) = \int_{a_1}^{b_1} \dots \int_{a_n}^{b_n} f(x_1, x_2, \dots, x_n) dx_1 \dots dx_n \quad (1.1)$$

*for all  $a_1, \dots, b_n$ .*

## EXPECTATION

**Defn 1.15** (Expectation of a Discrete Random Variable). *Let  $X$  be a discrete rv with pmf  $f$ . The **expectation** (or expected value) of  $X$ , denoted  $\mathbf{E}(X)$ , is defined as*

$$\mathbf{E}(X) = \sum_x xP(X = x) = \sum_x xf(x),$$

*provided that this sum is absolutely convergent.*

*(That is,  $\sum_x |x|f(x) < \infty$ .)*

**Example 7.** Considering the binomial distribution again, we have

$$\mathbf{E}(X) = \sum_{x=0}^n x \binom{n}{x} p^x (1-p)^{n-x} = np.$$

One way to interpret the expectation is as a long-run average. If we repeatedly draw random numbers from the pmf  $f$ , then eventually the average will converge to  $\mathbf{E}(X)$ .

**Defn 1.16** (Expectation of a Continuous Random Variable). *Let  $X$  be a continuous rv with pdf  $f$ . The **expectation** (or expected value) of  $X$ , denoted as  $\mathbf{E}(X)$ , is defined as*

$$\mathbf{E}(X) = \int_{-\infty}^{\infty} xf(x)dx,$$

*provided this integral is absolutely convergent.*

*(That is,  $\int_{-\infty}^{\infty} |x|f(x)dx < \infty$ .)*

**Example 8.** Considering the exponential distribution again, we have

$$\mathbf{E}(X) = \int_0^{\infty} x\lambda e^{-\lambda x} dx = \frac{1}{\lambda}.$$

Note, if  $X$  is a rv then a function of  $X$ , such as  $X^2$  or  $\log(X)$ , is also a rv.

**Theorem 1.4** (Expectation of a Function of a rv). *If  $X$  is discrete with pmf  $f$ , then for any real-valued function  $g$ ,*

$$\mathbf{E}(g(X)) = \sum_x g(x)f(x).$$

*Similarly, if  $X$  is continuous with pdf  $f$ , then*

$$\mathbf{E}(g(X)) = \int_{-\infty}^{\infty} g(x)f(x)dx,$$

*provided the sum/integral is absolutely convergent.*

**Theorem 1.5** (Properties of Expectation). *For any real numbers  $a$  and  $b$ , and functions  $g$  and  $h$ ,*

$$\mathbf{E}(aX + b) = a\mathbf{E}(X) + b,$$

$$\mathbf{E}(g(X) + h(X)) = \mathbf{E}(g(X)) + \mathbf{E}(h(X)).$$

**Theorem 1.6** (Jensen's Inequality). *Let  $h(x)$  be a convex function and  $X$  a random variable. Then*

$$\mathbf{E}(h(X)) \geq h(\mathbf{E}(X)).$$

*Proof.* Assignment 1.

□

### §1.3 CONVEX AND QUASI-CONVEX FUNCTIONS, AND NONLINEAR OPTIMISATION

#### CONVEX COMBINATIONS AND SETS

**Defn 1.17.** Given any finite collection of points  $\mathbf{x}_1, \dots, \mathbf{x}_m \in \mathbb{R}^n$ , a point  $\mathbf{z} \in \mathbb{R}^n$  is said to be a **convex combination** of the points  $\{\mathbf{x}_1, \dots, \mathbf{x}_m\}$  if there is some  $\lambda \in \mathbb{R}^m$  satisfying

1.  $\lambda_i \geq 0, i = 1, \dots, m$

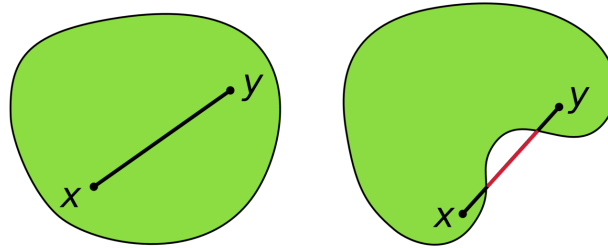
2.  $\sum_{i=1}^m \lambda_i = 1$

such that  $\mathbf{z} = \sum_{i=1}^m \lambda_i \mathbf{x}_i$ .

**Defn 1.18.** A subset  $\mathcal{S}$  of  $\mathbb{R}^n$  is **convex** if the convex combination of any two points in  $\mathcal{S}$  is also in  $\mathcal{S}$ . That is, for any  $\mathbf{x}_1, \mathbf{x}_2 \in \mathcal{S}$  and  $\lambda \in (0, 1)$ , then

$$\mathbf{x} = \lambda \mathbf{x}_1 + (1 - \lambda) \mathbf{x}_2 \in \mathcal{S}.$$

**Example 9.** The following Figure 1.3 gives an example of a convex set and a non-convex set.

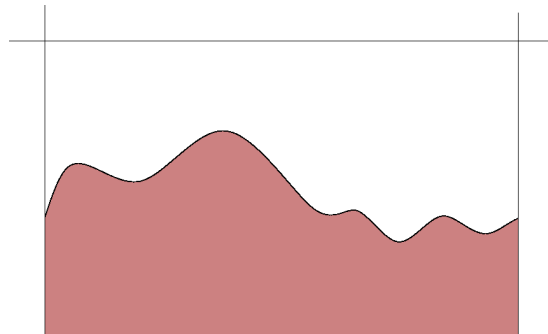


**Figure 1.3:** A convex and non-convex set, respectively.

### CONCAVE AND CONVEX FUNCTIONS

**Defn 1.19.** Let  $\mathcal{S}$  be a convex subset of  $\mathbb{R}^n$  and let  $f : \mathcal{S} \rightarrow \mathbb{R}$  be a function. The **subgraph** of  $f$ , denoted  $\text{sub}(f)$ , is the set

$$\text{sub}(f) = \{(x, y) \in \mathcal{S} \times \mathbb{R} : f(x) \geq y\}.$$

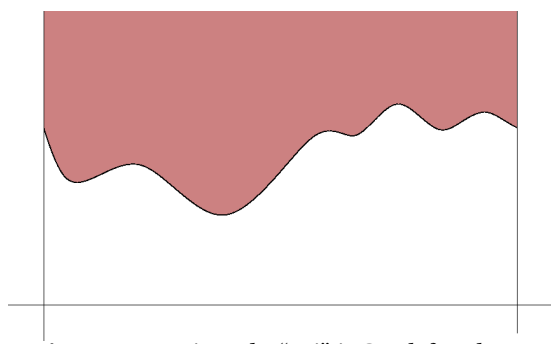


**Figure 1.4:** Subgraph

That is, the area lying below the graph of a function is its subgraph.

**Defn 1.20.** Let  $\mathcal{S}$  be a convex subset of  $\mathbb{R}^n$  and let  $f : \mathcal{S} \rightarrow \mathbb{R}$  be a function. The **epigraph** of  $f$ , denoted  $\text{epi}(f)$ , is the set

$$\text{epi}(f) = \{(x, y) \in \mathcal{S} \times \mathbb{R} : f(x) \leq y\}.$$



**Figure 1.5:** Epigraph (“epi” is Greek for above)

That is, the area lying above the graph of a function is its epigraph.

**Defn 1.21.** Let  $\mathcal{S}$  be a convex subset of  $\mathbb{R}^n$  and let  $f : \mathcal{S} \rightarrow \mathbb{R}$  be a function. We say that  $f$  is **concave** on  $\mathcal{S}$  if  $\text{sub}(f)$  is a convex set.

**Defn 1.22.** Let  $\mathcal{S}$  be a convex subset of  $\mathbb{R}^n$  and let  $f : \mathcal{S} \rightarrow \mathbb{R}$  be a function. We say that  $f$  is **convex** on  $\mathcal{S}$  if  $\text{epi}(f)$  is a convex set.

An alternative definition of concave and convex functions is available.



**Defn 1.23.** Let  $\mathcal{S}$  be a convex subset of  $\mathbb{R}^n$  and let  $f : \mathcal{S} \rightarrow \mathbb{R}$  be a function. We say that  $f$  is **concave** on  $\mathcal{S}$  if and only if  $\forall x, y \in \mathcal{S}$  and  $\lambda \in (0, 1)$  we have

$$f(\lambda x + (1 - \lambda)y) \geq \lambda f(x) + (1 - \lambda)f(y).$$

**Defn 1.24.** Let  $\mathcal{S}$  be a convex subset of  $\mathbb{R}^n$  and let  $f : \mathcal{S} \rightarrow \mathbb{R}$  be a function. We say that  $f$  is **convex** on  $\mathcal{S}$  if and only if  $\forall x, y \in \mathcal{S}$  and  $\lambda \in (0, 1)$  we have

$$f(\lambda x + (1 - \lambda)y) \leq \lambda f(x) + (1 - \lambda)f(y).$$

We also have *strict* versions of these definitions.

**Defn 1.25.** Let  $\mathcal{S}$  be a convex subset of  $\mathbb{R}^n$  and let  $f : \mathcal{S} \rightarrow \mathbb{R}$  be a function. We say that  $f$  is **strictly concave** on  $\mathcal{S}$  if and only if  $\forall x, y \in \mathcal{S}$ , with  $x \neq y$ , and  $\lambda \in (0, 1)$  we have

$$f(\lambda x + (1 - \lambda)y) > \lambda f(x) + (1 - \lambda)f(y).$$

**Defn 1.26.** Let  $\mathcal{S}$  be a convex subset of  $\mathbb{R}^n$  and let  $f : \mathcal{S} \rightarrow \mathbb{R}$  be a function. We say that  $f$  is **strictly convex** on  $\mathcal{S}$  if and only if  $\forall x, y \in \mathcal{S}$ , with  $x \neq y$ , and  $\lambda \in (0, 1)$  we have

$$f(\lambda x + (1 - \lambda)y) < \lambda f(x) + (1 - \lambda)f(y).$$

**Theorem 1.7.** *Let  $\mathcal{S}$  be a convex subset of  $\mathbb{R}^n$  and let  $f : \mathcal{S} \rightarrow \mathbb{R}$  be a function. Then*

1.  *$f$  is concave if and only if the function  $(-f)$  is convex.*
2.  *$f$  is strictly concave if and only if the function  $(-f)$  is strictly convex.*

This allows us to easily apply all results about concave functions to convex functions, and vice-versa.

There is also important consequences of convexity for optimisation, which can be essentially summarised as follows:

- if  $f$  is convex, then a local minimum is a global minimum; furthermore, it is either unique or there are infinitely many which are 'grouped together'.
- If the function  $f$  is strictly convex, then a local minimum is a global minimum and uniqueness is guaranteed.

### QUASI-CONCAVE AND QUASI-CONVEX FUNCTIONS

Whilst powerful, convexity is unfortunately a very restrictive assumption. It is often not satisfied in real-world applications. We will consider a slight generalisation of convexity, called **quasi-convexity**.

**Defn 1.27.** *Let  $\mathcal{S}$  be a convex subset of  $\mathbb{R}^n$  and let  $f : \mathcal{S} \rightarrow \mathbb{R}$  be a function. The **upper contour set** of  $f$  at  $a \in \mathbb{R}$ , denoted  $U_f(a)$ , is the set*

$$U_f(a) = \{x \in \mathcal{S} : f(x) \geq a\}.$$

**Defn 1.28.** Let  $\mathcal{S}$  be a convex subset of  $\mathbb{R}^n$  and let  $f : \mathcal{S} \rightarrow \mathbb{R}$  be a function. The **lower contour set** of  $f$ , denoted  $L_f(a)$ , is the set

$$L_f(a) = \{x \in \mathcal{S} : f(x) \leq a\}.$$

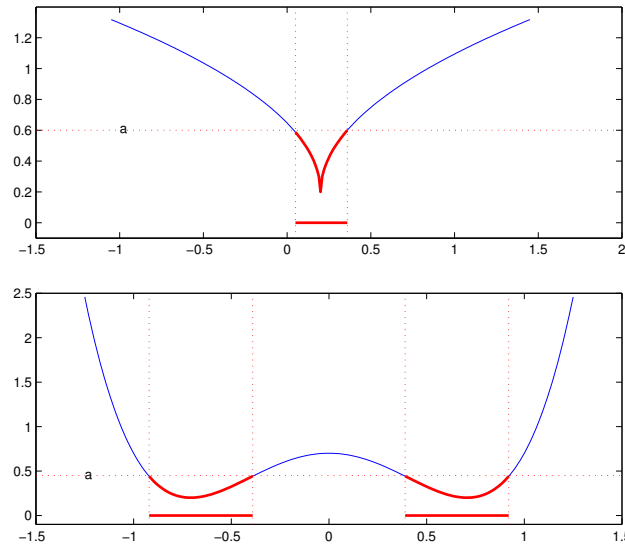
**Defn 1.29.** Let  $\mathcal{S}$  be a convex subset of  $\mathbb{R}^n$  and let  $f : \mathcal{S} \rightarrow \mathbb{R}$  be a function. We say that  $f$  is **quasi-concave** on  $\mathcal{S}$  if  $U_f(a)$  is a convex set  $\forall a \in \mathbb{R}$ .

**Defn 1.30.** Let  $\mathcal{S}$  be a convex subset of  $\mathbb{R}^n$  and let  $f : \mathcal{S} \rightarrow \mathbb{R}$  be a function. We say that  $f$  is **quasi-convex** on  $\mathcal{S}$  if  $L_f(a)$  is a convex set  $\forall a \in \mathbb{R}$ .

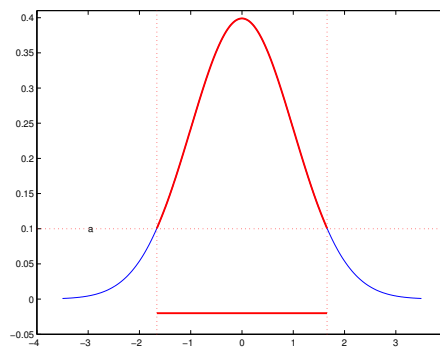
**Example 10.** Below in figure (1.6) is an example,  $f(x) = \sqrt{|x - 0.2|} + 0.2$ , of a quasi-convex function which is not convex, and also an example of a non-quasi-convex function,  $f(x) = 2(x^4 - x^2 + 0.35)$ .

An example of a quasi-concave, but not concave function, is the Gaussian pdf.

Just as was the case for concave and convex functions, a relationship exists between the value of a function at two points and the value of the function at a convex combination:



**Figure 1.6:** Top: An example of a quasi-convex, but not convex, function. Bottom: a non-quasi-convex function, respectively.



**Figure 1.7:** The Gaussian pdf, which is quasi-concave but not concave.

**Theorem 1.8.** Let  $\mathcal{S}$  be a convex subset of  $\mathbb{R}^n$  and let  $f : \mathcal{S} \rightarrow \mathbb{R}$  be a function. Then the following statements are equivalent.

1.  $f$  is quasi-concave on  $\mathcal{S}$ .
2. for all  $x, y \in \mathcal{S}$  and all  $\lambda \in (0, 1)$ ,

$$f(x) \geq f(y) \implies f(\lambda x + (1 - \lambda)y) \geq f(y).$$

3. or all  $x, y \in \mathcal{S}$  and all  $\lambda \in (0, 1)$ ,

$$f(\lambda x + (1 - \lambda)y) \geq \min\{f(x), f(y)\}.$$

**Defn 1.31.** Let  $\mathcal{S}$  be a convex subset of  $\mathbb{R}^n$  and let  $f : \mathcal{S} \rightarrow \mathbb{R}$  be a function. We say that  $f$  is **strictly quasi-concave** on  $\mathcal{S}$  if  $\forall x, y \in \mathcal{S}$ , with  $x \neq y$ , and  $\lambda \in (0, 1)$  we have

$$f(\lambda x + (1 - \lambda)y) > \min\{f(x), f(y)\}.$$

**Defn 1.32.** Let  $\mathcal{S}$  be a convex subset of  $\mathbb{R}^n$  and let  $f : \mathcal{S} \rightarrow \mathbb{R}$  be a function. We say that  $f$  is **strictly quasi-convex** on  $\mathcal{S}$  if  $\forall x, y \in \mathcal{S}$ , with  $x \neq y$ , and  $\lambda \in (0, 1)$  we have

$$f(\lambda x + (1 - \lambda)y) < \max\{f(x), f(y)\}.$$

Part of the flexibility of quasi-convexity over convexity is portrayed in the following theorem.

**Theorem 1.9.** *Let  $\mathcal{S}$  be a convex subset of  $\mathbb{R}^n$  and let  $f : \mathcal{S} \rightarrow \mathbb{R}$  be a quasi-concave function.*

1. *If  $\phi : \mathbb{R} \rightarrow \mathbb{R}$  is an increasing function, then the composition  $\phi \circ f$  is a quasi-concave function from  $\mathcal{S} \rightarrow \mathbb{R}$ .*
2. *In particular, any increasing transform of a concave function results in a quasi-concave function.*

The relationship between quasi-convexity and optimisation is summarised in the following theorem.

**Theorem 1.10.** *Let  $\mathcal{S}$  be a convex subset of  $\mathbb{R}^n$  and let  $f : \mathcal{S} \rightarrow \mathbb{R}$  be a quasi-concave function. Then*

1. *Any local maximum of  $f$  is a global maximum of  $f$ .*
2. *The set  $\operatorname{argmax}\{f(x) : x \in \mathcal{S}\}$  of maximisers of  $f$  on  $\mathcal{S}$  is either empty or a singleton.*

#### §1.4 DISCRETE-TIME MARKOV CHAINS

A discrete-time Markov chain (DTMC) is a stochastic process – a sequence of random variables, evolving in time, corresponding to a sequence of random variables  $X_n \in S$ , where the random variable  $X_n$  is the state of the Markov chain at time  $n$ , and  $\mathcal{S}$  is the state space of all values the Markov chain may adopt. A DTMC evolves in discrete time steps,  $n = 0, 1, 2, \dots$

The use of *chain* (as opposed to ‘process’) implies that  $X_n$  is a discrete random variable.

**Defn 1.33.**  $\mathcal{X} = \{X_n, n \geq 0\}$  is called a *discrete-time Markov chain* provided that for all  $x_0, x_1, \dots, x_n, j \in S$  and  $n \in \mathbb{N}$

$$\begin{aligned} \Pr(X_{n+1} = j | X_0 = x_0, X_1 = x_1, \dots, X_n = x_n) \\ = \Pr(X_{n+1} = j | X_n = x_n) \quad - \text{Markov property,} \end{aligned}$$

with  $\mathcal{S}$  being referred to as the *state space*, consisting of all values the Markov chain may adopt.

We will assume  $\mathcal{S}$  is countable (and often furthermore finite) so that it can be put in one-to-one correspondence with a subset of the integers  $S \subseteq \mathbb{Z}$ .

Given that the chain is current in state  $i$ , say  $X_n = i$ , then we specify the probabilities of it being in state  $j$  at the next time step  $n + 1$ :

$$\Pr(X_{n+1} = j | X_n = i) = p_{ij}(n).$$

Note that these probabilities depend on three quantities: the current state,  $i$ ; the future state,  $j$ ; and the time step  $n$ .

We will focus on *time-homogeneous* Markov chains where the transition probabilities only depend upon  $i$  and  $j$ :

$$\Pr(X_{n+1} = j | X_n = i) = p_{ij} = \Pr(X_1 = j | X_0 = i).$$

Considering these probabilities for each  $i, j \in S$ , we can form a (one-step) *transition (probability) matrix*

$$P = (p_{ij})$$

of size  $|S| \times |S|$ .

**Example 11 (Tossing a coin  $n$  times).** Let  $X_n$  be the number of heads after  $n$  tosses,  $n = 0, 1, \dots, T$ . Then we can form a *one-step probability transition*

matrix  $P$  as follows

$$P = \begin{pmatrix} \frac{1}{2} & \frac{1}{2} & 0 & \dots & 0 \\ 0 & \frac{1}{2} & \frac{1}{2} & \ddots & 0 \\ 0 & 0 & \frac{1}{2} & \ddots & 0 \\ \vdots & \vdots & \ddots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & 1 \end{pmatrix},$$

which is of dimension  $|S| \times |S|$  where  $S = \{0, 1, \dots, N\}$  being the set of possible numbers of heads from  $N$  coin tosses.

We now specify a (row) vector

$$\mathbf{p}(n) = (\Pr(X_n = 0), \Pr(X_n = 1), \dots, \Pr(X_n = N)),$$

specifying the pmf of the DTMC at time step  $n$ . Then, the  $j$ th entry of the vector  $\mathbf{p}(n+1)$  is given by

$$p_j(n+1) = \Pr(X_{n+1} = j) = \sum_i p_i(n) p_{ij} = [\mathbf{p}(n)P]_j$$

and hence  $\mathbf{p}(n+1) = \mathbf{p}(n)P$ . This specifies the evolution of the DTMC. It follows that

$$\mathbf{p}(n) = \mathbf{p}(0)P^n,$$

where  $\mathbf{p}(0)$  is the initial distribution of the DTMC.

**Example 12.** So, considering our coin example again and assuming we start at time  $n = 0$ , then we have

$$\mathbf{p}(0) = (1, 0, \dots, 0) = \mathbf{e}_1^{N+1}$$

where  $\mathbf{e}_b^a$  is a (row) vector of length  $a$  with a 1 in the  $b$ th entry and 0s elsewhere. Now, at time  $n = 1$  we have

$$\mathbf{p}(1) = \mathbf{p}(0)P = \mathbf{e}_1^{N+1} \begin{pmatrix} \frac{1}{2} & \frac{1}{2} & 0 & \dots & 0 \\ 0 & \frac{1}{2} & \frac{1}{2} & \ddots & 0 \\ 0 & 0 & \frac{1}{2} & \ddots & 0 \\ \vdots & \vdots & \ddots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & 1 \end{pmatrix} = \left( \frac{1}{2}, \frac{1}{2}, 0, \dots, 0 \right).$$



What do you think  $p(N)$  would look like, for  $N = 10$ ? (Check!)

### CLASSIFICATION OF STATES

**Defn 1.34.** If states  $i$  and  $j$  are such that  $[P^n]_{ij} > 0$  for some  $n > 0$ , we say that  $i$  **has access to**  $j$  and write  $i \rightarrow j$ .

**Defn 1.35.** We say that  $i$  and  $j$  **communicate** if  $i \rightarrow j$  and  $j \rightarrow i$ . Then, we write  $i \leftrightarrow j$ .

The relation  $i \leftrightarrow j$  is an *equivalence relation*:

- $i \leftrightarrow i$
- $i \leftrightarrow j \iff j \leftrightarrow i$
- $i \leftrightarrow j, j \leftrightarrow k \iff i \leftrightarrow k$ .

We can therefore divide  $\mathcal{S}$  into *equivalence (or communicating) classes* such that states of a communicating class communicate with all other states within that class but not with any state outside that class.

**Defn 1.36.** If there is only one communicating class ( $= S$ ), the DTMC is said to be **irreducible**.

**Defn 1.37.** A communicating class is classified as **transient** if it is possible for the process to leave that class so that it can never return.

On the other hand, a communicating class is classified as **recurrent** if once the process enters that class, it can never leave.

**Example 13** (Toss a coin  $n$  times). In the coin tossing example, we can draw a state diagram (see Figure 1.8) and use the notation of access to see what type of states and communicating classes we have.

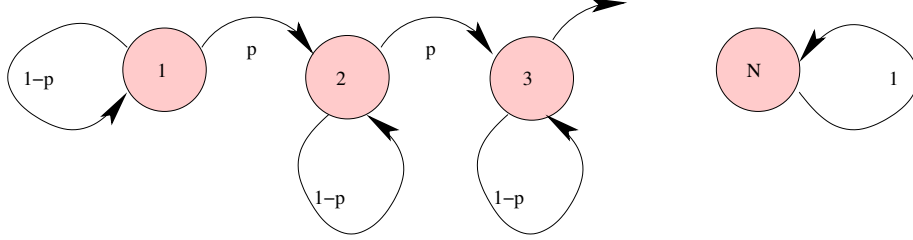


Figure 1.8: State diagram for the coin tossing example.

By Definition 1.37, each of the states  $\{0, 1, \dots, N-1\}$  is their own communicating classes which are transient; the state  $N$ , on the other hand, is a communicating class that is recurrent or absorbing.

### EQUILIBRIUM BEHAVIOUR OF DTMCs

We can easily show that the probability entries  $p_{ij}^{(m)}$  of the  $m$ -step transition probability matrix  $P^{(m)}$ , for  $m \geq 1$ , are given by

$$p_{ij}^{(m)} = \sum_{k \in S} p_{ik}^{(m-1)} p_{kj} \quad \text{for } i, j \in S. \quad (1.2)$$

For irreducible finite-state DTMCs, there is an **equilibrium behaviour**, that is,

$$\lim_{m \rightarrow \infty} p_{ij}^{(m)} = \pi_j. \quad (1.3)$$

In other words, the probability of being in state  $j$  after many steps converges to some constant value  $\pi_j$  and hence is independent of the initial state  $i$ . (Note that this property does not always occur with all Markov Chains that are not irreducible finite-state.)

Equation (1.2) and the assumption (1.3) implies that

$$\pi_j = \sum_{k \in S} \pi_k p_{kj}, \quad \text{for all } j \in S. \quad (1.4)$$

Equations (1.4) are known as the **equilibrium equations** for a Markov Chain and can be written in vector/matrix form as

$$\boldsymbol{\pi} = \boldsymbol{\pi} P, \quad (1.5)$$

where the row vector  $\boldsymbol{\pi} = (\pi_1, \pi_2, \dots)$ . If in addition,  $\boldsymbol{\pi}$  can be scaled such that  $\sum_i \pi_i = 1$ , then  $\boldsymbol{\pi}$  is known as the **equilibrium probability distribution** of the DTMC, which, as can be seen from (1.5), is a left-eigenvector of  $P$  associated with an eigenvalue of 1.

**Example 14.** Consider an irreducible four-state DTMC defined on  $\mathcal{S} = \{0, 1, 2, 3\}$  with the following transition probability matrix:

$$P = \begin{pmatrix} \frac{1}{2} & \frac{1}{2} & 0 & 0 \\ \frac{1}{4} & 0 & \frac{3}{4} & 0 \\ 0 & \frac{1}{4} & 0 & \frac{3}{4} \\ 0 & \frac{1}{4} & \frac{3}{4} & 0 \end{pmatrix}.$$

The equilibrium probability equations  $\boldsymbol{\pi} = \boldsymbol{\pi} P$  are

$$\pi_0 = \frac{1}{2}\pi_0 + \frac{1}{4}\pi_1 \quad \Rightarrow \quad \pi_0 = \frac{1}{2}\pi_1, \quad (1.6)$$

$$\pi_1 = \frac{1}{2}\pi_0 + \frac{1}{4}\pi_2 + \frac{1}{4}\pi_3, \quad (1.7)$$

$$\pi_2 = \frac{3}{4}\pi_1 + \frac{3}{4}\pi_3, \quad (1.8)$$

$$\pi_3 = \frac{3}{4}\pi_2. \quad (1.9)$$

Substituting (1.9) into (1.8), we have

$$\pi_2 = \frac{3}{4}\pi_1 + \frac{3}{4}\frac{3}{4}\pi_2 \Rightarrow \pi_2 = \frac{12}{7}\pi_1. \quad (1.10)$$

Similarly, substituting (1.10) into (1.9), we have

$$\pi_3 = \frac{3}{4}\pi_2 = \frac{3}{4}\frac{12}{7}\pi_1 = \frac{9}{7}\pi_1. \quad (1.11)$$

Then using the fact that  $\sum_{i=0}^3 \pi_i = 1$ , we obtain

$$\boldsymbol{\pi} = (\pi_0, \pi_1, \pi_2, \pi_3) = \frac{14}{63} \left( \frac{1}{2}, 1, \frac{12}{7}, \frac{9}{7} \right) = \left( \frac{7}{63}, \frac{14}{63}, \frac{24}{63}, \frac{18}{63} \right).$$

**Defn 1.38.** If a set of states  $C$  is such that

$$\sum_{j \in C} p_{ij} = 1$$

for all  $i \in C$ , then  $C$  is called a **closed set**.

**Defn 1.39.** A state  $i$  is called an **absorbing state** if  $\{i\}$  is closed.

**Defn 1.40** (Absorbing class). A class  $C$  is called an **absorbing class** if  $C$  is closed.

An irreducible DTMC effectively has a single absorbing class.

### §1.5 LINEAR PROGRAMMING

In section I of the course, on Stochastic Linear Programming, we will formulate so-called *deterministic equivalent problems* (DEPs) of stochastic linear programs. Such DEPs are commonly general (nonlinear) mathematical programs, and occasionally will be linear programs. For this reason, we review (briefly) **Linear Programming**.

A mathematical programming problem (or mathematical program):

$$\begin{aligned} &\text{minimise (or maximise)} && z = g_0(x) \\ &\text{such that} && g_i(x) = (\text{or } \leq \text{ or } \geq) b_i, i = 1, \dots, m \\ &&& x \in X \subset \mathbb{R}^n \end{aligned}$$

is called a **linear program** if  $g_0(x) = c^\top x$  and  $\{g_i(x)\} \equiv Ax$ . That is, the objective function and constraints are all linear.

All linear programs can be expressed in the *standard form*:

$$\begin{aligned} &\text{minimise} && z = \mathbf{c}^\top \mathbf{x} \\ &\text{such that} && A\mathbf{x} = \mathbf{b}, \\ &&& \mathbf{x} \geq \mathbf{0} \end{aligned}$$

where  $A \in \mathbb{R}^{m \times n}$ ,  $\mathbf{x} \in \mathbb{R}^n$ ,  $\mathbf{b} \in \mathbb{R}^m$ ,  $\text{rank}(A) = m \leq n$  (not essential) and  $b_i \geq 0$  for  $i = 1, \dots, m$ .

**Example 15 (A scheduling problem).** A manufacturing company makes three types of items: desks, chairs, and bed-frames, using metal and wood.

A single desk requires 2 hours of labour, 1 unit of metal, 3 units of wood. One chair requires 1 hour of labour, 1 unit of metal, and 3 units of wood. Making one bed frame requires 2 hours of labour, 1 unit of metal and 4 units of wood. In a given time period, there are 225 hours of labour available, 117 units of metal, and 420 units of wood. The profit on one desk is \$13, one chair is \$12, and one bed frame is \$17.

The company wants a manufacturing schedule designed, which maximises profits without violating any constraint on resource availability during that time period.

To formulate the problem, we start by reading the problem statement carefully and then tabulating the data. The data for this example is listed in Table 1.2.

	Labour (hrs)	Metal	Wood	Profit
Desk	2	1	3	13
Chair	1	1	3	12
Bedframe	2	1	4	17
Total	225	117	420	

Table 1.2: Example 15 data

Then, we identify **decision variables**, which in this case are  $x_1$ ,  $x_2$ , and  $x_3$ , where  $x_1$  is the number of desks,  $x_2$  the number of chairs, and  $x_3$  the number of bed frames, made per time period. Next, we need to determine the **objective function**, which here is to maximise the profit,  $z$ , (in dollars); so, we have:

$$\max z = 13x_1 + 12x_2 + 17x_3.$$

The constraints are

$$2x_1 + x_2 + 2x_3 \leq 225 \quad (\text{labour})$$

$$x_1 + x_2 + x_3 \leq 117 \quad (\text{metal})$$

$$3x_1 + 3x_2 + 4x_3 \leq 420 \quad (\text{wood})$$

$$x_1, x_2, x_3 \geq 0.$$

The feasible region for this LP is depicted in Figure 1.9

Should we include any other constraints on the  $x_i$ 's? If integer constraints were added (whole hours or whole units of base products), it would become an Integer Linear Program (ILP).

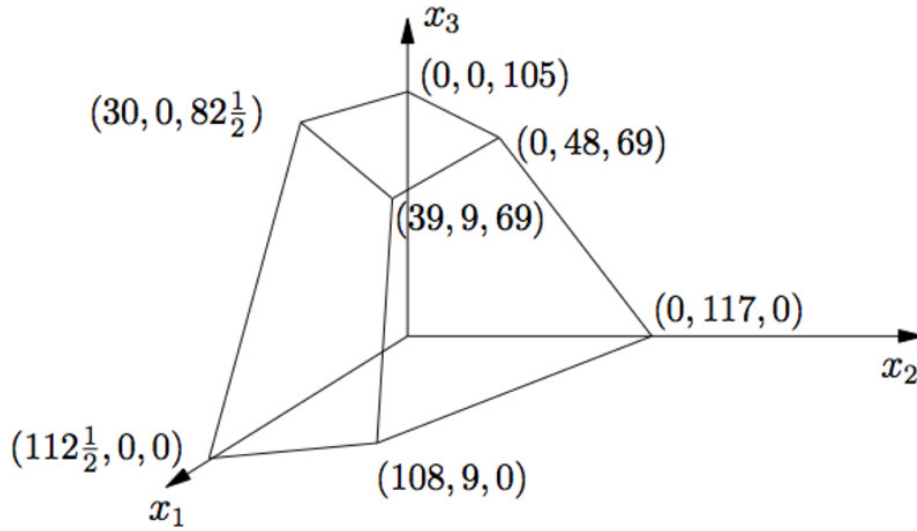


Figure 1.9: Feasible region for the LP in Example 15.

In summary, the *linear program* (LP) is

$$\begin{aligned}
 \max \quad & z = 13x_1 + 12x_2 + 17x_3 \\
 \text{subject to} \quad & 2x_1 + x_2 + 2x_3 \leq 225 \\
 & x_1 + x_2 + x_3 \leq 117 \\
 & 3x_1 + 3x_2 + 4x_3 \leq 420 \\
 & x_1, x_2, x_3 \geq 0.
 \end{aligned}$$

To write in *standard form*, we add **slack variables**  $x_4, x_5, x_6 \geq 0$  as follows

$$\begin{aligned}
 \min \quad & -z = -13x_1 - 12x_2 - 17x_3 \\
 \text{subject to} \quad & 2x_1 + x_2 + 2x_3 + x_4 = 225 \\
 & x_1 + x_2 + x_3 + x_5 = 117 \\
 & 3x_1 + 3x_2 + 4x_3 + x_6 = 420 \\
 & x_1, x_2, x_3, x_4, x_5, x_6 \geq 0.
 \end{aligned}$$

The variables  $x_4, x_5$  and  $x_6$  are known as slack variables because they actually take up the slack to convert an inequality into an equality.

We will revisit this example later to see a solution of the above LP.

**Defn 1.41.** Variables  $x$  satisfying the constraints

$$Ax = b, \quad x \geq 0$$

are called **feasible solutions**.

**Defn 1.42.** Suppose  $\text{rank}(A) = m$ , then there exists  $m$  columns of  $A$  which are linearly independent, say the first  $m$  columns. Define

$$B := (A_{:,1}, A_{:,2}, \dots, A_{:,m}),$$

then  $B$  is non-singular. We say  $B$  is a **basis**.



**Defn 1.43.** If  $B$  is a basis, then the linear system

$$B\mathbf{x}_B = \mathbf{b}$$

has a unique solution,  $\mathbf{x}_B = B^{-1}\mathbf{b}$ . Let  $\mathbf{x} = (\mathbf{x}_B, \mathbf{0})$ , then  $\mathbf{x}$  satisfies

$$A\mathbf{x} = \mathbf{b},$$

and  $\mathbf{x}$  is called a **basic solution** with respect to  $B$ .

Furthermore, if  $\mathbf{x}_B \geq \mathbf{0}$  it is called a **basic feasible solution**.

The  $\mathbf{x}_B$  are called **basic variables**; the other variables corresponding to the 0s are called **non-basic variables**.

**Example 16.** Consider the linear program

$$\begin{aligned} \text{maximise} \quad & z = 9x_1 + 12x_2 - 6x_1 - 10x_2 \\ \text{such that} \quad & x_1 \leq 10 \\ & x_2 \leq 10 \\ & 2x_1 + 3x_2 \leq 50 \\ & x_1 + 2x_2 \leq 40 \\ & x_1, x_2 \geq 0 \end{aligned}$$

where  $x_i, i = 1, 2$  is the number of product  $i$  produced. In standard form we have

$$\begin{aligned} \text{minimise} \quad & z = -3x_1 - 2x_2 \\ \text{such that} \quad & x_1 + x_3 = 10 \\ & x_2 + x_4 = 10 \\ & 2x_1 + 3x_2 + x_5 = 50 \\ & x_1 + 2x_2 + x_6 = 40 \\ & x_1, x_2, \dots, x_6 \geq 0. \end{aligned}$$

Here,

$$A = \begin{pmatrix} 1 & 0 & 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 1 & 0 & 0 \\ 2 & 3 & 0 & 0 & 1 & 0 \\ 1 & 2 & 0 & 0 & 0 & 1 \end{pmatrix},$$

and  $\mathbf{b} = (10, 10, 50, 40)^\top$ . By inspection, we can choose

$$B = (A_{:,3}, A_{:,4}, A_{:,5}, A_{:,6}) = I_{4 \times 4},$$

as a basis, and hence  $\mathbf{x} = (0, 0, 10, 10, 50, 40)^\top$  is a **basic feasible solution**;  $\{x_3, x_4, x_5, x_6\}$  is the set of **basic variables**.

**Defn 1.44.** *The set of all the feasible solutions is called the **feasible region**.*

There are three possible cases for the feasible region:

- Empty set – no solution;
- Unbounded convex set (an infinite collection of solutions); or,
- A **bounded** convex polytope (convex bounded set, or convex hull) – bounded collection of solutions.

**Defn 1.45.** *Any point  $z$  in a convex set  $\mathcal{S}$  which cannot be expressed as  $z = \lambda x + (1 - \lambda)y$  for  $x, y \in S$  and  $\lambda \in (0, 1)$  is called an **extreme point**.*

*The **vertices** of a convex set are called **extreme points**.*

**Theorem 1.11** (Optimum solution of an LP). *An optimum of a linear objective function  $z = \mathbf{c}^\top \mathbf{x}$ , with convex polytope feasible region, occurs at an extreme point.*

**Theorem 1.12.** *A point in the feasible region is an extreme point if and only if it is a basic feasible solution.*

The **Simplex Method** is a procedure to solve linear programs in standard form. Essentially, it transfers from one basis (corresponding to a basic feasible solution) to another basis (also corresponding to a basic feasible solution) such that the objective function value never increases.

#### MATLAB LINPROG.M

In MATLAB, the function `linprog.m` allows one to solve linear programming problems. Please familiarise yourself with this algorithm; it will be assumed that you can solve linear programming problems using `linprog.m`.

To begin with, Matlab has a help command which enables you to find out about all the inbuilt commands. For example the command `linprog.m` can be investigated at the prompt as follows.

```
>> help linprog
```

```
linprog Linear programming.
```

`X = linprog(f,A,b)` attempts to solve the linear programming problem:

$$\begin{array}{ll} \min & \mathbf{f}'\mathbf{x} \\ \text{subject to:} & \mathbf{A}\mathbf{x} \leq \mathbf{b} \\ & \mathbf{x} \end{array}$$

`X = linprog(f,A,b,Aeq,beq)` solves the problem above while additionally satisfying the equality constraints  $Aeq \cdot x = beq$ .

`X = linprog(f,A,b,Aeq,beq,LB,UB)` defines a set of lower and upper bounds on the design variables,  $X$ , so that the solution is in the range  $LB \leq X \leq UB$ . Use empty matrices for  $LB$  and  $UB$  if no bounds exist. Set  $LB(i) = -\text{Inf}$  if  $X(i)$  is unbounded below; set  $UB(i) = \text{Inf}$  if  $X(i)$  is unbounded above.

`X = linprog(f,A,b,Aeq,beq,LB,UB,X0)` sets the starting point to  $X0$ . This option is only available with the active-set algorithm. The default interior point algorithm will ignore any non-empty starting point.

**Example 17.** For Example 15, using the original inequality form of the LP but rewriting it as a minimisation problem as

$$\min \quad -z = -13x_1 - 12x_2 - 17x_3 = \mathbf{c}^T \mathbf{x} = (-13, -12, -17) \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix},$$

in Matlab we do the following

```
>> f=[-13,-12,-17];
>> A=[2 1 2 ; 1 1 1; 3 3 4]
>> b=[225 ; 117 ; 420]
>> LB=[0;0;0]
>> [x,fval]=linprog(f,A,b,[],[],LB)
Optimization terminated.
x =
    30.0000
     0.0000
    82.5000
```

```
fval =
-1.7925e+03
```

This means we obtain the optimal value  $z^* = 1792.5$  with

$$\mathbf{x} = (x_1, x_2, x_3)^\top = (30, 0, 82.5)^\top.$$

Also for Example 15, we could use the standard equality form of the LP, which is rewritten as a minimisation problem. In that case, in MATLAB we do the following

```
>> f=[-13,-12,-17,0,0,0];
>> Aeq=[2 1 2 1 0 0 ; 1 1 1 0 1 0 ; 3 3 4 0 0 1]
>> beq=[225 ; 117 ; 420]
>> LB=[0;0;0;0;0;0]
>> [x,fval]=linprog(f,[],[],Aeq,beq,LB)
Optimization terminated.
```

```
x =
30.0000
0.0000
82.5000
0.0000
4.5000
0.0000
```

```
fval =
-1.7925e+03
```

Once again, we obtain  $z^* = 1792.5$ , but with the extended version

$$\mathbf{x} = (30, 0, 82.5, 0, 4.5, 0)^\top.$$

**Example 18.** For Example 16, using the original inequality form of the LP but re-writing as a minimisation problem as

$$\min \quad -z = -3x_1 - 2x_2 = \mathbf{c}^\top \mathbf{x} = (-3, -2) \begin{pmatrix} x_1 \\ x_2 \end{pmatrix},$$

in MATLAB we do the following

```
>> f=[-3;-2];
>> A=[1 0;0 1;2 3;1 2]
>> b=[10;10;50;40]
>> LB=[0;0]
>> [x,fval]=linprog(f,A,b,[],[],LB)
Optimization terminated.
```

```
x =
    10
    10
```

```
fval =
   -50
```

So we get the optimal value  $z^* = 50$  with

$$\mathbf{x} = (x_1, x_2)^\top = (10, 10)^\top.$$

Also for Example 16, we could use the standard equality form of the LP, which is also rewritten as a minimisation problem, so that in MATLAB we do the following

```
>> f=[-3;-2;0;0;0;0];
>> Aeq=[1 0 1 0 0 0;0 1 0 1 0 0;2 3 0 0 1 0;1 2 0 0 0 1]
>> beq=[10;10;50;40]
>> LB=[0;0;0;0;0;0]
>> [x,fval]=linprog(f,[],[],Aeq,beq,LB)
```

Optimization terminated.

```
x =
  10.0000
  10.0000
   0.0000
   0.0000
   0.0000
   0.0000
   5.0000
```

```
fval =
 -50.0000
```

We again get  $z^* = 50$  but with  $\mathbf{x} = (10, 10, 0, 0, 0, 5)^\top$ .

### DUALITY THEORY

Finally, we revise the **dual** of a linear program and **duality** theory. Given the linear program

$$\begin{aligned} (P) \text{ minimise } & z = \mathbf{c}^\top \mathbf{x} \\ \text{such that } & A\mathbf{x} = \mathbf{b}, \\ & \mathbf{x} \geq \mathbf{0} \end{aligned}$$

then the *dual* is

$$\begin{aligned} (D) \text{ maximise } & \mathbf{w} = \mathbf{y}^\top \mathbf{b} \\ \text{subject to } & A^\top \mathbf{y} \leq \mathbf{c} \\ & \mathbf{y} \text{ free.} \end{aligned}$$

Why is this of interest? Recall the following theorem, which shows that a feasible solution of (D), if exists, would provide a **lower bound** for a feasible solution of (P).

**Theorem 1.13.** *For a feasible solution of (P) with value  $z$ , and a feasible solution of (D) with value  $w$ , we have*

$$z \geq w \quad (\text{weak duality}).$$

**Theorem 1.14.** *If (P) has an optimal solution then (D) has an optimal solution and*

$$\min z = \max w \quad (\text{strong duality}).$$

Hence, if (P) has no optimal solution ( $z \rightarrow -\infty$ ) then (D) cannot have any feasible solution, and if (D) has no optimal solution ( $w \rightarrow \infty$ ) then (P) cannot have any feasible solution.

**Additional reading:** The book *Operations Research: An Introduction*, by Taha [7] (available in UoA library, Chapters 2–4), and Matt Roughan’s online lecture notes for *Optimisation and Operations Research II* (available online).



**Generalised Primal/Dual pairs**

	<b>Primal (Dual)</b> $\max z = \sum_{j=1}^n c_j x_j + z_0$	<b>Dual (Primal)</b> $\min w = \sum_{i=1}^m y_i b_i + z_0$
1.	$\sum_{j=1}^n a_{ij} x_j = b_i$	$y_i$ free
2.	$\sum_{j=1}^n a_{ij} x_j \leq b_i$	$y_i \geq 0$
3.	$\sum_{j=1}^n a_{ij} x_j \geq b_i$	$y_i \leq 0$
4.	$x_j \geq 0$	$\sum_{i=1}^m y_i a_{ij} \geq c_j$
5.	$x_j \leq 0$	$\sum_{i=1}^m y_i a_{ij} \leq c_j$
6.	$x_j$ free	$\sum_{i=1}^m y_i a_{ij} = c_j$

## Chapter 2

# Stochastic Linear Programming

This section follows **very, very closely** the reference [3]. Another very good reference is [2].

### §2.1 MOTIVATION

As we have seen in Optimisation and Operations Research II, many real-world decision problems can be modelled as *linear programs*. That is, as a mathematical program of the form:

$$\begin{aligned} &\text{minimise} && z = \mathbf{c}^\top \mathbf{x} \\ &\text{such that} && A\mathbf{x} = \mathbf{b}, \\ &&& \mathbf{x} \geq \mathbf{0}. \end{aligned} \tag{2.1}$$

In (2.1), the *cost* coefficients ( $c_j$ ), the *constraint* coefficients ( $a_{ij}$ ), and the *demand/supply* coefficients ( $b_i$ ) are all assumed to have fixed known values. Furthermore, the model (2.1) is only appropriate when the objective function and constraints are (at least close to) linear in the decision variables ( $x_i$ ).

More generally, we have mathematical programs of the form:

$$\begin{aligned} &\text{minimise} && z = g_0(\mathbf{x}) \\ &\text{such that} && g_i(\mathbf{x}) \leq 0, i = 1, \dots, m \\ &&& \mathbf{x} \in \mathcal{X} \subset \mathbb{R}^n \end{aligned} \tag{2.2}$$

The set  $\mathcal{X}$  and functions  $g_i : \mathbb{R}^n \rightarrow \mathbb{R}, i = 0, \dots, m$ , are defined by the modelling process (i.e., they are problem specific). Depending upon the properties of the set  $\mathcal{X}$  and functions ( $g_i$ ), the problem (2.2) has different names:

1. **linear**, if the set  $\mathcal{X}$  is a convex polytope (bounded or unbounded) and the functions  $g_i$ 's are linear;
2. **nonlinear**, if at least one of the functions  $g_i$ 's is nonlinear, or  $\mathcal{X}$  is not a convex polytope. Nonlinear programs can be further classified as:
  - a) **convex**, if  $\mathcal{X} \cap \{x : g_i(x) \leq 0, i = 1, \dots, m\}$  is a convex set and  $g_0$  is a convex function — for example, if  $g_i$ 's are all convex functions and  $\mathcal{X}$  is a convex set; or,
  - b) **nonconvex**, if either  $\mathcal{X} \cap \{x : g_i(x) \leq 0, i = 1, \dots, m\}$  is not a convex set or the objective function  $g_0$  is not convex.

In many modelling scenarios, it is unreasonable to expect the entries of the vectors  $\mathbf{b}, \mathbf{c}$ , and the matrix  $\mathbf{A}$  of the LP (2.1) to be deterministically fixed. Similarly, we would not always expect the functions  $g_i$ 's or the set  $\mathcal{X}$  in the mathematical program (2.2) to be deterministically fixed.

For instance, the future productivities in a production problem, inflows into a reservoir connected to a hydro power station, and demands at various nodes in a transportation network, are often appropriately modelled as uncertainty parameters, which are best characterised by discrete or continuous probability distributions.

The uncertainty about the realised values of those parameters cannot always be ignored by just inserting their mean values or some other (fixed) estimates during the modelling process. That is, depending on the practical

situation under consideration, problems like those specified by (2.1) may not be the appropriate models for describing what we want to solve. = We will demonstrate this using a basic idealised yet illustrative example of a refinery, which We will manipulate in various ways to show a variety of stochastic linear/non-linear models that you may encounter in the real world.

**The goal** is to essentially find a solution, which is feasible and in some sense, optimal. Such models are generally formulated, solved analytically or numerically, and analysed in order to provide useful information to a decision-maker.

**Example 19 (Refinery).** A refinery produces STANDARD and PREMIUM from two raw crude oils, raw<sub>1</sub> and raw<sub>2</sub>, which cost \$200 and \$300 per barrel respectively, giving the production cost for the products. The demand in gallons for (STANDARD, PREMIUM) per day is given by  $\mathbf{h}^\top = (1800, 1620)$  and the refinery can process a maximum of 100 barrels of raw materials each day. The productivity (or output) of STANDARD is 20 gallons per barrel of raw<sub>1</sub> and 60 gallons per barrel of raw<sub>2</sub>, whereas the productivity (output) of PREMIUM is 30 gallons per barrel of raw<sub>1</sub> and 30 gallons per barrel of raw<sub>2</sub>.

We can write this down as a simple LP and solve for the minimum cost to the refinery to satisfy demand, by letting  $x_i$  be the number of barrels of raw <sub>$i$</sub> , for  $i = 1, 2$  processed by the refinery per day.

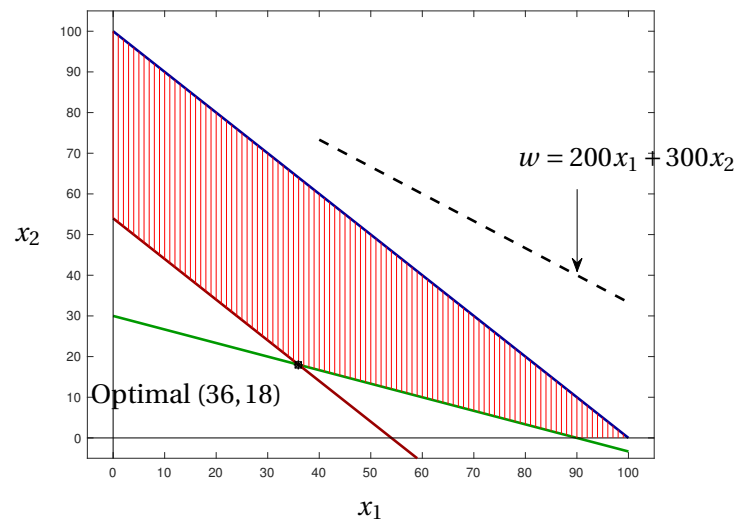
$$\left. \begin{array}{ll} \min & w = 200x_1 + 300x_2 \\ \text{(supply) s.t.} & x_1 + x_2 \leq 100 \\ \text{(standard demand)} & 20x_1 + 60x_2 \geq 1800 \\ \text{(premium demand)} & 30x_1 + 30x_2 \geq 1620 \\ & x_1 \geq 0, x_2 \geq 0 \end{array} \right\} \quad (2.3)$$

We could then solve this LP graphically or use the Simplex algorithm.

**Graphical Solution:**

As can be seen from Figure 2.1, the optimal solution is

$$\mathbf{x}^* = (x_1^*, x_2^*) = (36, 18),$$



**Figure 2.1:** Feasible solution space for Refinery LP, represented by the shaded red area. The thick blue line corresponds to  $x_1 + x_2 = 100$ , the thick green line to  $20x_1 + 60x_2 = 1800$ , and the thick red line to  $30x_1 + 30x_2 = 1620$ . The slope of the objective function  $w = 200x_1 + 300x_2$  is  $-2/3$ .

with the optimal value  $w^* = \$12,600$ .

The **MATLAB solution** agrees with the graphical solution.

```
>> f=[200,300];
>> A=[1,1;-20,-60;-30,-30];
>> b=[100,-1800,-1620];
>> LB=[0;0];
>> [x,fval]=linprog(f,A,b,[],[],LB)
Optimization terminated.
```

```
x =
    36.0000
    18.0000
```

```
fval =
    1.2600e+04
```

Our production problem in Example 19 is properly described by a linear program, and its correct solution is as given above provided the productivities, the unit costs, the demands and the capacity are **fixed data** and **known to us prior to making our decision** on the production plan. However, these are obviously not always realistic assumptions. It may happen that at least some of the data, such as the productivities and demands vary within certain limits (for our discussion, randomly) and that we have to make our decision on the production plan before knowing the exact values of those data.

To be more specific, let us assume in the next example some specific variations that follow a known distributional form.

**Example 20 (Refinery with Random Inputs).** Consider now the following, where some of the data varies randomly and we have to make our decision before we know the exact values:

- demand for STANDARD and PREMIUM varies randomly,
- productivities (outputs) of the fuels from raw products also vary randomly, and
- the outputs of fuel can only be observed during production.

At the same time, as with the case when we had all fixed and known data, the conditions are:

- daily production must be fixed in advance and cannot be changed,
- customers expect their actual demands to be met on a daily basis.

First, we assume that the demand in gallons for (STANDARD, PREMIUM) is now given by

$$\mathbf{h}^\top = (h_s, h_p) = (1800 + Z_1, 1620 + Z_2),$$

where  $Z_1 \sim \mathcal{N}(0, 120)$  and  $Z_2 \sim \mathcal{N}(0, 90)$ . ‘ $\mathcal{N}(\mu, \sigma)$ ’ indicates a Normal distribution with mean  $\mu$  and variance  $\sigma^2$ . In this case, the average demands remain equivalent to the deterministic demands of the original linear program; that is,

$$\mathbf{E}(1800 + Z_1) = 1800 \quad \text{and} \quad \mathbf{E}(1620 + Z_2) = 1620.$$

Let us also assume that the output of STANDARD from raw<sub>1</sub>,  $\alpha$ , and the output of PREMIUM from raw<sub>2</sub>,  $\beta$ , are given by

$$\alpha = 20 + Y_1, \quad \beta = 34 - Y_2,$$

where  $Y_1 \sim U[-8, 8]$  and  $Y_2 \sim \exp(0.25)$ . ‘ $U[a, b]$ ’ indicates uniform distribution on  $[a, b]$ , and ‘ $\exp(\lambda)$ ’ indicates exponential distribution with parameter  $\lambda$ . With these distributions, the average productivities remain equivalent to the deterministic productivities of the original linear program; that is,

$$\mathbf{E}(20 + Y_1) = 20 \quad \text{and} \quad \mathbf{E}(34 - Y_2) = 30.$$

Note that all random variables are assumed to be mutually independent. If we define  $\zeta := (Z_1, Z_2, Y_1, Y_2)^\top$  and let

$$\alpha(\zeta) = 20 + Y_1, \quad \beta(\zeta) = 34 - Y_2, \quad h_s(\zeta) = 1800 + Z_1, \quad h_p(\zeta) = 1620 + Z_2,$$

we can write the ‘stochastic’ linear program

$$\left. \begin{array}{ll} \text{‘min’} & w = 200x_1 + 300x_2 \\ \text{s.t.} & \begin{array}{l} x_1 + x_2 \leq 100 \\ \alpha(\zeta)x_1 + 60x_2 \geq h_s(\zeta) \\ 30x_1 + \beta(\zeta)x_2 \geq h_p(\zeta) \\ x_1 \geq 0, \quad x_2 \geq 0. \end{array} \end{array} \right\} \quad (2.4)$$

The ‘stochastic’ linear program given in (2.4) is not well defined, as the meaning of ‘min’ is not clear until we know the realisation  $\epsilon = (z_1, z_2, y_1, y_2)^\top$  of the random vector  $\zeta = (Z_1, Z_2, Y_1, Y_2)^\top$ . The random variables  $Z_1, Z_2$  and  $Y_2$

are also unbounded and so we restrict them to their 99% confidence intervals, to get the following ‘allowable’ realisations:

$$Z_1 \in [-309.1, 309.1], Z_2 \in [-231.8, 231.8], Y_1 \in [-8, 8], Y_2 \in [0.0, 18.42].$$

**How should we deal with the randomness of  $Z_1$ ,  $Z_2$ ,  $Y_1$ , and  $Y_2$ ?**

**Option 1.** Suppose we fix  $Y_1 = 0$  and  $Y_2 = 4$  to their mean values, and use the 99% confidence interval limits of  $Z_1$  and  $Z_2$ . The extreme values for each random variable yield

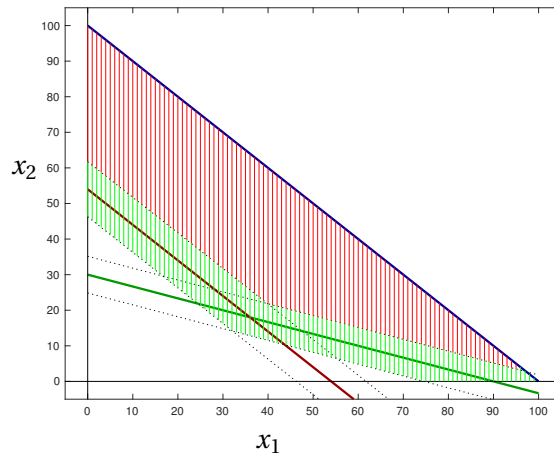
$$\begin{aligned} \text{‘min’ } w &= 200x_1 + 300x_2 \\ \text{s.t.} \quad & x_1 + x_2 \leq 100 \\ & 20x_1 + 60x_2 \geq 1800 - 309.1 \\ & 20x_1 + 60x_2 \geq 1800 + 309.1 \\ & 30x_1 + 30x_2 \geq 1620 - 231.8 \\ & 30x_1 + 30x_2 \geq 1620 + 231.8 \\ & x_1 \geq 0, \quad x_2 \geq 0, \end{aligned}$$

or, equivalently,

$$\begin{aligned} \text{‘min’ } w &= 200x_1 + 300x_2 \\ \text{s.t.} \quad & x_1 + x_2 \leq 100 \\ & 20x_1 + 60x_2 \geq 1490.9 \\ & 20x_1 + 60x_2 \geq 2109.1 \\ & 30x_1 + 30x_2 \geq 1388.2 \\ & 30x_1 + 30x_2 \geq 1851.8 \\ & x_1 \geq 0, \quad x_2 \geq 0. \end{aligned}$$

This results in a parallel shift in the averaged value facets of the original feasible region, shown in solid (red and green) lines, to the dotted lines in Figure 2.2. Note that the shaded red area is feasible for all values of  $Z_1$  and  $Z_2$  (in their 99% confidence intervals).



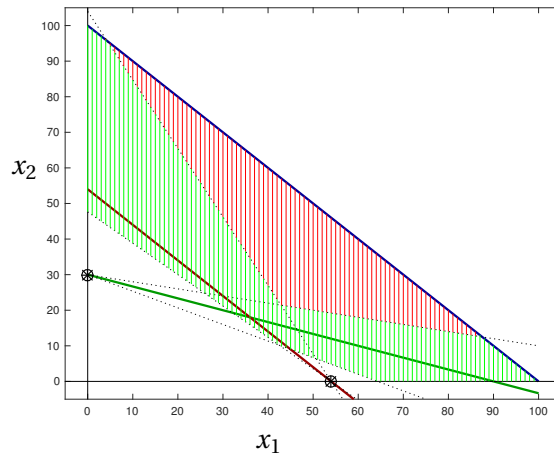


**Figure 2.2:** Feasible solution space varying with demands

**Option 2:** Alternatively, fix  $Z_1 = 0$  and  $Z_2 = 0$  to their mean values, and use the full distribution of  $Y_1$  and the 99% confidence interval limits of  $Y_2$ . The extreme values for each random variable yield

$$\begin{aligned}
 \text{'min' } w &= 200x_1 + 300x_2 \\
 \text{s.t.} \quad &x_1 + x_2 \leq 100 \\
 &12x_1 + 60x_2 \geq 1800 \\
 &28x_1 + 60x_2 \geq 1800 \\
 &30x_1 + 34x_2 \geq 1620 \\
 &30x_1 + 15.58x_2 \geq 1620 \\
 &x_1 \geq 0, \quad x_2 \geq 0.
 \end{aligned}$$

This results in rotations (about the pivot points  $\otimes$ ) of the facets of the original feasible region, shown in solid (red and green) lines, to the dotted lines in Figure 2.3. Note that the shaded red area is feasible for all values of  $Y_1$  and  $Y_2$ .

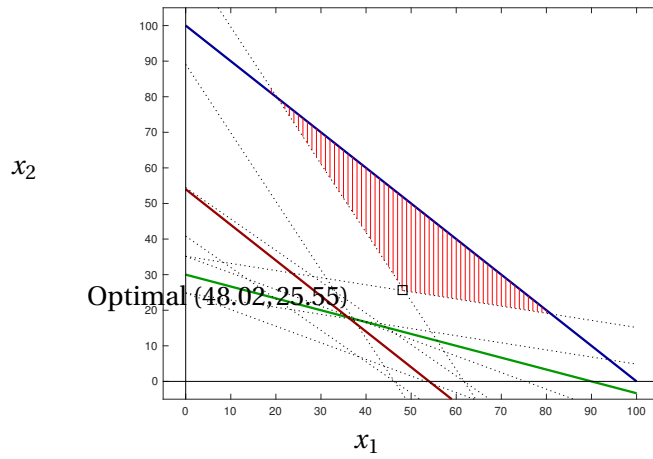


**Figure 2.3:** Feasible solution space varying with productivities

**Option 3:** Consider all possible changes in demands and productivities together, yielding

$$\begin{aligned}
 \text{'min' } w &= 200x_1 + 300x_2 \\
 \text{s.t. } & x_1 + x_2 \leq 100 \\
 & 12x_1 + 60x_2 \geq 1490.9 \\
 & 12x_1 + 60x_2 \geq 2109.1 \\
 & 28x_1 + 60x_2 \geq 1490.9 \\
 & 28x_1 + 60x_2 \geq 2109.1 \\
 & 30x_1 + 34x_2 \geq 1388.2 \\
 & 30x_1 + 34x_2 \geq 1851.8 \\
 & 30x_1 + 15.58x_2 \geq 1388.2 \\
 & 30x_1 + 15.58x_2 \geq 1851.8 \\
 & x_1 \geq 0, \quad x_2 \geq 0.
 \end{aligned}$$

This results in both parallel translations and rotational movements, shown as dotted lines in Figure 2.4. The facets of the average value feasible region are shown in solid black.



**Figure 2.4:** Feasible solution space varying with productivities and demands

Note that the feasible solution space depicted in Figure 2.4 is much contracted, being risk averse, and yields an **expensive optimal** solution of

$$\mathbf{x}^* = (x_1^*, x_2^*) = (48.02, 25.55)$$

with the corresponding optimal value  $w^* = \$17,269.00$ . This optimal solution is a ‘fat’ solution, which provides a ‘safe’ production program where the solution remains feasible under all possible realisations of the random variables.

**The variation in the feasible region is substantial even for such a small problem, so what can we do?**

## §2.2 GENERAL FORMULATION

**Defn 2.1.** In a *stochastic program*, we have that some of the coefficients are random variables. We might think of the problem as

$$\begin{aligned} \text{'min' } & z = g_0(\mathbf{x}, \boldsymbol{\zeta}) \\ \text{such that } & g_i(\mathbf{x}, \boldsymbol{\zeta}) \leq 0, i = 1, \dots, m \\ & \mathbf{x} \in \mathcal{X} \subset \mathbb{R}^n \end{aligned} \tag{2.5}$$

where  $\boldsymbol{\zeta}$  is a random vector. Note that  $g_i(\mathbf{x}, \boldsymbol{\zeta})$  are random variables for all  $\mathbf{x}$  and  $i = 0, \dots, m$ .

**Assumptions:** The distribution function of  $\boldsymbol{\zeta}$ ,  $F_{\boldsymbol{\zeta}}(\cdot)$ , is known, and  $\boldsymbol{\zeta}$  is independent of the decision variables,  $\mathbf{x}$ .

The ‘mathematical program’ (2.5) is not well defined. The objective function  $g_0(\mathbf{x}, \boldsymbol{\zeta})$  is a random variable, so what is the meaning of ‘minimisation’? Given a realisation,  $\boldsymbol{\varepsilon}$  of  $\boldsymbol{\zeta}$ , we could then choose decision variables  $\mathbf{x}$  to minimise  $g_0(\mathbf{x}, \boldsymbol{\varepsilon})$  — a standard mathematical programming problem — but what about the situation where we wish, or need, to make a decision now, before the realisation  $\boldsymbol{\zeta}$ ? Similarly, what is the interpretation of the constraints  $g_i(\mathbf{x}, \boldsymbol{\zeta}) \leq 0$ ?

**How might you handle such a problem?**

One approach is to translate this problem into a *deterministic* problem, so the mathematical program is well defined. This leads to so-called **deterministic equivalent problems** (DEPs).

## §2.3 NAÏVE DEP

**Defn 2.2.** One DEP, which we will call the **naïve DEP**, is

$$\begin{aligned} \min \quad & z = g_0(\mathbf{x}, \boldsymbol{\zeta}) \\ \text{subject to} \quad & g_i(\mathbf{x}, \boldsymbol{\epsilon}) \leq 0, \quad i = 1, \dots, m, \text{ and for all realisations } \boldsymbol{\epsilon} \text{ of } \boldsymbol{\zeta}, \\ & \mathbf{x} \in \mathcal{X} \subset \mathbb{R}^n, \end{aligned} \quad (2.6)$$

where the objective function  $z$  is evaluated for the resultant  $\mathbf{x}$  when every constraint under  $\boldsymbol{\zeta}$  is satisfied.

**Do you think this is sensible? Why, or why not?**

## §2.4 RECOURSE DEP

## RECOURSE DEP: EXAMPLES

**Example 21 (Refinery with Recourse).** Consider again our Refinery problem, where we now assume that the refinery has made the following arrangement with their customers.

The customers expect the refinery to satisfy their daily demands. However, due to the random nature of demands and productivity, the demands may not be covered. This will cause a ‘penalty’ cost on the refinery, as the shortfall must be purchased from the market. Assume that the penalty is proportional to the shortfall and that the price per gallon of the undeliverable fuels are as follows

$$\mathbf{q}^\top = (q_s, q_p) = (\$70, \$120).$$

These costs are due to shortage or in general due to a violation of the constraints and must be determined after the observation of the random data. They are referred to as **recourse costs**.

In the case of repeated execution of the production program (daily in our Refinery problem), it is meaningful to apply an **expected value criterion**.

That is, we may want to find a production plan that minimises the original **first stage costs** and the **expected recourse costs**.

For each of the two stochastic constraints in the model of equations (2.4), we introduce a recourse variable  $y_i(\zeta)$  for  $i = 1, 2$ , which measures the shortfall in production (it is like a slack variable), where recall that  $\zeta := (Z_1, Z_2, Y_1, Y_2)^\top$ . Note that because the production shortfall depends on  $\zeta$ , so do the corresponding recourse variables  $y_i(\zeta)$ . This makes the recourse variables themselves **random**.

We can now replace our vague stochastic program with a well defined stochastic program with recourse. That is, we include the recourse variables  $y_1(\zeta)$  and  $y_2(\zeta)$  in the inequalities and add the expected recourse cost w.r.t.  $\zeta$ , given by

$$\mathbf{E}_\zeta[70y_1(\zeta) + 120y_2(\zeta)]$$

to the objective function to get the **two-stage linear recourse problem** given by:

$$\left. \begin{array}{ll} \text{'min'} & w = 200x_1 + 300x_2 + \mathbf{E}_\zeta(70y_1(\zeta) + 120y_2(\zeta)) \\ \text{s.t.} & \begin{array}{rcl} x_1 & + & x_2 & \leq & 100 \\ \alpha(\zeta)x_1 & + & 60x_2 & + & y_1(\zeta) & \geq & h_s(\zeta) \\ 30x_1 & + & \beta(\zeta)x_2 & + & y_2(\zeta) & \geq & h_p(\zeta) \\ x_1 \geq 0, & x_2 \geq 0, & y_1(\zeta) \geq 0, & y_2(\zeta) \geq 0 \end{array} \end{array} \right\} \quad (2.7)$$

The problem given in (2.7) is an example of a deterministic equivalent problem known as a **Recourse DEP**. In general the DEP given in (2.7) must hold almost surely (a.s., or with probability 1).

**Example 22 (Refinery with Recourse: Finite discrete variations).** If the random vector  $\zeta$  has a finite discrete distribution

$$\{(\epsilon_k, p_k), k = 1, \dots, r\},$$

where  $r < \infty$  and  $p_k > 0$  for all  $k$ , we get a linear program written in **expanded form** as

$$\left. \begin{array}{ll} \min w &= 200x_1 + 300x_2 + \sum_{k=1}^r p_k [70y_1(\epsilon_k) + 120y_2(\epsilon_k)] \\ \text{s.t.} & \begin{array}{l} x_1 + x_2 \leq 100 \\ \alpha(\epsilon_k)x_1 + 60x_2 + y_1(\epsilon_k) \geq h_s(\epsilon_k), \forall k \\ 30x_1 + \beta(\epsilon_k)x_2 + y_2(\epsilon_k) \geq h_p(\epsilon_k), \forall k \\ x_1 \geq 0, x_2 \geq 0, y_1(\epsilon_k) \geq 0, y_2(\epsilon_k) \geq 0, \forall k. \end{array} \end{array} \right\} \quad (2.8)$$

Because  $\zeta$  has a finite discrete distribution  $\{(\epsilon_k, p_k), k = 1, \dots, r\}$ , where  $p_k > 0$  for all  $k$ , we noted that the LP (2.8) is just an ordinary linear program, but it has what is called a **Dual Decomposition Structure**. Depending on the number  $r$  of realisations  $\epsilon$  of  $\zeta$ , this linear program may become very large in scale (thus making it difficult or even impossible to solve in this form), but its particular block structure is amenable to specially designed algorithms. We will return to this later in the course.

Let's continue our example with continuous distributions  $\zeta$ , where we will form discrete approximations.

**Example 23** (Refinery with Recourse: Random Demands and Discretisation).

Initially, let's consider the situation where the demands  $\mathbf{d}$  vary randomly as the normal distributions  $Z_1$  and  $Z_2$  and the productivities are fixed, similar to the situation as displayed in Figure 2.2. Even in this situation the non-linear LP (NLP) can present numerical difficulties as evaluation of the expectation in the objective function requires

- multivariate numerical integration,
- implicit definition of the functions  $y_1(\epsilon)$  and  $y_2(\epsilon)$  for every possible realisation  $\epsilon$  of  $\zeta$ .

To deal with these difficulties, We will approximate the normal distributions by **finite discrete** distributions, using a technique that can be applied to most continuous distributional forms.

In order to approximate the normal distributions we can

- generate samples  $z_1^\mu$  and  $z_2^\mu$ , from the random variables  $Z_1$  and  $Z_2$  respectively for  $\mu = 1, 2, \dots, K$ , restricted to the 99% confidence intervals (CIs) for some large value of  $K$ ;
- create equi-spaced partitions of the 99% CIs into a finite number — say,  $r_i$  — of subintervals  $I_{i,v}$  for  $i = 1, 2$ ;
- calculate for each subinterval the mean  $\overline{z_{i,v}}$  of sample values  $z_i^\mu \in I_{i,v}$  for  $i = 1, 2$  and  $v = 1, 2, \dots, r_i$ , to use as an estimate of

$$\mathbf{E}(Z_i | Z_i \in I_{i,v});$$

- calculate for every subinterval  $I_{i,v}$  the relative frequency  $p_{i,v}$  for  $z_i^\mu \in I_{i,v}$ , which yields an estimate for

$$P(Z_i \in I_{i,v}).$$

The following figure displays a result of using  $K = 100,000$  and  $r_1 = r_2 = 15$  in this approximation method for the two normal distributions  $\mathcal{N}(1800, 120)$  and  $\mathcal{N}(1620, 90)$ .

**NOTE:** We must always remain aware that using discrete approximations inevitably leads to some discretisation error.

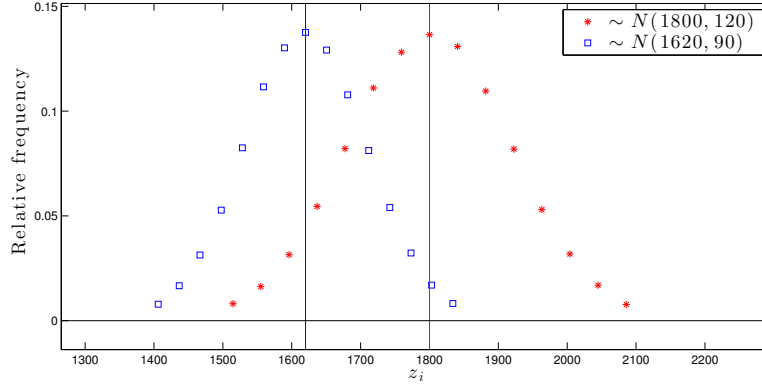
With  $K = 100,000$ ,  $r_1 = r_2 = 15$ , we get  $15^2 = 225$  realisations for the joint distribution of  $\zeta$ . Using these realisations in the expanded form LP (2.8), we obtain

$$\begin{aligned} \tilde{w} &= 14084.3728 \\ \tilde{\mathbf{x}} = (\tilde{x}_1, \tilde{x}_2) &= (38.5490, 20.5478) \\ &\text{with first stage costs of} \\ 200\tilde{x}_1 + 300\tilde{x}_2 &= 13874.1504. \end{aligned}$$

The **empirical reliability** is defined to be the probability that  $\tilde{\mathbf{x}}$  is feasible such that the recourse variables are zero:

$$P(\alpha(\zeta)\tilde{x}_1 + 60\tilde{x}_2 \geq h_s(\zeta), 30\tilde{x}_1 + \beta(\zeta)\tilde{x}_2 \geq h_p(\zeta)),$$





**Figure 2.5:** Discrete approximation of demand distributions.

where we use the discrete approximation of  $\zeta$ .

In this case, **empirical reliability** is

$$\rho(\tilde{\mathbf{x}}) = 0.9516.$$

If we force our original LP solution  $\mathbf{x}^* = (36, 18)$  (that is, if we substitute the original LP solution into the expanded form LP), we see that

$$w^* = 19980.0372$$

with first stage costs of

$$\mathbf{x}^*(200, 300)' = 12600.$$

The empirical reliability of this solution is  $\rho(\mathbf{x}^*) = 0.3206$ , which in this instance, far exceeds the theoretical expected reliability of 0.25 (half of the time each of the normal distributions are  $\geq 0$ ).

This indicates to some extent the discretisation errors that arise due to the approximation of the continuous distributions.

In the next example We will consider the same technique used when the productivities are random and the demands are fixed.

**Example 24** (Refinery with Recourse: Random Productivities and Discretisation). In the second instance, the demands are fixed and the productivities vary randomly as per the distributions of  $Y_1$  and  $Y_2$ , similar to the situation as displayed in Figure 2.3. Similarly to before, we

- generate samples  $y_1^\mu$  and  $y_2^\mu$ , from the random variables  $Y_1$  and  $Y_2$  respectively for  $\mu = 1, 2, \dots, K$ , for some large value of  $K$ , where  $Y_2$  is restricted to its 99% CI,
- create equi-spaced partitions of the data range intervals into a finite number—say,  $r_i$ —of subintervals  $I_{i,v}$  for  $i = 1, 2$ ,
- calculate for each subinterval the mean  $\overline{y_{i,v}^\mu}$  of sample values  $y_i^\mu \in I_{i,v}$  for  $i = 1, 2$  and  $v = 1, 2, \dots, r_i$ , to use as an estimate of

$$\mathbf{E}(Y_i | Y_i \in I_{i,v}),$$

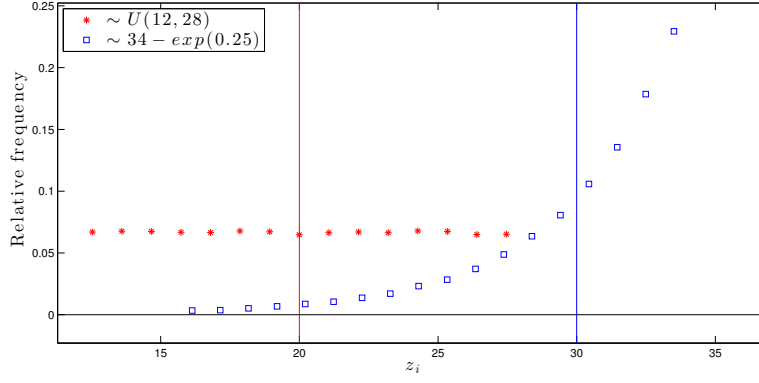
- calculate for every subinterval  $I_{i,v}$  the relative frequency  $p_{i,v}$  for  $y_i^\mu \in I_{i,v}$ , which yields an estimate for

$$P(Y_i \in I_{i,v}).$$

The following figure displays a result of using  $K = 100,000$ ,  $r_1 = 15$  and  $r_2 = 18$  in this approximation method for the two distributions  $U(12, 28)$  and  $34 - \exp(0.25)$ .

Using these discrete distributions as approximations to the given uniform and exponential distributions, with  $K = 100,000$ ,  $r_1 = 15$  and  $r_2 = 18$  we get  $15 \times 18 = 270$  realisations for the joint distribution of  $\zeta$ . Using this in the expanded form LP (2.8), yields

$$\begin{aligned} \tilde{w} &= 14451.7938 \\ \tilde{\mathbf{x}} = (x_1, x_2) &= (37.5610, 22.1520) \\ &\quad \text{with first stage costs of} \\ 200\tilde{x}_1 + 300\tilde{x}_2 &= 14157.7893. \end{aligned}$$



**Figure 2.6:** Discrete approximation of productivities.

The empirical reliability  $\rho(\tilde{\mathbf{x}}) = 0.9617$ .

If we again force our original LP solution  $\mathbf{x}^* = (36, 18)$  we see that

$$w^* = 20477.7807$$

with first stage costs of

$$\mathbf{x}^*(200, 300)' = 12600.$$

The empirical reliability of solution is  $\rho(\mathbf{x}^*) = 0.3020$ .

**Example 25** (Refinery with Recourse: Random Inputs and their Discrete Approximations). Approximating all of the four random demand and productivity distributions, with  $K = 100,000$ ,  $r_1 = 5$ ,  $r_2 = 9$ ,  $r_3 = 7$  and  $r_4 = 11$ , we get  $5 \times 9 \times 7 \times 11 = 3465$  realisations of the joint distribution of  $\zeta$ .

Using these realisations in the expanded form LP (2.8), yields

$$\tilde{w} = 15058.2884$$

$$\tilde{\mathbf{x}} = (\tilde{x}_1, \tilde{x}_2) = (37.8024, 23.6010)$$

with first stage costs of

$$200x_1 + 300x_2 = 14640.7874.$$

The empirical reliability is  $\rho(\tilde{\mathbf{x}}) = 0.9441$ .

Again, if we force our original LP solution  $\mathbf{x}^* = (36, 18)$  we see that

$$\begin{aligned} w^* &= 23183.6708 \\ &\quad \text{with first stage costs of} \\ \mathbf{x}^*(200, 300)' &= 12600. \end{aligned}$$

The empirical reliability of this solution is  $\rho(\mathbf{x}^*) = 0.2548$ .

It is worth noting that for a relatively small increase in first stage costs, we have retrospectively calculated the empirical reliability for our recourse problems, which appear to have increased from 0.25 up to around 0.95.

Although these calculations may not provide a true indication of the size of the constraint violations, there are many situations where this form of reliability is a prime consideration, such as in medical and technical applications.

## RECOURSE DEP: GENERAL FORMULATION

**Defn 2.3.** For a given stochastic linear program, let us start by defining for  $i = 1, \dots, m$  and for each realisation  $\epsilon$  of  $\zeta$  the following

$$g_i^+(\mathbf{x}, \epsilon) := \begin{cases} 0 & \text{if } g_i(\mathbf{x}, \epsilon) \leq 0, \\ g_i(\mathbf{x}, \epsilon) & \text{otherwise.} \end{cases}$$

This implies the  $i$ th constraint of (2.5) is not satisfied if and only if  $g_i^+(\mathbf{x}, \epsilon) > 0$  for a given decision  $\mathbf{x}$  and realisation  $\epsilon$ . Furthermore, the extent of the violation in the  $i$ th constraint is  $g_i^+(\mathbf{x}, \epsilon)$ .

If such violations occur for a given *a priori* choice of  $\mathbf{x}$  and realisation  $\epsilon$ , it might be possible to undertake a **second-stage activity**, or **recourse**, in order to satisfy the constraint.

That is, we could provide for each constraint violation, a recourse or **second-stage** activity  $y_i(\epsilon)$  that, after observing the realisation  $\epsilon$ , is chosen such as to compensate for each constraint violation, if they exist, by satisfying

$$g_i^+(\mathbf{x}, \epsilon) - y_i(\epsilon) \leq 0,$$

or, equivalently,

$$y_i(\epsilon) \geq g_i^+(\mathbf{x}, \epsilon).$$

If this extra effort is assumed to cause an extra cost or penalty of  $q_i (\geq 0)$  per unit of violation, then our additional costs – that is, the **recourse function** – is

$$Q(\mathbf{x}, \epsilon) = \min_y \left\{ \sum_{i=1}^m q_i y_i(\epsilon) : y_i(\epsilon) \geq g_i^+(\mathbf{x}, \epsilon), i = 1, \dots, m \right\}, \quad (2.9)$$

(which is a linear program itself,) as we wish to satisfy the constraints with minimum expense. Hence, the *total cost* – composed of first stage costs plus the recourse costs – is

$$f_0(\mathbf{x}, \epsilon) = g_0(\mathbf{x}, \epsilon) + Q(\mathbf{x}, \epsilon). \quad (2.10)$$

Note that here we have one recourse variable per constraint.

Rather than (2.9), we might think of a more general linear recourse program with a different recourse vector

$$\mathbf{y} = (y_1, y_2, \dots, y_{\bar{n}}) \in \mathcal{Y} \subset \mathbb{R}^{\bar{n}},$$

where  $\mathcal{Y}$  is some given polytope (such as  $\{\mathbf{y} | \mathbf{y} \geq \mathbf{0}\}$ ). We also have an arbitrary, fixed  $m \times \bar{n}$  matrix  $W$  – the **recourse matrix** – and a corresponding unit cost vector  $\mathbf{q} = (q_1, q_2, \dots, q_{\bar{n}}) \in \mathbb{R}^{\bar{n}}$ , yielding for (2.10) the recourse function:

$$Q(\mathbf{x}, \boldsymbol{\varepsilon}) = \min_{\mathbf{y}(\boldsymbol{\varepsilon})} \{ \mathbf{q}^\top \mathbf{y}(\boldsymbol{\varepsilon}) : W \mathbf{y}(\boldsymbol{\varepsilon}) \geq \mathbf{g}^+(\mathbf{x}, \boldsymbol{\varepsilon}), \mathbf{y}(\boldsymbol{\varepsilon}) \in \mathcal{Y} \}, \quad (2.11)$$

where  $\mathbf{g}^+(\mathbf{x}, \boldsymbol{\varepsilon}) = (g_1^+(\mathbf{x}, \boldsymbol{\varepsilon}), \dots, g_m^+(\mathbf{x}, \boldsymbol{\varepsilon}))^\top$  is the vector of constraint violations (recall that these may be 0).

The recourse vector here no longer directly coincides with one variable per constraint with an associated direct cost, but rather intends to satisfy the constraint violations using some other process.

The difference between the recourse functions (2.9) and (2.11) is in the mechanism by which we use recourse to satisfy any violations:

- The recourse function (2.9) could for instance be interpreted as buying the shortage of products at the market.
- Alternatively, the recourse function in (2.11) gives us the flexibility to model a *second-stage*, or an *emergency* production program, which is effected with *inputs*  $\mathbf{y}$  and a *technology* (a procedure for converting inputs into products) represented by the matrix  $W$ .

Note that the recourse function (2.9) appears as a special case of (2.11); that is, the recourse function (2.9) is readily expressible in the form (2.11), with  $W = I_{m \times m}$ .

Finally, note that the recourse function could also be non-linear:

$$Q(\mathbf{x}, \boldsymbol{\varepsilon}) = \min_{\mathbf{y}} \left\{ q(\mathbf{y}) \mid \begin{array}{l} H_i(\mathbf{y}) \geq g_i^+(\mathbf{x}, \boldsymbol{\varepsilon}), i = 1, \dots, m \\ \mathbf{y} \in \mathcal{Y} \subset \mathbb{R}^{\bar{n}} \end{array} \right\}, \quad (2.12)$$

where  $q(\mathbf{y}) : \mathbb{R}^{\bar{n}} \rightarrow \mathbb{R}$  and  $H_i(\mathbf{y}) : \mathbb{R}^{\bar{n}} \rightarrow \mathbb{R}$  are assumed to be given.

In any case, if it is acceptable and meaningful to the decision maker to minimise the expected value of the total costs, instead of problem (2.5), we consider its deterministic equivalent, the **two-stage stochastic program with recourse**:

$$\min_{\mathbf{x} \in \mathcal{X}} \mathbf{E}_{\zeta}(f_0(\mathbf{x}, \zeta)) = \min_{\mathbf{x} \in \mathcal{X}} \mathbf{E}_{\zeta}(g_0(\mathbf{x}, \zeta) + Q(\mathbf{x}, \zeta)) \quad (2.13)$$

where  $\mathbf{E}_{\zeta}(\cdot)$  means the expectation with respect to  $\zeta$ .

### RECOURSE DEP: MULTI-STAGE RECOURSE PROBLEM

This approach can be extended to a **multi-stage recourse problem**. That is, instead of two decisions  $\mathbf{x}$  and  $\mathbf{y}$ , to be taken at stages 1 and 2, we would have  $K + 1$  sequential decisions  $\mathbf{x}_0, \mathbf{x}_1, \dots, \mathbf{x}_K$ , where  $\mathbf{x}_\tau \in \mathbb{R}^{\bar{n}_\tau}$  for  $\tau = 1, \dots, K$ , to be taken at the subsequent stages  $\tau = 0, 1, \dots, K$ .

Let's assume that the objective function from (2.5) is deterministic, so that  $g_0(\mathbf{x}, \zeta) \equiv g_0(\mathbf{x})$ . At stage  $\tau \geq 1$ , we know the realisations  $\boldsymbol{\varepsilon}_1, \dots, \boldsymbol{\varepsilon}_\tau$  of the random vectors  $\zeta_1, \dots, \zeta_\tau$ , as well as the previous decision vectors  $\mathbf{x}_1, \dots, \mathbf{x}_{\tau-1}$ , and we have to decide on  $\mathbf{x}_\tau$  such that the constraint(s) with vector-valued constraint functions

$$\mathbf{g}_K(\mathbf{x}_0, \mathbf{x}_1, \dots, \mathbf{x}_\tau, \boldsymbol{\varepsilon}_1, \dots, \boldsymbol{\varepsilon}_\tau)$$

are satisfied. At this stage, these can only be achieved by the proper choice of  $\mathbf{x}_\tau$ , based on the knowledge of the previous decisions and realizations.

Hence, given a cost function  $q_\tau(\mathbf{x}_\tau)$  at stage  $\tau \geq 1$ , we have a recourse function

$$\begin{aligned} & Q_\tau(\mathbf{x}_0, \mathbf{x}_1, \dots, \mathbf{x}_{\tau-1}, \boldsymbol{\varepsilon}_1, \dots, \boldsymbol{\varepsilon}_\tau) \\ &= \min_{\mathbf{x}_\tau} \left\{ q_\tau(\mathbf{x}_\tau) \mid \mathbf{g}_\tau(\mathbf{x}_0, \mathbf{x}_1, \dots, \mathbf{x}_\tau, \boldsymbol{\varepsilon}_1, \dots, \boldsymbol{\varepsilon}_\tau) \leq 0 \right\}, \end{aligned} \quad (2.14)$$

where  $\mathbf{g}_\tau(\cdot)$  is an  $m$ -vector of constraints at stage  $\tau$  such that the term “stages” can, but need not, be interpreted as “time periods”.

The recourse function (2.14) indicates that the optimal recourse action  $\hat{\mathbf{x}}_\tau$  at time  $\tau$  depends on the previous decisions and on the realizations observed until stage  $\tau$ . That is,

$$\hat{\mathbf{x}}_\tau = \hat{\mathbf{x}}_\tau(\mathbf{x}_0, \mathbf{x}_1, \dots, \mathbf{x}_{\tau-1}, \boldsymbol{\varepsilon}_1, \dots, \boldsymbol{\varepsilon}_\tau), \quad \tau \geq 1.$$

Hence, taking into account the multiple stages, we get the total costs for the multistage problem as

$$f_0(\mathbf{x}_0, \boldsymbol{\varepsilon}_1, \boldsymbol{\varepsilon}_2, \dots, \boldsymbol{\varepsilon}_K) = g_0(\mathbf{x}_0) + \sum_{\tau=1}^K Q_\tau(\mathbf{x}_0, \hat{\mathbf{x}}_1, \dots, \hat{\mathbf{x}}_{\tau-1}, \boldsymbol{\varepsilon}_1, \dots, \boldsymbol{\varepsilon}_\tau).$$

This yields the deterministic equivalent for the described dynamic decision problem, **the multistage stochastic program with recourse**, as

$$\min_{\mathbf{x}_0 \in \mathcal{X}} \mathbb{E}_\zeta \left( g_0(\mathbf{x}_0) + \sum_{\tau=1}^K Q_\tau(\mathbf{x}_0, \hat{\mathbf{x}}_1, \dots, \hat{\mathbf{x}}_{\tau-1}, \zeta_1, \dots, \zeta_\tau) \right),$$

which is a straight generalisation of the former (two-stage) stochastic program with recourse given in (2.13).



## STOCHASTIC LINEAR PROGRAMS WITH RECOURSE

Let us now focus on *stochastic linear programs with recourse*. Assume that we are given the following stochastic linear program

$$\begin{aligned}
 &\text{minimise} && w = \mathbf{c}^\top \mathbf{x} \\
 &\text{such that} && A\mathbf{x} = \mathbf{b}, \\
 &&& T(\boldsymbol{\zeta})\mathbf{x} = \mathbf{h}(\boldsymbol{\zeta}), \\
 &&& \mathbf{x} \geq \mathbf{0}.
 \end{aligned} \tag{2.15}$$

The  $m_1 \times n$  matrix  $T(\boldsymbol{\zeta})$  and the vector  $\mathbf{h}(\boldsymbol{\zeta})$  can depend on vector  $\boldsymbol{\zeta}$  and hence can have random entries. The  $m_0 \times n$  matrix  $A$  and the vector  $\mathbf{b}$ , on the other hand, are deterministic.

In general, we will assume that the dependence on  $\boldsymbol{\zeta} \in \mathcal{S}_{\boldsymbol{\zeta}} \subset \mathbb{R}^k$  is of the form

$$\begin{aligned}
 T(\boldsymbol{\zeta}) &= \hat{T}_0 + \sum_{i=1}^k 1_{\{\boldsymbol{\varepsilon}_i\}} \hat{T}_i, \\
 \mathbf{h}(\boldsymbol{\zeta}) &= \hat{\mathbf{h}}_0 + \sum_{i=1}^k 1_{\{\boldsymbol{\varepsilon}_i\}} \hat{\mathbf{h}}_i,
 \end{aligned} \tag{2.16}$$

with deterministic matrices  $\hat{T}_0, \hat{T}_1, \dots, \hat{T}_k$  and deterministic vectors  $\hat{\mathbf{h}}_0, \hat{\mathbf{h}}_1, \dots, \hat{\mathbf{h}}_k$ , where  $1_{\{\boldsymbol{\varepsilon}_i\}}$  is the indicator function that realisation  $\boldsymbol{\varepsilon}_i$  happens. We now have a **stochastic linear program with fixed recourse** defined as follows.

**Defn 2.4.** A *two-stage stochastic linear program with fixed recourse* has the following form:

$$\begin{aligned} \min \quad & w = \mathbf{c}^\top \mathbf{x} + \mathbb{E}_\zeta[Q(\mathbf{x}, \zeta)] \\ \text{such that} \quad & A\mathbf{x} = \mathbf{b}, \\ & \mathbf{x} \geq \mathbf{0}, \end{aligned} \quad (2.17)$$

where for each realisation  $\epsilon$  of  $\zeta$ ,

$$Q(\mathbf{x}, \epsilon) = \min_{\mathbf{y}} \{ \mathbf{q}^\top \mathbf{y} : W\mathbf{y} = \mathbf{h}(\epsilon) - T(\epsilon)\mathbf{x}, \mathbf{y} \geq \mathbf{0} \}. \quad (2.18)$$

Note that for each realisation  $\epsilon$  of  $\zeta$ ,  $Q(\mathbf{x}, \epsilon)$  is just a linear program

$$\begin{aligned} \min \quad & \mathbf{q}^\top \mathbf{y} \\ \text{subject to} \quad & W\mathbf{y} = \mathbf{h}(\epsilon) - T(\epsilon)\mathbf{x}, \\ & \mathbf{y} \geq \mathbf{0}. \end{aligned}$$

Recall that if  $\zeta$  has a finite discrete distribution  $\{(\epsilon_k, p_k), k = 1, \dots, r\}$ , where  $p_k > 0$  for all  $k$ , we can write it out in expanded form as follows.

**Defn 2.5.** The *expanded deterministic equivalent problem form* is:

$$\left. \begin{aligned} \min \quad w &= \mathbf{c}^\top \mathbf{x} + \sum_{k=1}^r p_k \mathbf{q}^\top \mathbf{y}_k \\ \text{subject to} \quad A\mathbf{x} &= \mathbf{b} \\ T(\epsilon_k)\mathbf{x} + W\mathbf{y}_k &= \mathbf{h}(\epsilon_k) \quad \text{for } k = 1, \dots, r \\ \mathbf{x}, \mathbf{y}_k &\geq \mathbf{0}, \quad \text{for } k = 1, \dots, r. \end{aligned} \right\} \quad (2.19)$$

Note that the above deterministic equivalent problem (DEP) is just a simple LP, which if small enough can be directly solved.

**Defn 2.6.** In particular, we have so-called **complete fixed recourse** if the fixed  $m_1 \times \bar{n}$  recourse matrix  $W$  satisfies

$$\{\mathbf{z} : \mathbf{z} = W\mathbf{y}, \mathbf{y} \geq \mathbf{0}\} = \mathbb{R}^{m_1},$$

where  $m_1$  is the number of stochastic constraints, and  $\bar{n}$  is the number of recourse variables (that is, the dimension of the recourse vector  $\mathbf{y}$ ).

This implies that, whatever the first-stage decision  $\mathbf{x}$  and realisation  $\boldsymbol{\varepsilon}$  of  $\boldsymbol{\zeta}$ , the second-stage program

$$Q(\mathbf{x}, \boldsymbol{\varepsilon}) = \min_{\mathbf{y}} \{\mathbf{q}^\top \mathbf{y} : W\mathbf{y} = \mathbf{h}(\boldsymbol{\varepsilon}) - T(\boldsymbol{\varepsilon})\mathbf{x}, \mathbf{y} \geq \mathbf{0}\}$$

will always be feasible.

**Defn 2.7.** A special case of complete fixed recourse is called **simple recourse** where

$$W = (I, -I)$$

with  $I$  the identity matrix of order  $m_1$ , and we have two non-negative second stage variables  $y^+, y^- \geq 0$  for each of the  $m_1$  first stage constraints.

In simple recourse, we also take account of the oversupply of commodities in the second stage by applying a penalty for such an occurrence as well as for the subsequent supply of the shortfall.

Here, the second stage reads as

$$Q(\mathbf{x}, \boldsymbol{\varepsilon}) = \begin{cases} \min & (\mathbf{q}^+)^T \mathbf{y}^+ + (\mathbf{q}^-)^T \mathbf{y}^- \\ \text{s.t.} & \mathbf{y}^+ - \mathbf{y}^- = \mathbf{h}(\boldsymbol{\varepsilon}) - T(\boldsymbol{\varepsilon})\mathbf{x}, \\ & \mathbf{y}^+, \mathbf{y}^- \geq \mathbf{0}, \end{cases}$$

where the unit cost of shortage for commodity  $i$  is given by  $q_i^+$  and the unit cost of oversupply for product  $j$  is given by  $q_j^-$ .

Then for  $\mathbf{q}^+ + \mathbf{q}^- \geq \mathbf{0}$  the recourse variables  $\mathbf{y}^+$  and  $\mathbf{y}^-$  can be chosen to measure (positively) the absolute deficiencies in the stochastic constraints. That is, the absolute measure of undersupply and oversupply.

**Example 26 (Simple recourse).** Consider

$$\begin{aligned} \min \quad & w = -x_2 \\ \text{subject to} \quad & x_1 + x_2 + x_3 = 2 \\ & a_1 x_1 + a_2 x_2 + x_4 = 2 \\ & -1 \leq x_1 \leq 1, x_i \geq 0, \text{ for } i = 2, 3, 4, \end{aligned}$$

where

$$(a_1, a_2) = \begin{cases} (1, 3/4) & \text{with probability } 1/2, \\ (-3, 5/4) & \text{with probability } 1/2. \end{cases}$$

Note that  $\mathbf{E}(a_1) = -1$ ,  $\mathbf{E}(a_2) = 1$ , and the deterministic LP with  $a_1 = -1$  and  $a_2 = 1$  has solution  $w^* = -2$ , with  $\mathbf{x}^* = (0, 2, 0, 0)^\top$ .

Suppose that we have a recourse policy that allows compensation (incurs a penalty ) of 5 per unit violation of the constraint

$$a_1 x_1 + a_2 x_2 + x_4 = 2$$

from its value 2. This penalty can be expressed as

$$5|2 - (a_1 x_1 + a_2 x_2 + x_4)|.$$

Thus, we can restate our problem in the following way

$$\begin{aligned} \min \quad & w = -x_2 + 5\mathbf{E}(|2 - (a_1 x_1 + a_2 x_2 + x_4)|) \\ \text{subject to} \quad & x_1 + x_2 + x_3 = 2 \\ & -1 \leq x_1 \leq 1, x_i \geq 0 \text{ for } i = 2, 3, 4, \end{aligned}$$

and the expectation may be written as

$$\frac{1}{2}(y_1^+ + y_1^-) + \frac{1}{2}(y_2^+ + y_2^-), \quad \text{where } y_j^+, y_j^- \geq 0 \text{ for } j = 1, 2.$$

Hence, in expanded form we may write

$$\begin{aligned}
& \min \quad w = -x_2 + \frac{5}{2}(y_1^+ + y_1^- + y_2^+ + y_2^-) \\
& \text{subject to} \quad x_1 + x_2 + x_3 = 2 \\
& \quad \quad \quad x_1 + 3/4x_2 + x_4 + y_1^+ - y_1^- = 2 \\
& \quad \quad \quad -3x_1 + 5/4x_2 + x_4 + y_2^+ - y_2^- = 2 \\
& \quad \quad \quad -1 \leq x_1 \leq 1, \quad x_i, y_j^+, y_j^- \geq 0, \text{ for } i = 2, 3, 4, \text{ and } j = 1, 2,
\end{aligned}$$

which has the optimal solution  $\mathbf{x}^* = (0.2222, 1.7778, 0.0000, 0.4444)^\top$  and the corresponding optimal value  $w^* = -1.7778$ .

As a two stage DEP, we may write

$$\begin{aligned}
& \min \quad w = -x_2 + 5\mathbf{E}(Q(\mathbf{x}, \boldsymbol{\zeta})) \\
& \text{s.t.} \quad x_1 + x_2 + x_3 = 2 \\
& \quad \quad -1 \leq x_1 \leq 1, x_i \geq 0 \text{ for } i = 2, 3, 4,
\end{aligned}$$

where

$$\mathbf{E}(Q(\mathbf{x}, \boldsymbol{\zeta})) = 1/2(Q(\mathbf{x}, \boldsymbol{\varepsilon}_1) + Q(\mathbf{x}, \boldsymbol{\varepsilon}_2)),$$

such that

$$\begin{aligned}
Q(\mathbf{x}, \boldsymbol{\varepsilon}_1) &: \min \quad y_1^+ + y_1^- \\
& \quad \text{s.t.} \quad y_1^+ - y_1^- = 2 - (x_1 + 3/4x_2 + x_4) \\
& \quad \quad \quad y_1^+, y_1^- \geq 0 \quad \text{and} \\
Q(\mathbf{x}, \boldsymbol{\varepsilon}_2) &: \min \quad y_2^+ + y_2^- \\
& \quad \text{s.t.} \quad y_2^+ - y_2^- = 2 - (-3x_1 + 5/4x_2 + x_4) \\
& \quad \quad \quad y_2^+, y_2^- \geq 0.
\end{aligned}$$

### RECOURSE DEP: PROPERTIES

A natural question which then arises is: what properties do these problems have, and in particular, are they of a class which can be solved? Some important properties to investigate further are *convexity*, and *smoothness*, such that

we have some chance to deal with them, and confidence in the results, using the toolkit of mathematical programming methods. We will now consider such issues.

*Convexity* for the recourse problem (2.6) can be shown under rather mild conditions, assuming the integrability of  $g_0 + Q$ .

**Proposition 2.1.** *If  $g_0(\mathbf{x}, \boldsymbol{\varepsilon})$  and  $Q(\mathbf{x}, \boldsymbol{\varepsilon})$  are convex in  $\mathbf{x}$  for each  $\boldsymbol{\varepsilon} \in \mathcal{S}_\zeta$ , where  $\mathcal{S}_\zeta$  is the support of  $\zeta$ , and if  $\mathcal{X}$  is a convex set, then (2.13),*

$$\min_{\mathbf{x} \in \mathcal{X}} \mathbf{E}_\zeta(f_0(\mathbf{x}, \zeta)) = \min_{\mathbf{x} \in \mathcal{X}} \mathbf{E}_\zeta(g_0(\mathbf{x}, \zeta) + Q(\mathbf{x}, \zeta))$$

*is a convex program.*

*Proof.* For  $\hat{\mathbf{x}}, \bar{\mathbf{x}} \in \mathcal{X}$  and  $\lambda \in (0, 1)$ , let  $\tilde{\mathbf{x}} := \lambda \hat{\mathbf{x}} + (1 - \lambda)\bar{\mathbf{x}}$ . We have

$$g_0(\tilde{\mathbf{x}}, \boldsymbol{\varepsilon}) + Q(\tilde{\mathbf{x}}, \boldsymbol{\varepsilon}) \leq \lambda [g_0(\hat{\mathbf{x}}, \boldsymbol{\varepsilon}) + Q(\hat{\mathbf{x}}, \boldsymbol{\varepsilon})] + (1 - \lambda) [g_0(\bar{\mathbf{x}}, \boldsymbol{\varepsilon}) + Q(\bar{\mathbf{x}}, \boldsymbol{\varepsilon})],$$

for all  $\boldsymbol{\varepsilon} \in \mathcal{S}_\zeta$ , since  $g_0(\mathbf{x}, \boldsymbol{\varepsilon})$  and  $Q(\mathbf{x}, \boldsymbol{\varepsilon})$  are convex in  $\mathbf{x}$  for all  $\boldsymbol{\varepsilon} \in \mathcal{S}_\zeta$ . This implies, by taking expectations, that

$$\begin{aligned} \mathbf{E}_\zeta [g_0(\tilde{\mathbf{x}}, \zeta) + Q(\tilde{\mathbf{x}}, \zeta)] &\leq \\ &\lambda \mathbf{E}_\zeta [g_0(\hat{\mathbf{x}}, \zeta) + Q(\hat{\mathbf{x}}, \zeta)] + (1 - \lambda) \mathbf{E}_\zeta [g_0(\bar{\mathbf{x}}, \zeta) + Q(\bar{\mathbf{x}}, \zeta)]. \end{aligned}$$

□

**Remark.** For  $\mathcal{Y} = \mathbb{R}_+^{\bar{n}}$  in the linear case (2.17) the convexity of  $Q(\mathbf{x}, \boldsymbol{\varepsilon})$  can immediately be asserted. For  $\mathcal{Y} = \mathbb{R}_+^{\bar{n}}$  in the nonlinear case (2.12) the convexity of  $Q(\mathbf{x}, \boldsymbol{\varepsilon})$  holds if the functions  $q(\mathbf{y})$  and  $g_i(\mathbf{x}, \boldsymbol{\varepsilon})$  are convex and the functions  $H_i(\mathbf{y})$  are concave.

A sketch of the latter argument is as follows. Assume that  $\hat{\mathbf{y}}$  and  $\bar{\mathbf{y}}$  solve (2.12) for  $\hat{\mathbf{x}}$  and  $\bar{\mathbf{x}}$  respectively, at some realisation  $\boldsymbol{\varepsilon} \in \mathcal{S}_\zeta$ . (That is,  $\hat{\mathbf{y}}$  and  $\bar{\mathbf{y}}$

are optimal solutions of (2.12) for  $\hat{\mathbf{x}}$  and  $\bar{\mathbf{x}}$  respectively, at some realisation  $\boldsymbol{\varepsilon} \in \mathcal{S}_\zeta$ .)

Then, by the convexity of  $g_i$  and concavity of  $H_i$ , we have, for any  $\lambda \in (0, 1)$ ,

$$\begin{aligned} g_i(\lambda \hat{\mathbf{x}} + (1 - \lambda) \bar{\mathbf{x}}, \boldsymbol{\varepsilon}) &\leq \lambda g_i(\hat{\mathbf{x}}, \boldsymbol{\varepsilon}) + (1 - \lambda) g_i(\bar{\mathbf{x}}, \boldsymbol{\varepsilon}) \\ &\leq \lambda H_i(\hat{\mathbf{y}}) + (1 - \lambda) H_i(\bar{\mathbf{y}}) \\ &\leq H_i(\lambda \hat{\mathbf{y}} + (1 - \lambda) \bar{\mathbf{y}}). \end{aligned}$$

Hence,  $\tilde{\mathbf{y}} = \lambda \hat{\mathbf{y}} + (1 - \lambda) \bar{\mathbf{y}}$  is feasible in (2.12) for  $\tilde{\mathbf{x}} = \lambda \hat{\mathbf{x}} + (1 - \lambda) \bar{\mathbf{x}}$ , and therefore by the convexity of  $q$ ,

$$\begin{aligned} Q(\tilde{\mathbf{x}}, \boldsymbol{\varepsilon}) &\leq q(\tilde{\mathbf{y}}) \\ &\leq \lambda q(\hat{\mathbf{y}}) + (1 - \lambda) q(\bar{\mathbf{y}}) \\ &= \lambda Q(\hat{\mathbf{x}}, \boldsymbol{\varepsilon}) + (1 - \lambda) Q(\bar{\mathbf{x}}, \boldsymbol{\varepsilon}). \end{aligned}$$

□

*Smoothness* (i.e. partial differentiability) of recourse problems may also be asserted under fairly general conditions. However, We will not consider this property in this course. See [3] for further details.

**Summary.** In summary, stochastic programs with recourse often have the properties of convexity, and given continuous distributions differentiability, both highly beneficial from a mathematical programming perspective.

### RECOURSE DEP: INDUCED CONSTRAINTS

Let  $\text{supp}(P)$  denote the support of the probability measure  $P$ , i.e. the smallest closed set  $\mathcal{S}$  such that  $P(\mathcal{S}) = 1$ . With the practical interpretation given to recourse problems in earlier lectures, and assuming  $\mathcal{S}_\zeta = \text{supp}(P_\zeta)$ , we should hope that for any first-stage decision  $\mathbf{x} \in \mathcal{X}$  there exists recourse variables  $\mathbf{y} \in \mathcal{Y}$  for all realisations  $\boldsymbol{\varepsilon} \in \mathcal{S}_\zeta$ . In other words, we hope that the

program

$$\begin{aligned} & \text{minimise} && Q(\mathbf{x}, \boldsymbol{\varepsilon}) = \mathbf{q}^\top \mathbf{y} \\ & \text{such that} && W\mathbf{y} = \mathbf{h}(\boldsymbol{\varepsilon}) - T(\boldsymbol{\varepsilon})\mathbf{x}, \\ & && \mathbf{y} \geq \mathbf{0} \end{aligned} \tag{2.20}$$

is feasible  $\forall \boldsymbol{\varepsilon} \in \mathcal{S}_\zeta$ .

However, depending upon the defined recourse matrix  $W$  and the given support  $\mathcal{S}_\zeta$ , there might not exist a feasible solution  $\mathbf{y} \in \mathcal{Y}$  for all first-stage decisions  $\mathbf{x} \in \mathcal{X}$ . Hence, it may be necessary to impose — in addition to  $\mathbf{x} \in \mathcal{X}$  — further restrictions on our first-stage decisions,  $\mathbf{x}$ , called **induced constraints**.

Let us assume that  $\mathcal{S}_\zeta$  is a convex polytope, i.e. the **convex hull of finitely many points**  $\boldsymbol{\varepsilon}_j \in \mathcal{S}_\zeta \subset \mathbb{R}^k$ ,  $j = 1, \dots, r$ :

$$\begin{aligned} \mathcal{S}_\zeta &= \text{conv}\{\boldsymbol{\varepsilon}_1, \dots, \boldsymbol{\varepsilon}_r\} \\ &= \left\{ \boldsymbol{\varepsilon} : \boldsymbol{\varepsilon} = \sum_{j=1}^r \lambda_j \boldsymbol{\varepsilon}_j, \sum_{j=1}^r \lambda_j = 1, \lambda_j \in [0, 1] \text{ for all } j = 1, \dots, r \right\}. \end{aligned}$$

It follows that  $\mathbf{x} \in \mathbb{R}^n$  allows for a feasible solution of the second stage program for all  $\boldsymbol{\varepsilon} \in \mathcal{S}_\zeta$  if and only if this is true for all  $\boldsymbol{\varepsilon}_j$  for  $j = 1, \dots, r$ . In other words, the corresponding induced first-stage feasibility set  $\mathcal{K}$  is given as

$$\mathcal{K} = \left\{ \mathbf{x} : T(\boldsymbol{\varepsilon}_j)\mathbf{x} + W\mathbf{y}_j = \mathbf{h}(\boldsymbol{\varepsilon}_j), \mathbf{y}_j \geq \mathbf{0}, j = 1, \dots, r \right\}.$$

**Theorem 2.1.** *If the support  $\mathcal{S}_\zeta$  of  $\zeta$  is either a finite set or a convex polytope, then the **induced first-stage feasibility set**  $\mathcal{K}$  is a convex polytope.*

The first-stage decision variables are then restricted to

$$\mathbf{x} \in \mathcal{X} \cap \mathcal{K}.$$



**Example 27.** Consider the following first-stage feasibility set

$$\mathcal{X} = \{\mathbf{x} \in \mathbb{R}_+^2 : x_1 - 2x_2 \geq -4, x_1 + 2x_2 \leq 8, 2x_1 - x_2 \leq 6\}.$$

For the second-stage constraints, we have

$$W = \begin{pmatrix} -1 & 3 & 5 \\ 2 & 2 & 2 \end{pmatrix}, \quad T(\boldsymbol{\varepsilon}) = T = \begin{pmatrix} 2 & 3 \\ 3 & 1 \end{pmatrix}, \quad \mathbf{h}(\boldsymbol{\varepsilon}) = \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \end{pmatrix},$$

and assume  $\boldsymbol{\zeta}$  has support  $\mathcal{S}_{\boldsymbol{\zeta}} = [4, 19] \times [13, 21]$ . Then, the constraints to be satisfied for all  $\boldsymbol{\varepsilon} \in \mathcal{S}_{\boldsymbol{\zeta}}$  are

$$W\mathbf{y} = \mathbf{h}(\boldsymbol{\varepsilon}) - T\mathbf{x}, \mathbf{y} \geq \mathbf{0}. \quad (2.21)$$

Note that the second column  $W_{:,2}$  of  $W$  is a linear combination of columns 1 and 3, that is,

$$W_{:,2} = \frac{1}{3}W_{:,1} + \frac{2}{3}W_{:,3},$$

and so (2.21) can be reduced to

$$\mathbf{h}(\boldsymbol{\varepsilon}) - T\mathbf{x} = \lambda W_{:,1} + \mu W_{:,3} \quad \text{for } \lambda, \mu \geq 0.$$

That is,

$$\begin{aligned} \varepsilon_1 - 2x_1 - 3x_2 &= -\lambda + 5\mu, \\ \varepsilon_2 - 3x_1 - x_2 &= 2\lambda + 2\mu, \\ \lambda, \mu &\geq 0, \quad \text{for all } \boldsymbol{\varepsilon} \in \mathcal{S}_{\boldsymbol{\zeta}}. \end{aligned} \quad (2.22)$$

Now, by multiplying the constraints using the transformation matrix

$$S = \begin{pmatrix} 2 & 1 \\ -2 & 5 \end{pmatrix},$$

we obtain

$$S\mathbf{h}(\boldsymbol{\varepsilon}) - ST\mathbf{x} = S \begin{pmatrix} -\lambda & 5\mu \\ 2\lambda & 2\mu \end{pmatrix},$$

corresponding to adding two times the first equation to the second and  $-2$  the first equations to 5 times the second, respectively. Then, we arrive at the *equivalent* system

$$\begin{aligned} 2\varepsilon_1 + \varepsilon_2 - 7x_1 - 7x_2 &= 12\mu, \\ -2\varepsilon_1 + 5\varepsilon_2 - 11x_1 + x_2 &= 12\lambda, \\ \lambda, \mu &\geq 0, \forall \varepsilon \in \mathcal{S}_\zeta. \end{aligned} \quad (2.23)$$

Furthermore, by noting that (2.23) has to hold for some  $\lambda, \mu \geq 0$ , we can rearrange this system of inequalities into an equivalent system for which we can propose some solutions. In particular, let  $\lambda = \mu = 0$ , we obtain

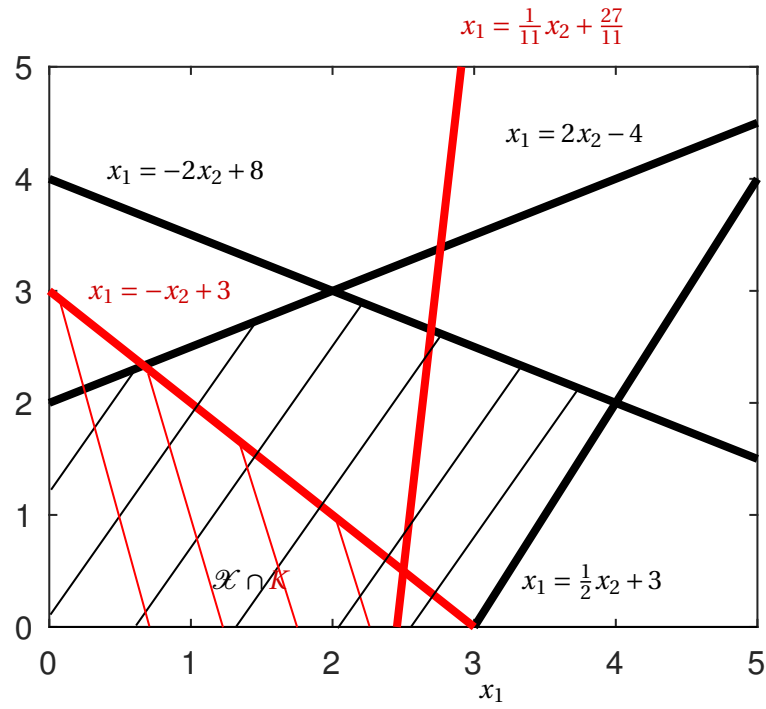
$$\begin{aligned} 7x_1 + 7x_2 &\leq 2\varepsilon_1 + \varepsilon_2, \\ 11x_1 - x_2 &\leq -2\varepsilon_1 + 5\varepsilon_2, \end{aligned}$$

for all  $\varepsilon \in \mathcal{S}_\zeta$ . Then, as these inequalities have to hold for all  $\varepsilon \in \mathcal{S}_\zeta$ , we can choose the minimal right-hand sides to give the induced constraints

$$\mathcal{K} = \{\mathbf{x} : 7x_1 + 7x_2 \leq 21, 11x_1 - x_2 \leq 27\}.$$

Graphically, the first-stage feasibility set  $\mathcal{X}$  together with the induced feasible set  $\mathcal{K}$  is shown in Figure 2.7. The restricted feasibility set for the problem is then  $\mathcal{X} \cap \mathcal{K}$ .

Finally, for *complete fixed recourse*, for which  $\mathcal{K} = \mathbb{R}^n$ , the problem of induced constraints does not exist, since  $(\mathcal{X} \cap \mathcal{K}) = \mathcal{X}$ . This leads to the question: when do we have a stochastic program with complete fixed recourse? The following proposition addresses this question.



**Figure 2.7:** The first-stage feasibility set  $\mathcal{X}$  together with the induced feasible set  $\mathcal{K}$ .

**Proposition 2.2.** *An  $m_1 \times \bar{n}$  matrix  $W$  is a complete recourse matrix if and only if*

1. *it has  $\text{rank}(W) = m_1$ , and*
2. *assuming without loss of generality that its first  $m_1$  columns  $W_{:,1}, \dots, W_{:,m_1}$  are linearly independent, the linear constraints*

$$\begin{aligned} W\mathbf{y} &= \mathbf{0} \\ y_i &\geq 1 \quad \text{for } i = 1, \dots, m_1 \\ \mathbf{y} &\geq \mathbf{0} \end{aligned} \tag{2.24}$$

*have a feasible solution.*

*Proof.*  $W$  is a complete recourse matrix if and only if

$$\{\mathbf{z} : \mathbf{z} = W\mathbf{y}, \mathbf{y} \geq \mathbf{0}\} = \mathbb{R}^{m_1}.$$

From this condition, it follows immediately that  $\text{rank}(W) = m_1$  necessarily has to hold. In addition, choosing

$$\hat{\mathbf{z}} = - \sum_{i=1}^{m_1} W_{:,i} \in \mathbb{R}^{m_1},$$

the second-stage constraints

$$W\mathbf{y} = \hat{\mathbf{z}}, \mathbf{y} \geq \mathbf{0},$$

by the property of the complete recourse matrix  $W$ , are guaranteed to have at

least one feasible solution, say,  $\hat{\mathbf{y}}$ , such that

$$\begin{aligned} \sum_{i=1}^{m_1} W_{:,i} \hat{y}_i + \sum_{i=m_1+1}^{\bar{n}} W_{:,i} \hat{y}_i &= \bar{\mathbf{z}} \\ &= - \sum_{i=1}^{m_1} W_{:,i}, \\ \hat{y}_i &\geq 0, \quad i = 1, \dots, \bar{n}. \end{aligned}$$

Thus, with

$$\tilde{y}_i = \begin{cases} \hat{y}_i + 1 & i = 1, \dots, m_1, \\ \hat{y}_i & i = m_1 + 1, \dots, \bar{n}, \end{cases}$$

the conditions (2.24) are *necessarily feasible*.

To show *sufficiency*, let us choose an arbitrary  $\bar{\mathbf{z}} \in \mathbb{R}^{m_1}$ . Since the columns  $W_{:,1}, \dots, W_{:,m_1}$  are linearly independent, the system of linear equations

$$\sum_{i=1}^{m_1} W_{:,i} y_i = \bar{\mathbf{z}},$$

has a unique solution  $\bar{y}_1, \dots, \bar{y}_{m_1}$ . If  $\bar{y}_i \geq 0$ ,  $i = 1, \dots, m_1$ , we are done. Otherwise, define

$$\gamma := \min\{\bar{y}_1, \dots, \bar{y}_{m_1}\}.$$

Then  $\mathbf{y}^*$  defined by

$$y_i^* = \begin{cases} \bar{y}_i - \gamma \tilde{y}_i & i = 1, \dots, m_1, \\ -\gamma \tilde{y}_i & i = m_1 + 1, \dots, \bar{n}, \end{cases}$$

solves  $W\mathbf{y} = \bar{\mathbf{z}}$ ,  $\mathbf{y} \geq \mathbf{0}$ . □

**Defn 2.8.** If the program (2.20) is feasible for all  $\boldsymbol{\varepsilon} \in \mathcal{S}_\zeta$  and at least **all**  $\mathbf{x} \in \mathcal{X} = \{\mathbf{x} | A\mathbf{x} = \mathbf{b}, \mathbf{x} \geq \mathbf{0}\}$ , then (2.17) is said to be of *relatively complete recourse*.

### §2.5 CHANCE CONSTRAINED DEP

We will now consider another class of DEPs for (2.5), called **chance constrained stochastic programs**.

Consider a DEP with a collection of constraints

$$g_i(\mathbf{x}, \boldsymbol{\zeta}) \leq 0, \quad i = 1, \dots, m, \quad \text{for } \mathbf{x} \in \mathcal{X} \subset \mathbb{R}^n.$$

For each realisation  $\boldsymbol{\varepsilon}$  of  $\boldsymbol{\zeta}$  and for some  $\alpha \in [0, 1]$ , we define a *payoff* function for all constraints as

$$\varphi(\mathbf{x}, \boldsymbol{\varepsilon}) := \begin{cases} 1 - \alpha & \text{if } g_i(\mathbf{x}, \boldsymbol{\varepsilon}) \leq 0, \forall i = 1, \dots, m, \\ -\alpha & \text{otherwise.} \end{cases} \quad (2.25)$$

Consequently, with each realisation  $\boldsymbol{\varepsilon}$  of  $\boldsymbol{\zeta}$ , for  $\mathbf{x}$  infeasible we have an *absolute loss* of  $\alpha$ , whereas for  $\mathbf{x}$  feasible we have *reward* of  $1 - \alpha$ . With this payoff function, it is natural to aim for decisions on  $\mathbf{x}$  that, at least on average, avoid an absolute loss. That is, we could specify the requirement

$$\mathbb{E}_{\boldsymbol{\zeta}}[\varphi(\mathbf{x}, \boldsymbol{\zeta})] \geq 0.$$

Define

$$\mathbf{g}(\mathbf{x}, \boldsymbol{\varepsilon}) := (g_1(\mathbf{x}, \boldsymbol{\varepsilon}), \dots, g_m(\mathbf{x}, \boldsymbol{\varepsilon}))^\top.$$

Note that we write ' $\mathbf{m} \leq \mathbf{0}$ ' to indicate that every component of a vector  $\mathbf{m}$  is less than or equal to 0, and ' $\mathbf{m} \not\leq \mathbf{0}$ ' to indicate that *at least one* component of vector  $\mathbf{m}$  is *not* less than or equal to 0.

If  $\zeta$  is a **discrete** random vector, we can write

$$\begin{aligned}
\mathbf{E}_\zeta[\varphi(\mathbf{x}, \zeta)] &= \sum_{\boldsymbol{\varepsilon} \in \mathcal{S}_\zeta} \varphi(\mathbf{x}, \boldsymbol{\varepsilon}) P(\boldsymbol{\varepsilon}) \\
&= \left( \sum_{\{\boldsymbol{\varepsilon} : \mathbf{g}(\mathbf{x}, \boldsymbol{\varepsilon}) \leq \mathbf{0}\}} (1 - \alpha) P(\boldsymbol{\varepsilon}) \right) - \left( \sum_{\{\boldsymbol{\varepsilon} : \mathbf{g}(\mathbf{x}, \boldsymbol{\varepsilon}) \not\leq \mathbf{0}\}} \alpha P(\boldsymbol{\varepsilon}) \right) \\
&= (1 - \alpha) \left( \sum_{\{\boldsymbol{\varepsilon} : \mathbf{g}(\mathbf{x}, \boldsymbol{\varepsilon}) \leq \mathbf{0}\}} P(\boldsymbol{\varepsilon}) \right) - \alpha \left( \sum_{\{\boldsymbol{\varepsilon} : \mathbf{g}(\mathbf{x}, \boldsymbol{\varepsilon}) \not\leq \mathbf{0}\}} P(\boldsymbol{\varepsilon}) \right) \\
&= (1 - \alpha) P(\{\boldsymbol{\varepsilon} : \mathbf{g}(\mathbf{x}, \boldsymbol{\varepsilon}) \leq \mathbf{0}\}) - \alpha P(\{\boldsymbol{\varepsilon} : \mathbf{g}(\mathbf{x}, \boldsymbol{\varepsilon}) \not\leq \mathbf{0}\}) \\
&= P(\{\boldsymbol{\varepsilon} : \mathbf{g}(\mathbf{x}, \boldsymbol{\varepsilon}) \leq \mathbf{0}\}) - \alpha \{P(\{\boldsymbol{\varepsilon} : \mathbf{g}(\mathbf{x}, \boldsymbol{\varepsilon}) \leq \mathbf{0}\}) + P(\{\boldsymbol{\varepsilon} : \mathbf{g}(\mathbf{x}, \boldsymbol{\varepsilon}) \not\leq \mathbf{0}\})\} \\
&= P(\{\boldsymbol{\varepsilon} : \mathbf{g}(\mathbf{x}, \boldsymbol{\varepsilon}) \leq \mathbf{0}\}) - \alpha.
\end{aligned}$$

Therefore, the condition  $\mathbf{E}_\zeta[\varphi(\mathbf{x}, \zeta)] \geq 0$  is equivalent to requiring

$$P(\{\boldsymbol{\varepsilon} : \mathbf{g}(\mathbf{x}, \boldsymbol{\varepsilon}) \leq \mathbf{0}\}) \geq \alpha.$$

Note that the same applies when  $\zeta$  is a **continuous** random vector. In particular, we can write

$$\begin{aligned}
\mathbf{E}_\zeta[\varphi(\mathbf{x}, \zeta)] &= \int_{\mathcal{S}_\zeta} \varphi(\mathbf{x}, \boldsymbol{\varepsilon}) dP(\boldsymbol{\varepsilon}) \\
&= \int_{\{\boldsymbol{\varepsilon} : \mathbf{g}(\mathbf{x}, \boldsymbol{\varepsilon}) \leq \mathbf{0}\}} (1 - \alpha) dP(\boldsymbol{\varepsilon}) - \int_{\{\boldsymbol{\varepsilon} : \mathbf{g}(\mathbf{x}, \boldsymbol{\varepsilon}) \not\leq \mathbf{0}\}} \alpha dP(\boldsymbol{\varepsilon}) \\
&= (1 - \alpha) P(\{\boldsymbol{\varepsilon} : \mathbf{g}(\mathbf{x}, \boldsymbol{\varepsilon}) \leq \mathbf{0}\}) - \alpha P(\{\boldsymbol{\varepsilon} : \mathbf{g}(\mathbf{x}, \boldsymbol{\varepsilon}) \not\leq \mathbf{0}\}) \\
&= P(\{\boldsymbol{\varepsilon} : \mathbf{g}(\mathbf{x}, \boldsymbol{\varepsilon}) \leq \mathbf{0}\}) - \alpha,
\end{aligned}$$

which would lead us to the same conclusion, that  $\mathbf{E}_\zeta[\varphi(\mathbf{x}, \zeta)] \geq 0$  is equivalent to requiring

$$P(\{\boldsymbol{\varepsilon} : \mathbf{g}(\mathbf{x}, \boldsymbol{\varepsilon}) \leq \mathbf{0}\}) \geq \alpha.$$

**Defn 2.9.** A *probabilistically constrained program*, also called a *chance constrained program* or a *problem with joint probabilistic constraints*, is defined as

$$\begin{aligned} \min \quad & w = \mathbf{E}_{\zeta}[f_0(\mathbf{x}, \zeta)] \\ \text{such that} \quad & P(\{\boldsymbol{\varepsilon} : g_i(\mathbf{x}, \boldsymbol{\varepsilon}) \leq 0, i = 1, \dots, m\}) \geq \alpha. \end{aligned} \quad (2.26)$$

Instead of the payoff function (2.25), we might choose to define  $\alpha_i \in [0, 1]$ ,  $i = 1, \dots, m$ , and analogous payoffs for every single constraint, thus allowing us to weigh certain constraints as more important than others. This results in

$$\varphi_i(\mathbf{x}, \boldsymbol{\varepsilon}) := \begin{cases} 1 - \alpha_i & \text{if } g_i(\mathbf{x}, \boldsymbol{\varepsilon}) \leq 0, \\ \alpha_i & \text{otherwise,} \end{cases} \quad (2.27)$$

for  $i = 1, \dots, m$ . Then we get the following DEP.

**Defn 2.10.** A *single probabilistic/chance constrained problem*, or *separate probabilistic/chance constrained problem*, is defined as follows:

$$\begin{aligned} \min_{\mathbf{x} \in \mathcal{X}} \quad & w = \mathbf{E}_{\zeta}[f_0(\mathbf{x}, \zeta)] \\ \text{such that} \quad & P(\{\boldsymbol{\varepsilon} : g_i(\mathbf{x}, \boldsymbol{\varepsilon}) \leq 0\}) \geq \alpha_i, i = 1, \dots, m. \end{aligned} \quad (2.28)$$

There are *many* other possibilities to generate types of DEPs for (2.5) by constructing the payoff functions  $\varphi_i$  in different ways, out of the objective function and the constraints. All the problems derived here are *mathematical programs*.



### CHANCE CONSTRAINED DEP: PROPERTIES

For chance constrained programs, things are generally more difficult than for recourse problems. Consider the constraint

$$P(\{\boldsymbol{\varepsilon} : \mathbf{g}(\mathbf{x}, \boldsymbol{\varepsilon}) \leq \mathbf{0}\}) \geq \alpha$$

of (2.26) with the functions  $g_i$  replaced by the vector-valued function

$$\mathbf{g}(\mathbf{x}, \boldsymbol{\varepsilon}) := (g_1(\mathbf{x}, \boldsymbol{\varepsilon}), \dots, g_m(\mathbf{x}, \boldsymbol{\varepsilon}))^\top.$$

A point  $\hat{\mathbf{x}}$  is feasible iff

$$P(\{\boldsymbol{\varepsilon} : \mathbf{g}(\hat{\mathbf{x}}, \boldsymbol{\varepsilon}) \leq \mathbf{0}\}) \geq \alpha;$$

that is,  $\hat{\mathbf{x}}$  is feasible iff the set (event)

$$\mathcal{S}(\hat{\mathbf{x}}) = \{\boldsymbol{\varepsilon} : \mathbf{g}(\hat{\mathbf{x}}, \boldsymbol{\varepsilon}) \leq \mathbf{0}\}$$

has a probability measure of at least  $\alpha$ .

Another way to define a feasible point  $\hat{\mathbf{x}}$  is as follows. Let  $\mathbb{G}$  be a collection of all events such that

$$P(G) \geq \alpha \text{ for all } G \in \mathbb{G}.$$

Then,  $\hat{\mathbf{x}}$  is feasible iff we have at least one event  $\tilde{G} \in \mathbb{G}$  such that for all  $\boldsymbol{\varepsilon} \in \tilde{G}$ ,  $\mathbf{g}(\hat{\mathbf{x}}, \boldsymbol{\varepsilon}) \leq \mathbf{0}$ . Formally,  $\hat{\mathbf{x}}$  is feasible iff there exists  $\tilde{G} \in \mathbb{G}$  such that

$$\hat{\mathbf{x}} \in \bigcap_{\boldsymbol{\varepsilon} \in \tilde{G}} \{\mathbf{x} : \mathbf{g}(\mathbf{x}, \boldsymbol{\varepsilon}) \leq \mathbf{0}\}. \quad (2.29)$$

Hence, the feasible set

$$B(\alpha) = \{\mathbf{x} : P(\{\boldsymbol{\varepsilon} : \mathbf{g}(\mathbf{x}, \boldsymbol{\varepsilon}) \leq \mathbf{0}\}) \geq \alpha\}$$

is the union of all those vectors  $\mathbf{x}$  feasible according to (2.29). Consequently,  $B(\alpha)$  may be written as

$$B(\alpha) = \bigcup_{G \in \mathbb{G}} \bigcap_{\boldsymbol{\varepsilon} \in G} \{\mathbf{x} : \mathbf{g}(\mathbf{x}, \boldsymbol{\varepsilon}) \leq \mathbf{0}\}.$$

Since a union of convex sets is not necessarily convex, this demonstrates that in general we may not expect  $B(\alpha)$  to be convex, even if  $\{\mathbf{x} : \mathbf{g}(\mathbf{x}, \boldsymbol{\varepsilon}) \leq \mathbf{0}\}$  is convex for all  $\boldsymbol{\varepsilon} \in \mathcal{S}_\zeta$ . Indeed, there are simple examples of nonconvex feasible sets.

**Example 28.** Assume that in our refinery example demands are random with the following joint discrete distribution:

$$\begin{aligned} P(h_s(\boldsymbol{\varepsilon}_1) = 1600, h_p(\boldsymbol{\varepsilon}_1) = 1350) &= 0.85, \\ P(h_s(\boldsymbol{\varepsilon}_2) = 1500, h_p(\boldsymbol{\varepsilon}_2) = 1950) &= 0.08, \\ P(h_s(\boldsymbol{\varepsilon}_3) = 2000, h_p(\boldsymbol{\varepsilon}_3) = 1200) &= 0.07. \end{aligned}$$

Recall the constraints  $x_1 + x_2 \leq 100, x_1, x_2 \geq 0$ , but this time with the probabilistic constraints

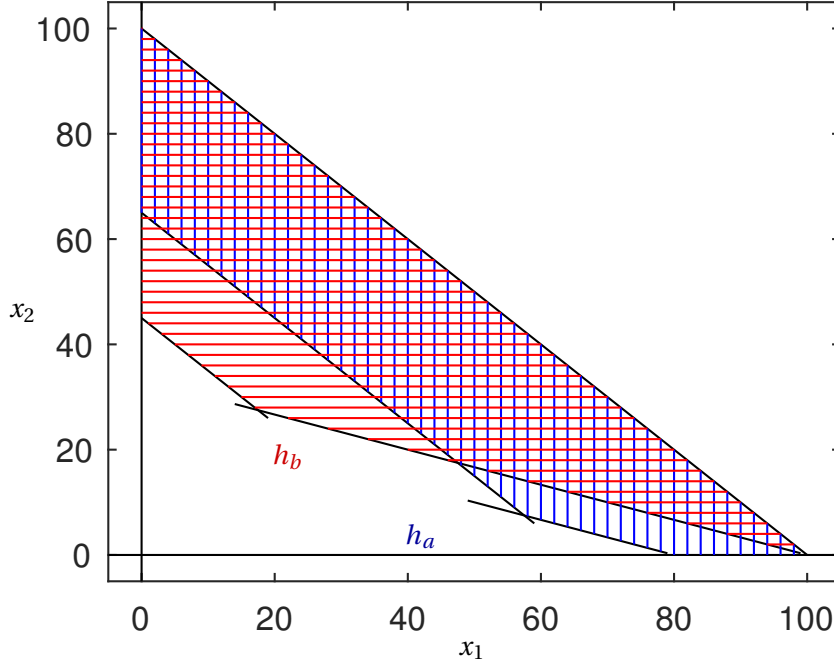
$$P(20x_1 + 60x_2 \geq h_s(\boldsymbol{\varepsilon}), 30x_1 + 30x_2 \geq h_p(\boldsymbol{\varepsilon})) \geq \alpha,$$

for any  $\alpha \in (0.85, 0.92]$ . These constraints, combined together, require that we either

1. either satisfy the demands  $h_s(\boldsymbol{\varepsilon}_1)$  and  $h_p(\boldsymbol{\varepsilon}_2)$ , thus enforcing a reliability of 93%, and hence choose a production program to cover a demand  $h_a = (1600; 1950)$ ;
2. or satisfy the demands  $h_s(\boldsymbol{\varepsilon}_3)$  and  $h_p(\boldsymbol{\varepsilon}_1)$ , thus enforcing a reliability of 92%, and hence choose a production program to cover a demand  $h_b = (2000; 1350)$ .

The above scenarios each has a convex solution space, but the feasible set for the union of the above constraints can clearly be seen to be nonconvex in Figure 2.8.

A natural question arises: when is  $B(\alpha)$  convex? The following proposition addresses this question.



**Figure 2.8:** Chance constraints can lead to a non-convex feasible set

**Proposition 2.3.** If  $\mathbf{g}(\mathbf{x}, \boldsymbol{\varepsilon})$  is jointly convex and the probability measure  $P$  is quasi-concave, then for all  $\alpha \in [0, 1]$ , the feasible set

$$B(\alpha) = \{\mathbf{x} : P(\{\boldsymbol{\varepsilon} : \mathbf{g}(\mathbf{x}, \boldsymbol{\varepsilon}) \leq \mathbf{0}\}) \geq \alpha\}$$

is convex.

Before we go into the proof, we need to define the concept of a probability measure  $P$  being quasi-concave. Suppose that the random vector  $\boldsymbol{\zeta}$  lives on the probability space  $(\Omega, \mathbb{F}, P)$ , so  $\mathbb{F}$  is the  $\sigma$ -field on which  $P$  is defined. Then,

$P$  is said to be quasi-concave if

$$P(\lambda \mathcal{S}_1 + (1 - \lambda) \mathcal{S}_2) \geq \min\{P(\mathcal{S}_1), P(\mathcal{S}_2)\} \quad \text{for all } \lambda \in (0, 1)$$

for all convex sets  $\mathcal{S}_i \in \mathbb{F}, i = 1, 2$ .

*Proof.* Define

$$\mathcal{S}(\mathbf{x}) := \{\boldsymbol{\varepsilon} : \mathbf{g}(\mathbf{x}, \boldsymbol{\varepsilon}) \leq \mathbf{0}\}.$$

For  $\mathbf{x}_1, \mathbf{x}_2 \in B(\alpha)$ ,  $\boldsymbol{\varepsilon}_1 \in \mathcal{S}(\mathbf{x}_1)$ ,  $\boldsymbol{\varepsilon}_2 \in \mathcal{S}(\mathbf{x}_2)$ , and  $\lambda \in (0, 1)$ , define

$$(\bar{\mathbf{x}}, \bar{\boldsymbol{\varepsilon}}) := \lambda(\mathbf{x}_1, \boldsymbol{\varepsilon}_1) + (1 - \lambda)(\mathbf{x}_2, \boldsymbol{\varepsilon}_2);$$

In other words, define  $\bar{\mathbf{x}} := \lambda \mathbf{x}_1 + (1 - \lambda) \mathbf{x}_2$  and  $\bar{\boldsymbol{\varepsilon}} := \lambda \boldsymbol{\varepsilon}_1 + (1 - \lambda) \boldsymbol{\varepsilon}_2$ . If  $\mathbf{g}(\mathbf{x}, \boldsymbol{\varepsilon})$  is jointly convex, then it follows that

$$\mathbf{g}(\bar{\mathbf{x}}, \bar{\boldsymbol{\varepsilon}}) \leq \lambda \mathbf{g}(\mathbf{x}_1, \boldsymbol{\varepsilon}_1) + (1 - \lambda) \mathbf{g}(\mathbf{x}_2, \boldsymbol{\varepsilon}_2) \leq \mathbf{0},$$

where the second inequality is due to the assumptions that  $\boldsymbol{\varepsilon}_1 \in \mathcal{S}(\mathbf{x}_1)$  and  $\boldsymbol{\varepsilon}_2 \in \mathcal{S}(\mathbf{x}_2)$ . The inequality  $\mathbf{g}(\bar{\mathbf{x}}, \bar{\boldsymbol{\varepsilon}}) \leq \mathbf{0}$  implies  $\bar{\boldsymbol{\varepsilon}} \in \mathcal{S}(\bar{\mathbf{x}})$ , and hence

$$\lambda \mathcal{S}(\mathbf{x}_1) + (1 - \lambda) \mathcal{S}(\mathbf{x}_2) \subset \mathcal{S}(\bar{\mathbf{x}}), \quad (2.30)$$

where the addition of two sets  $\mathcal{A}$  and  $\mathcal{B}$  is performed as follows:

$$\rho \mathcal{A} + \sigma \mathcal{B} := \{\boldsymbol{\varepsilon} := \rho \boldsymbol{\varepsilon}_1 + \sigma \boldsymbol{\varepsilon}_2 : \boldsymbol{\varepsilon}_1 \in \mathcal{A}, \boldsymbol{\varepsilon}_2 \in \mathcal{B}\}.$$

Equation (2.30) implies

$$P(\mathcal{S}(\bar{\mathbf{x}})) \geq P(\lambda \mathcal{S}(\mathbf{x}_1) + (1 - \lambda) \mathcal{S}(\mathbf{x}_2)).$$

By the joint convexity on  $\mathbf{g}(\mathbf{x}, \boldsymbol{\varepsilon})$ , any set  $\mathcal{S}(\mathbf{x})$  is convex. (For any  $\mathbf{x}$ , take  $\boldsymbol{\varepsilon}_1, \boldsymbol{\varepsilon}_2 \in \mathcal{S}(\mathbf{x})$  and show that  $\lambda \boldsymbol{\varepsilon}_1 + (1 - \lambda) \boldsymbol{\varepsilon}_2 \in \mathcal{S}(\mathbf{x})$ .) As by assumption  $P$  is quasi-concave, we have

$$P(\lambda \mathcal{S}(\mathbf{x}_1) + (1 - \lambda) \mathcal{S}(\mathbf{x}_2)) \geq \min\{P(\mathcal{S}(\mathbf{x}_1)), P(\mathcal{S}(\mathbf{x}_2))\} \geq \alpha \quad \text{for all } \lambda \in (0, 1)$$

for all convex sets  $\mathcal{S}_i \in \mathbb{F}, i = 1, 2$ .

Thus,  $P(\mathcal{S}(\bar{\mathbf{x}})) \geq \alpha$ , which implies  $\bar{\mathbf{x}} \in B(\alpha)$  and therefore  $B(\alpha)$  is convex for all  $\alpha \in (0, 1)$   $\square$

Therefore, we see that  $P$  being quasi-concave assists us in establishing a nice property of chance constrained DEPs. Hence, we would like to be able to classify when a probability measure  $P$  is quasi-concave.

**Theorem 2.2.** *If the probability measure  $P$  is quasi-concave, then the corresponding distribution function  $F_\zeta$  is quasi-concave.*

This follows from the definition of distribution functions:

$$F_\zeta(\epsilon_i) = P(\mathcal{S}_i)$$

with  $\mathcal{S}_i = \{\epsilon : \epsilon \leq \epsilon_i\}$ ,  $i = 1, 2$ , and from the fact that for  $\hat{\epsilon} = \lambda\epsilon_1 + (1 - \lambda)\epsilon_2$ ,  $\lambda \in (0, 1)$  we have  $\widehat{\mathcal{S}} = \{\epsilon : \epsilon \leq \hat{\epsilon}\} = \lambda\mathcal{S}_1 + (1 - \lambda)\mathcal{S}_2$ .

When  $P$  is quasi-concave, this gives

$$F_\zeta(\hat{\epsilon}) = P(\widehat{\mathcal{S}}) \geq \min\{P(\mathcal{S}_1), P(\mathcal{S}_2)\} = \min\{F_\zeta(\epsilon_1), F_\zeta(\epsilon_2)\}.$$

Unfortunately, the converse is not true:  $F_\zeta$  being quasi-concave does not imply, in general, that the corresponding probability measure  $P$  is quasi-concave. **[Example in class]**

## §2.6 LAGRANGE MULTIPLIERS, LPS, AND DUALITY

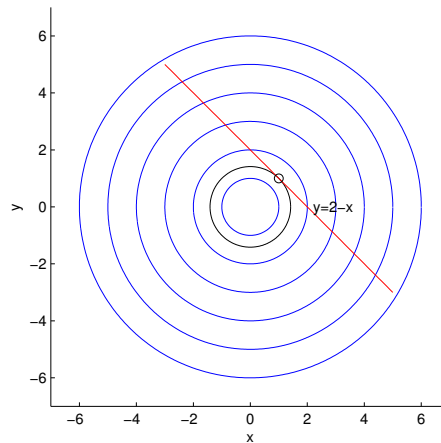
Before looking at the Dual Decomposition Method, We will revisit the relationship between Linear Programs, Duality and Lagrange multipliers. **Lagrange multipliers** provide another method for solving constrained optimisation problems and because of this, We will discover a direct relationship between them and the Dual of an LP.

**Example 29.** Consider the following non-linear optimisation problem:

$$\begin{aligned} \min \quad & f(x, y) = x^2 + y^2 \\ \text{such that} \quad & g(x, y) = x + y - 2 = 0. \end{aligned} \tag{2.31}$$

The solution is clearly  $x = y = 1$  with  $f(x, y) = 2$ .

The contour of the objective function  $f(x, y)$  is tangent to the constraint surface at the solution, as can be seen in Figure 2.9.



**Figure 2.9:** Contours of the objective function and constraint

The simplest version of the **Lagrange Multiplier Theorem** says that this is always true with equality constraints: *at the constraint optimum, if it exists,  $\nabla f$  will be a multiple of  $\nabla g$ .*

The Lagrange Multiplier Theorem enables us translate the original constrained optimisation problem into an ordinary system of simultaneous equations at the cost of introducing a new variable; that is,

$$g(x, y) = 0, \quad (2.32)$$

$$\nabla f(x, y) = -p \nabla g(x, y), \quad (2.33)$$

where Equation (2.32) says we must satisfy the original constraint (2.31), and Equation (2.33) adds that  $\nabla f$  must be parallel to  $\nabla g$ . The new variable  $p$  is called the **Lagrange multiplier**.

This system can be written more compactly by writing it in terms of the Lagrangian  $L$ , where

$$\begin{aligned} L(x, y, p) &= f(x, y) + pg(x, y) \\ &= x^2 + y^2 + p(x + y - 2). \end{aligned}$$

We then set  $\nabla L = \mathbf{0}$ , that is,

$$\begin{aligned} \nabla_p L(x, y, p) &= x + y - 2 = 0, \\ \nabla_x L(x, y, p) &= 2x + p = 0, \\ \nabla_y L(x, y, p) &= 2y + p = 0. \end{aligned}$$

The unique solution of the above system of equations is  $x = y = 1$  as before, with  $p = -2$ .

This technique applies to more than problems with single constraints, by introducing more than one Lagrange multiplier.

**Example 30.** Consider the following optimisation.

$$\begin{aligned} \min \quad & f(x, y, z) = x^2 + y^2 + z^2, \\ \text{such that} \quad & x + y - 2 = 0, \\ & x + z - 2 = 0. \end{aligned}$$

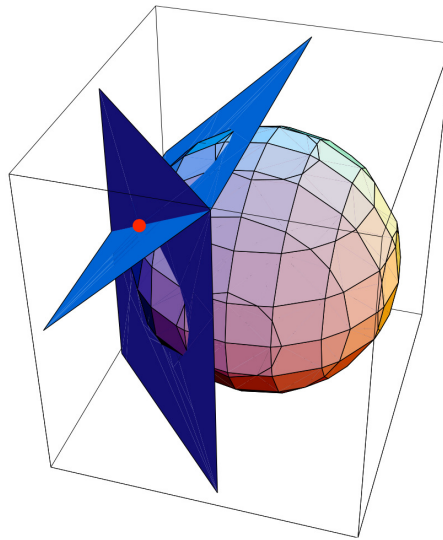
Then,  $L(x, y, z, p, q) = x^2 + y^2 + z^2 + p(x + y - 2) + q(x + z - 2)$ , and setting  $\nabla L = \mathbf{0}$  gives

$$\begin{aligned} \nabla_p L(x, y, z, p, q) &= x + y - 2 = 0, \\ \nabla_q L(x, y, z, p, q) &= x + z - 2 = 0, \\ \nabla_x L(x, y, z, p, q) &= 2x + p + q = 0, \\ \nabla_y L(x, y, z, p, q) &= 2y + p = 0, \\ \nabla_z L(x, y, z, p, q) &= 2z + q = 0. \end{aligned}$$

The unique solution, which is a little more difficult to find here, is

$$x = \frac{4}{3}, y = z = \frac{2}{3}, p = q = -\frac{4}{3}.$$

The level surface of the objective function  $f(x, y, z)$  is a sphere, and the constraints are no longer tangent to the level surface at the solution (see Figure 2.10).



**Figure 2.10:** Objective function (sphere) and two constraints (planes)

The normal to the surface is however a linear combination of the normals to the two constraint surfaces, with coefficients  $p$  and  $q$ .

This is an illustration of another version of the **Lagrange Multiplier Theorem**, which says that *if a solution exists, it is a saddle point of the Lagrangian such that no change in the original variables can decrease the Lagrangian, while no change in the Lagrange multipliers can increase it.*

What about inequality constraints?



**Example 31.**

$$\begin{aligned} \min \quad & f(x, y) = x^2 + y^2 \\ \text{such that} \quad & x + y - 2 \geq 0. \end{aligned}$$

Here,  $L(x, y, p) = x^2 + y^2 + p(x + y - 2)$ , but now  $p \leq 0$ . Unfortunately when there are inequality constraints,  $\nabla L = \mathbf{0}$  is neither necessary or sufficient to guarantee a solution to a constrained problem. We need the following theorem.

**Theorem 2.3** (Lagrange multipliers). *Solving*

$$\begin{aligned} \min \quad & f(\mathbf{x}) \\ \text{such that} \quad & A\mathbf{x} + \mathbf{b} \leq \mathbf{0} \end{aligned}$$

*is equivalent to solving*

$$\min_{\mathbf{x}} \left( \max_{\mathbf{p} \in \mathbb{R}_+^n} L(\mathbf{x}, \mathbf{p}) \right),$$

*where **the Lagrangian** is defined as*

$$L(\mathbf{x}, \mathbf{p}) = f(\mathbf{x}) + \mathbf{p}^\top (A\mathbf{x} + \mathbf{b}),$$

*for a vector of **Lagrange multipliers**  $\mathbf{p}$ .*

This theorem uses properties of convex cones and duality to transform the original problem to a problem that uses only the very simple cone  $\mathbb{R}_+^n$ .

**Example 32.** Consider the following linear program.

$$\begin{aligned} \min \quad & f(\mathbf{x}) = \mathbf{c}^\top \mathbf{x} \\ \text{such that} \quad & \left. \begin{aligned} A\mathbf{x} + \mathbf{b} &\leq \mathbf{0} \\ \mathbf{x} &\geq \mathbf{0} \end{aligned} \right\} \end{aligned} \tag{2.34}$$

with

$$\begin{aligned} L(\mathbf{x}, \mathbf{p}) &= \mathbf{c}^\top \mathbf{x} + \mathbf{p}^\top (A\mathbf{x} + \mathbf{b}) \\ &= (\mathbf{c}^\top + \mathbf{p}^\top A) \mathbf{x} + \mathbf{p}^\top \mathbf{b}. \end{aligned} \quad (2.35)$$

Note that (2.35) is the Lagrangian for the dual LP of (2.34). and the Lagrange multipliers are in fact dual variables.

That this is true becomes useful as many of the algorithms for solution of Dual Decomposition problems are written in terms of Lagrange multipliers. These may be simply found in MATLAB as the solution of the Dual of the associated LPs.

## §2.7 DUAL DECOMPOSITION METHOD

We now return to the DEP given by (2.17):

$$\begin{aligned} \min \quad & w = \mathbf{c}^\top \mathbf{x} + \mathbf{E}_\zeta[Q(\mathbf{x}, \zeta)] \\ \text{such that} \quad & A\mathbf{x} = \mathbf{b} \\ & \mathbf{x} \geq \mathbf{0}, \end{aligned}$$

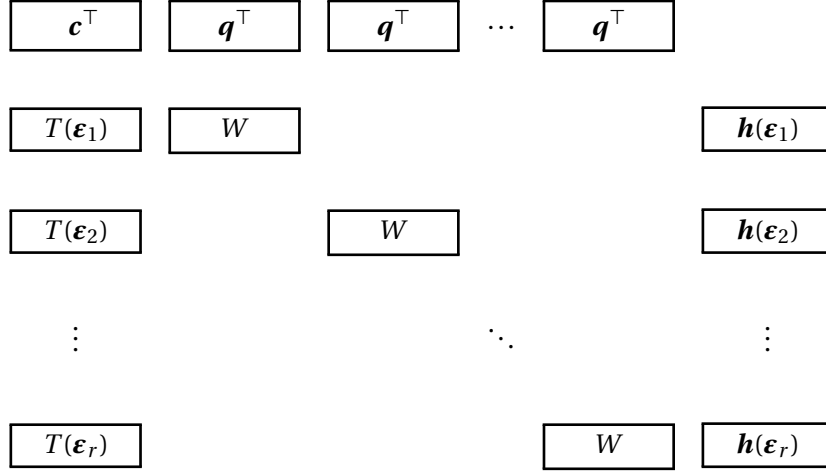
where for each realisation  $\epsilon$  of  $\zeta$ ,

$$Q(\mathbf{x}, \epsilon) = \min_{\mathbf{y}} \left\{ \mathbf{q}^\top \mathbf{y} \mid W\mathbf{y} = \mathbf{h}(\epsilon) - T(\epsilon)\mathbf{x}, \mathbf{y} \geq \mathbf{0} \right\}$$

or, equivalently,

$$\begin{aligned} \min \quad & \mathbf{q}^\top \mathbf{y} \\ \text{such that} \quad & W\mathbf{y} = \mathbf{h}(\epsilon) - T(\epsilon)\mathbf{x} \\ & \mathbf{y} \geq \mathbf{0}. \end{aligned}$$

This DEP has the so-called **Dual Decomposition Structure**:

**Figure 2.11:** Dual Decomposition Structure

We will now consider a solution method that exploits this structure. For simplicity and clarity of exposition, we restrict our attention to the case where  $S_\zeta$  contains just one realisation (the random vector  $\boldsymbol{\zeta}$  has just one realisation); hence, the problem reduces to the LP:

$$\left. \begin{array}{l} \min_{\mathbf{x} \in \mathcal{X}} \quad z = \mathbf{c}^\top \mathbf{x} + \mathbf{q}^\top \mathbf{y} \\ \text{such that} \quad \mathbf{A}\mathbf{x} = \mathbf{b}, \\ \quad \quad \quad T\mathbf{x} + W\mathbf{y} = \mathbf{h}, \\ \quad \quad \quad \mathbf{x}, \mathbf{y} \geq \mathbf{0}. \end{array} \right\} \quad (2.36)$$

Furthermore, we will assume that the problem is *solvable* and that the first-stage feasible set  $\{\mathbf{x} : \mathbf{A}\mathbf{x} = \mathbf{b}, \mathbf{x} \geq \mathbf{0}\}$  is bounded. Then, the LP (2.36) may

be restated as a non-linear program:

$$\left. \begin{array}{ll} \min & z = \mathbf{c}^\top \mathbf{x} + Q(\mathbf{x}) \\ \text{such that} & A\mathbf{x} = \mathbf{b} \\ & \mathbf{x} \geq \mathbf{0}, \end{array} \right\} \quad (2.37)$$

with

$$Q(\mathbf{x}) := \min_{\mathbf{y}} \{ \mathbf{q}^\top \mathbf{y} : W\mathbf{y} = \mathbf{h} - T\mathbf{x}, \mathbf{y} \geq \mathbf{0} \}.$$

The recourse function  $Q(\mathbf{x})$  is piecewise linear and convex. The problem (2.37) is also equivalent to the problem

$$\left. \begin{array}{ll} \min & z = \mathbf{c}^\top \mathbf{x} + \theta \\ \text{such that} & A\mathbf{x} = \mathbf{b} \\ & \theta - Q(\mathbf{x}) \geq 0 \\ & \mathbf{x} \geq \mathbf{0}, \end{array} \right\} \quad (2.38)$$

where the minimisation is over  $\mathbf{x}$  and  $\theta \in \mathbb{R}$ . However, (2.38) requires knowledge of the recourse function  $Q(\mathbf{x})$  in advance. Such knowledge will not be common in practice; that is, an explicit form will typically not be available.

Therefore, we might try to construct a sequence of additional linear constraints that can be used to define a monotonically decreasing feasible set  $\mathcal{B}_1$  of  $(n+1)$ -vectors,  $(x_1, x_2, \dots, x_n, \theta)^\top$ , such that eventually, with

$$\mathcal{B}_0 := \{(\mathbf{x}^\top, \theta)^\top : A\mathbf{x} = \mathbf{b}, \mathbf{x} \geq \mathbf{0}, \theta \in \mathbb{R}\}$$

the problem

$$\min_{(\mathbf{x}, \theta) \in \mathcal{B}_0 \cap \mathcal{B}_1} \mathbf{c}^\top \mathbf{x} + \theta$$

yields a first-stage solution  $\mathbf{x}$  of (2.36).

### Dual Decomposition Method

**Step 1.** Let  $\theta_0$  be a lower bound<sup>1</sup> for

$$\min_y \{ \mathbf{q}^\top \mathbf{y} : A\mathbf{x} = \mathbf{b}, T\mathbf{x} + W\mathbf{y} = \mathbf{h}, \mathbf{x} \geq \mathbf{0}, \mathbf{y} \geq \mathbf{0} \}.$$

Solve the linear program

$$\min_{\mathbf{x}, \theta} \{ \mathbf{c}^\top \mathbf{x} + \theta : A\mathbf{x} = \mathbf{b}, \theta \geq \theta_0, \mathbf{x} \geq \mathbf{0} \},$$

and denote its solution<sup>2</sup>  $(\hat{\mathbf{x}}, \hat{\theta})$ . Let  $\mathcal{B}_1 := \{\mathbb{R}^n \times \{\theta\} : \theta \geq \theta_0\}$ .

**Step 2.** Using the last first-stage solution  $\hat{\mathbf{x}}$ , evaluate the recourse function (i.e. solve the following LP):

$$\begin{aligned} Q(\hat{\mathbf{x}}) &= \min \{ \mathbf{q}^\top \mathbf{y} : W\mathbf{y} = \mathbf{h} - T\hat{\mathbf{x}}, \mathbf{y} \geq \mathbf{0} \} \\ &= \max \{ (\mathbf{h} - T\hat{\mathbf{x}})^\top \mathbf{u} : W^\top \mathbf{u} \leq \mathbf{q}, \mathbf{u} \text{ free} \}. \end{aligned}$$

where the second equality follows from *duality*.

Two cases may arise: (a)  $Q(\hat{\mathbf{x}})$  is unbounded, or (b)  $Q(\hat{\mathbf{x}})$  is bounded.

- a) If  $Q(\hat{\mathbf{x}})$  is unbounded, then  $\hat{\mathbf{x}}$  is not feasible with respect to all the constraints of (2.36) (i.e.  $\hat{\mathbf{x}}$  does not satisfy the induced constraints, as discussed in Theorem 2.1), for any values of  $\mathbf{q} \geq \mathbf{0}$ .

Then, we have a  $\tilde{\mathbf{u}}$  such that

$$\begin{aligned} W^\top \tilde{\mathbf{u}} &\leq \mathbf{0} \quad (\text{suppose we choose } \mathbf{q} = \mathbf{0}), \\ (\mathbf{h} - T\hat{\mathbf{x}})^\top \tilde{\mathbf{u}} &> 0. \end{aligned}$$

<sup>1</sup>How do we know that  $\theta_0$  exists? Because of the assumption that the original problem (2.36) is solvable; that is, there exists at least one solution to the first-stage problem,  $\mathbf{x}$ , such that  $Q(\mathbf{x})$  has a solution.

<sup>2</sup>The existence of this solution is due to the second assumption: that the feasible region of the first-stage problem,  $\{\mathbf{x} : A\mathbf{x} = \mathbf{b}, \mathbf{x} \geq \mathbf{0}\}$ , is bounded.

However, for any  $\mathbf{x}$  that is feasible to the original SLP (2.36), there must exist a  $\mathbf{y} \geq \mathbf{0}$  such that  $W\mathbf{y} = \mathbf{h} - T\mathbf{x}$ , and therefore scalar multiplication of this equation by  $\tilde{\mathbf{u}}$  yields

$$\tilde{\mathbf{u}}^\top (\mathbf{h} - T\mathbf{x}) = \underbrace{\tilde{\mathbf{u}}^\top W}_{\leq 0} \underbrace{\mathbf{y}}_{\geq 0} \leq 0, \quad (2.39)$$

which in turn implies that  $\tilde{\mathbf{u}}^\top \mathbf{h} \leq \tilde{\mathbf{u}}^\top T\mathbf{x}$  must also hold for any feasible  $\mathbf{x}$ . (Inequality (2.39) obviously does not hold for  $\hat{\mathbf{x}}$ , since, as noted above,  $\tilde{\mathbf{u}}^\top (\mathbf{h} - T\hat{\mathbf{x}}) > 0$ .) Therefore, we introduce the **feasibility cut**, cutting off the infeasible solution  $\hat{\mathbf{x}}$ :

$$\tilde{\mathbf{u}}^\top (\mathbf{h} - T\mathbf{x}) \leq 0.$$

Then, redefine

$$\mathcal{B}_1 := \mathcal{B}_1 \cap \{(\mathbf{x}, \theta) : \tilde{\mathbf{u}}^\top (\mathbf{h} - T\mathbf{x}) \leq 0\}$$

and proceed to **Step 3**.

- b) Otherwise, if  $Q(\hat{\mathbf{x}})$  is bounded, we have for the recourse problem, simultaneously, a primal optimal solution  $\hat{\mathbf{y}}$  and a dual optimal solution  $\hat{\mathbf{u}}$  (with  $W^\top \hat{\mathbf{u}} \leq \mathbf{q}$ ). From the dual formulation of the recourse problem, we have

$$Q(\hat{\mathbf{x}}) = (\mathbf{h} - T\hat{\mathbf{x}})^\top \hat{\mathbf{u}},$$

whereas for any  $\mathbf{x}$  we have

$$\begin{aligned} Q(\mathbf{x}) &= \sup_{\mathbf{u}} \{(\mathbf{h} - T\mathbf{x})^\top \mathbf{u} : W^\top \mathbf{u} \leq \mathbf{q}\} \\ &\geq (\mathbf{h} - T\mathbf{x})^\top \hat{\mathbf{u}} \\ &= \hat{\mathbf{u}}^\top (\mathbf{h} - T\mathbf{x}). \end{aligned}$$

The intended constraint  $\theta \geq Q(\mathbf{x})$  in (2.38) implies the linear constraint

$$\theta \geq \hat{\mathbf{u}}^\top (\mathbf{h} - T\mathbf{x}),$$

which is violated by  $(\hat{\mathbf{x}}, \hat{\theta})$  if and only if  $(\mathbf{h} - T\hat{\mathbf{x}})^\top \hat{\mathbf{u}} > \hat{\theta}$ ; in this case we introduce the **optimality cut**, cutting off the non-optimal solution  $(\hat{\mathbf{x}}, \hat{\theta})$ :

$$\theta \geq \hat{\mathbf{u}}^\top (\mathbf{h} - T\mathbf{x}).$$

Correspondingly, we redefine

$$\mathcal{B}_1 := \mathcal{B}_0 \cap \{(\mathbf{x}, \theta) : \theta \geq \hat{\mathbf{u}}^\top (\mathbf{h} - T\mathbf{x})\},$$

and continue to **Step 3**. Otherwise, if  $Q(\hat{\mathbf{x}}) \leq \hat{\theta}$ , we stop.

**Step 3.** Solve the updated (linear programming) problem

$$\min \{ \mathbf{c}^\top \mathbf{x} + \theta : (\mathbf{x}^\top, \theta) \in \mathcal{B}_0 \cap \mathcal{B}_1 \},$$

yielding the optimal solution  $(\tilde{\mathbf{x}}, \tilde{\theta})$ .

With  $(\hat{\mathbf{x}}, \hat{\theta}) := (\tilde{\mathbf{x}}, \tilde{\theta})$ , return to **Step 2**.

This has dealt with the scenario where there is only one set of constraints, because of the assumption that the random vector  $\boldsymbol{\zeta}$  has only one realisation. In summary, in the Dual Decomposition Method, we first look at the induced constraints problem and make appropriate feasibility cuts to the original first stage LP. Once we have a feasible second stage solution, we optimise that solution, which can also require the inclusion of further constraints in the first stage LP. This process is iterated until (assuming solvability) the optimal solution is determined.

When  $\boldsymbol{\zeta}$  has more than one realisation, we develop some pseudo-code, which naturally follows what we have just seen and gives us a handle on these types of problem. The computational burden with respect to the DEP is the solution of all the second-stage recourse linear programs, which essentially depends on the number of realisations of  $\boldsymbol{\zeta}$ . When this number is a finite value  $r$ , with probabilities  $p_k$  for  $1 \leq k \leq r$ , the computational burden can be significantly reduced by associating one set of second-stage decisions  $\mathbf{y}_k$  to each realisation  $\boldsymbol{\varepsilon}_k$ .

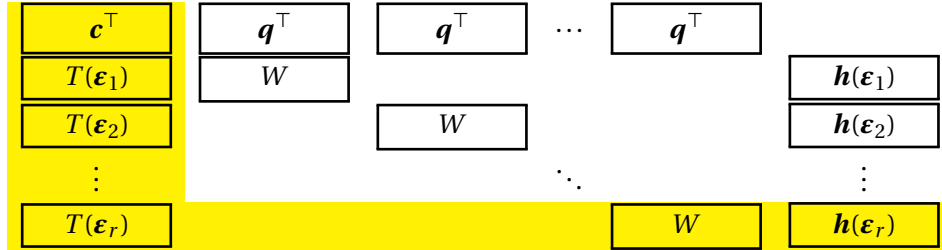
In other words, consider the expanded form of (2.17) given by

$$\begin{aligned} \min_{\mathbf{x} \in \mathcal{X}} \quad & w = \mathbf{c}^\top \mathbf{x} + \sum_{k=1}^r p_k \mathbf{q}_k^\top \mathbf{y}_k \\ \text{such that} \quad & \mathbf{A}\mathbf{x} = \mathbf{b} \\ & T_k \mathbf{x} + W \mathbf{y}_k = \mathbf{h}_k \text{ for } 1 \leq k \leq r \\ & \mathbf{x}, \mathbf{y}_k \geq \mathbf{0} \text{ for } 1 \leq k \leq r. \end{aligned}$$

For each realisation  $\boldsymbol{\varepsilon}_k$  of  $\boldsymbol{\zeta}$ ,  $1 \leq k \leq r$ , we have an extra row in the Decomposition Structure.

**Proposition 2.4.** *If the program (2.36) is solvable and  $\{\mathbf{x} : \mathbf{A}\mathbf{x} = \mathbf{b}, \mathbf{x} \geq \mathbf{0}\}$  is bounded, then the Dual Decomposition Method converges to an optimal solution after finitely many steps.*

We will now introduce the **L-shaped algorithm** for solving stochastic linear programs with the dual decomposition structure. The **L-shaped algorithm** (name derivative of the below structure) takes advantage of this structure, making feasibility and optimality cuts (essentially adding constraints) during an iterative process.



**Figure 2.12:** Dual Decomposition Structure



Initially, we will state the algorithm in pseudo-code and then demonstrate by way of example its processes of creating feasibility and optimality cuts.

### The L-shaped Algorithm

**Step 0. Initialisation.** Set

$$t = s = v := 0, \quad D_0 := \mathbf{0}, \quad d_0 := 0, \quad E_0 := \mathbf{0}, \quad e_0 := 0.$$

Here,  $v$  is the counter for iterations,  $t$  for the feasibility cuts, and  $s$  for optimality cuts.

**Step 1.** Set  $v = v + 1$  and solve the LP

$$\min \quad z = \mathbf{c}^\top \mathbf{x} + \theta \quad (2.40)$$

$$\text{such that} \quad A\mathbf{x} = \mathbf{b} \quad (2.41)$$

$$D_\ell \mathbf{x} \geq d_\ell \quad \text{for } \ell = 1, \dots, t \quad (2.42)$$

$$E_\ell \mathbf{x} + \theta \geq e_\ell \quad \text{for } \ell = 1, \dots, s \quad (2.43)$$

$$\mathbf{x} \geq \mathbf{0}, \quad \theta \in \mathbb{R}. \quad (2.44)$$

If no constraints (2.43) exist, we set  $\theta = 0$ , to solve the LP. The optimal solution at iteration  $v$  is denoted  $(\mathbf{x}^{(v)}, \theta^{(v)})$ , where we set  $\theta^{(v)} = -\infty$  if no constraints (2.43) exist.

Note that the constraints (2.42) correspond to feasibility cuts, and the constraints (2.43) correspond to optimality cuts.

**Step 2. Feasibility cuts.** For  $k = 1, \dots, r$  (consider all realisations of  $\xi$ ) and while  $\tilde{w} > 0$ , solve the LP

$$\min \quad \tilde{w} = \mathbf{e}^\top \mathbf{v}^+ + \mathbf{e}^\top \mathbf{v}^- \quad (2.45)$$

$$\text{such that} \quad W\mathbf{y}_k^{(v)} + \mathbf{v}^+ - \mathbf{v}^- = \mathbf{h}_k - T_k \mathbf{x}^{(v)} \quad (2.46)$$

$$\mathbf{y}_k^{(v)} \geq \mathbf{0}, \quad \mathbf{v}^+, \mathbf{v}^- \geq \mathbf{0}, \quad (2.47)$$

where  $\mathbf{e}^\top = (1, \dots, 1)$  is a vector of ones of the appropriate dimension, which is used to sum the vectors  $\mathbf{v}^+$  and  $\mathbf{v}^-$ . (The inequality  $\tilde{w} > 0$  means that there is no feasible  $\mathbf{y} \geq \mathbf{0}$ .)

If at any stage there is a value  $k$  with  $\tilde{w} > 0$ , let  $\boldsymbol{\sigma}_k^{(v)}$  be the associated vector of Lagrange multipliers (which we obtain by solving the dual of the above LP, (2.45)–(2.47)), and define

$$\mathbf{D}_{t+1} = (\boldsymbol{\sigma}_k^{(v)})^\top T_k, \quad (2.48)$$

$$d_{t+1} = (\boldsymbol{\sigma}_k^{(v)})^\top \mathbf{h}_k, \quad (2.49)$$

set  $t = t + 1$ , go back to **Step 1**.

Otherwise,  $\tilde{w} = 0$  for all  $k \in \{1, 2, \dots, r\}$ , which means that given  $\mathbf{x}^{(v)}$ , for every realisation of  $\boldsymbol{\xi}$  the second-stage program has a feasible solution  $\mathbf{y}_k^{(v)}$ .

This in turn implies we can proceed to **Step 3**, which considers the optimality of the solution since everything will now be feasible.

**Step 3. Optimality cuts.** For  $k = 1, \dots, r$  (consider all realisations of  $\boldsymbol{\xi}$ ), solve the LP

$$\min \quad \hat{w} = \mathbf{q}^\top \mathbf{y}_k^{(v)} \quad (2.50)$$

$$\text{such that} \quad W \mathbf{y}_k^{(v)} = \mathbf{h}_k - T_k \mathbf{x}^{(v)} \quad (2.51)$$

$$\mathbf{y}_k^{(v)} \geq \mathbf{0}. \quad (2.52)$$

Let  $\boldsymbol{\pi}_k^{(v)}$  be the vector of Lagrange multipliers associated with an optimal solution for a given  $k$ , and define

$$\mathbf{E}_{s+1} = \sum_{k=1}^r p_k (\boldsymbol{\pi}_k^{(v)})^\top T_k, \quad (2.53)$$

$$e_{s+1} = \sum_{k=1}^r p_k (\boldsymbol{\pi}_k^{(v)})^\top \mathbf{h}_k, \quad (2.54)$$

and set

$$J^{(v)} = e_{s+1} - \mathbf{E}_{s+1} \mathbf{x}^{(v)}.$$

If  $\theta^{(v)} \geq J^{(v)}$ , stop the iteration procedure, as  $\mathbf{x}^{(v)}$  is optimal. Otherwise, set  $s = s + 1$  and go to **Step 1**.

**Example 33** (Optimality cuts).

$$\left. \begin{array}{ll} \min z = & 100x_1 + 150x_2 + \mathbb{E}_{\xi}(q_1y_1 + q_2y_2) \\ \text{such that} & x_1 + x_2 \leq 120 \\ & 6y_1 + 10y_2 \leq 60x_1 \\ & 8y_1 + 5y_2 \leq 80x_2 \\ & y_1 \leq d_1 \\ & y_2 \leq d_2 \\ & x_1 \geq 40 \\ & x_2 \geq 20 \end{array} \right\} \quad (2.55)$$

where the random vector  $\xi = (d_1, d_2, q_1, q_2)^\top$  has two vector values

$$\boldsymbol{\varepsilon}_1 = (500, 100, -24, -28)^\top \quad \text{with probability } 0.4,$$

$$\boldsymbol{\varepsilon}_2 = (300, 300, -28, -32)^\top \quad \text{with probability } 0.6.$$

We first observe that the second stage is feasible for every  $\mathbf{x}$  that is feasible to the first-stage problem, because  $\mathbf{x} > \mathbf{0}$  and  $\mathbf{d} = (d_1, d_2)^\top > \mathbf{0}$ . (In other words,  $W$  is a relatively complete recourse matrix.) We write the first stage LP as

$$\begin{array}{ll} \min z = & 100x_1 + 150x_2 + \theta \\ \text{s.t.} & x_1 + x_2 \leq 120 \\ & x_1 \geq 40 \\ & x_2 \geq 20, \end{array}$$

where initially we set  $\theta = 0$ .

The second stage consists of two LPs based on the two realisations of  $\xi$  and the constraint

$$W\mathbf{y} = \mathbf{h}(\boldsymbol{\varepsilon}_i) - T(\boldsymbol{\varepsilon}_i)\mathbf{x},$$

where

$$W = \begin{pmatrix} 6 & 10 \\ 8 & 5 \\ 1 & 0 \\ 0 & 1 \end{pmatrix}, \quad \mathbf{h}(\boldsymbol{\varepsilon}_1) = \begin{pmatrix} 0 \\ 0 \\ 500 \\ 100 \end{pmatrix}, \quad \mathbf{h}(\boldsymbol{\varepsilon}_2) = \begin{pmatrix} 0 \\ 0 \\ 300 \\ 300 \end{pmatrix}$$

$$T(\boldsymbol{\varepsilon}_1) = T(\boldsymbol{\varepsilon}_2) = \begin{pmatrix} -60 & 0 \\ 0 & -80 \\ 0 & 0 \\ 0 & 0 \end{pmatrix}, \quad \mathbf{q}(\boldsymbol{\varepsilon}_1) = \begin{pmatrix} -24 \\ -28 \end{pmatrix}, \quad \mathbf{q}(\boldsymbol{\varepsilon}_2) = \begin{pmatrix} -28 \\ -32 \end{pmatrix}.$$

**Step 0.** We initialise by setting

$$t = s = v = 0, \quad \mathbf{D}_0 = \mathbf{0}, \quad d_0 = 0, \quad \mathbf{E}_0 = \mathbf{0}, \quad e_0 = 0.$$

**First Iteration:**

**Step 1.** We set  $v = v + 1 = 1$  and solve

$$\begin{aligned} \min \quad & z = 100x_1 + 150x_2 \\ \text{s.t.} \quad & x_1 + x_2 \leq 120 \\ & x_1 \geq 40 \\ & x_2 \geq 20 \end{aligned}$$

to get  $\mathbf{x}^{(1)} = (40, 20)^\top$ , with  $z^{(1)} = 7000$ . As there are no constraints (2.43), we set  $\theta^{(1)} = -\infty$ .

**Step 2.** We bypass this, as the second-stage is always feasible for every  $\mathbf{x}$  that is feasible for the first-stage problem.

**Step 3.** For  $\boldsymbol{\varepsilon}_1 = (500, 100, -24, -28)^\top$ , we solve

$$\begin{aligned} \min \quad & w = -24y_1 - 28y_2 \\ \text{s.t.} \quad & 6y_1 + 10y_2 \leq 60x_1 = 2400 \\ & 8y_1 + 5y_2 \leq 80x_2 = 1600 \\ & y_1 \leq 500 \\ & y_2 \leq 100 \\ & y_1, y_2 \geq 0, \end{aligned}$$

to obtain  $\mathbf{y}_1^{(1)} = (137.5, 100)^\top$  with  $w_1^{(1)} = -6100$ . The dual LP here yields solution

$$\boldsymbol{\pi}_1^{(1)} = (0, -3, 0, -13)^\top.$$

(Note that the dimensions, bounds and signs for the dual LP variables  $\boldsymbol{\pi}$  depend on the form of the primal LP)

For  $\boldsymbol{\epsilon}_2 = (300, 300, -28, -32)^\top$ , we solve

$$\begin{array}{llllll} \min & w & = & -28y_1 & - & 32y_2 \\ \text{s.t.} & & & 6y_1 & + & 10y_2 & \leq & 60x_1 = 2400 \\ & & & 8y_1 & + & 5y_2 & \leq & 80x_2 = 1600 \\ & & & & & y_1 & \leq & 300 \\ & & & & & y_2 & \leq & 300 \\ & & & & & \textcolor{red}{y_1, y_2} & \geq & \textcolor{red}{0}, \end{array}$$

to get  $\mathbf{y}_2^{(1)} = (80, 192)^\top$  with  $w_2^{(1)} = -8384$ . The dual LP has the solution

$$\boldsymbol{\pi}_2^{(1)} = (-2.32, -1.76, 0, 0)^\top.$$

Now, using  $\mathbf{h}(\boldsymbol{\epsilon}_1)$ ,  $\mathbf{h}(\boldsymbol{\epsilon}_2)$ ,  $T(\boldsymbol{\epsilon}_1)$  and  $T(\boldsymbol{\epsilon}_2)$ , we get

$$\begin{aligned} e_1 &= 0.4\boldsymbol{\pi}_1^{(1)} \mathbf{h}(\boldsymbol{\epsilon}_1) + 0.6\boldsymbol{\pi}_2^{(1)} \mathbf{h}(\boldsymbol{\epsilon}_2) = -520, \\ \mathbf{E}_1 &= 0.4\boldsymbol{\pi}_1^{(1)} T(\boldsymbol{\epsilon}_1) + 0.6\boldsymbol{\pi}_2^{(1)} T(\boldsymbol{\epsilon}_2) = (83.52, 180.48)^\top, \\ J^{(1)} &= e_1 - \mathbf{E}_1 \mathbf{x}^{(1)} = -520 - (83.52, 180.48) \mathbf{x}^{(1)} = -7470.4. \end{aligned}$$

As  $J^{(1)} = -7470.4 > \theta^{(1)} = -\infty$ , we add the optimality cut

$$\textcolor{blue}{83.52x_1 + 180.48x_2 + \theta \geq -520}$$

to the original first stage LP and perform the second iteration.

**Second Iteration:**

**Step 1:** We set  $v = v + 1 = 2$  and solve

$$\begin{array}{llllll} \min & z & = & 100x_1 & + & 150x_2 & + & \theta \\ \text{s.t.} & & & x_1 & + & x_2 & & \leq 120 \\ & & & 83.52x_1 & + & 180.48x_2 & + & \theta \geq -520 \\ & & & & & x_1 & \geq & 40 \\ & & & & & x_2 & \geq & 20, \end{array}$$

to get  $\mathbf{x}^{(2)} = (40, 80)^\top$ ,  $\theta^{(2)} = -18299.2$ , with  $z^{(2)} = -2299.2$ .

**Step 2:** Again, since we will always have a feasible  $\mathbf{y}$  for all first stage solutions  $\mathbf{x}$ , we skip this step.

**Step 3:**

$$\begin{aligned} \mathbf{y}_1^{(2)} &= (400, 0)^\top, w_1^{(2)} = -9600, \text{ with } \boldsymbol{\pi}_1^{(2)} = (-4, 0, 0, 0)^\top \\ \mathbf{y}_2^{(2)} &= (300, 60)^\top, w_2^{(2)} = -10320 \text{ with } \boldsymbol{\pi}_2^{(2)} = (-3.2, 0, -8.8, 0)^\top \\ e_2 &= 0.4\boldsymbol{\pi}_1^{(2)}\mathbf{h}(\boldsymbol{\epsilon}_1) + 0.6\boldsymbol{\pi}_2^{(2)}\mathbf{h}(\boldsymbol{\epsilon}_2) = -1584, \\ \mathbf{E}_2 &= 0.4\boldsymbol{\pi}_1^{(2)}T(\boldsymbol{\epsilon}_1) + 0.6\boldsymbol{\pi}_2^{(2)}T(\boldsymbol{\epsilon}_2) = (211.2, 0)^\top, \\ J^{(2)} &= e_2 - \mathbf{E}_2\mathbf{x}^{(2)} = -10032 > \theta^{(2)} = -18299.2. \end{aligned}$$

Hence, we add the following optimality cut to the first stage LP

$$211.2x_1 + 0x_2 + \theta \geq -1584.$$

**Third Iteration:**

**Step 1:** We set  $v = v + 1 = 3$  and solve

$$\begin{array}{llllll} \min & z & = & 100x_1 & + & 150x_2 & + & \theta \\ \text{s.t.} & & & x_1 & + & x_2 & & \leq 120 \\ & & & 83.52x_1 & + & 180.48x_2 & + & \theta \geq -520 \\ & & & 211.2x_1 & & & + & \theta \geq -1584 \\ & & & & & x_1 & \geq & 40 \\ & & & & & x_2 & \geq & 20 \end{array}$$

to get  $\mathbf{x}^{(3)} = (66.828, 53.172)^\top$ ,  $\theta^{(3)} = -15697.993$ , with  $z^{(3)} = -1039.375$ .

**Step 2:** Skipping this step, again.

**Step 3.**

$$\begin{aligned} \mathbf{y}_1^{(3)} &= (469.224, 100)^\top, w_1^{(3)} = -14061.37, \quad \text{with } \boldsymbol{\pi}_1^{(3)} = (0, -3, 0, -13)^\top, \\ \mathbf{y}_2^{(3)} &= (300, 220.966)^\top, w_2^{(3)} = -15470.903, \quad \text{with } \boldsymbol{\pi}_2^{(3)} = (-3.2, 0, -9.8, 0)^\top, \\ e_3 &= 0.4\boldsymbol{\pi}_1^{(3)} \mathbf{h}(\boldsymbol{\varepsilon}_1) + 0.6\boldsymbol{\pi}_2^{(3)} \mathbf{h}(\boldsymbol{\varepsilon}_2) = -2104, \\ \mathbf{E}_3 &= 0.4\boldsymbol{\pi}_1^{(3)} T(\boldsymbol{\varepsilon}_1) + 0.6\boldsymbol{\pi}_2^{(3)} T(\boldsymbol{\varepsilon}_2) = (115.2, 96)^\top, \\ J^{(3)} &= e_3 - \mathbf{E}_3 \mathbf{x}^{(3)} = -10032 > \theta^{(3)} = -14907.09. \end{aligned}$$

Hence, we add the cut

$$115.2x_1 + 96x_2 + \theta \geq -2104.$$

**Fourth Iteration:**

**Step 1.** We set  $v = v + 1 = 4$  and solve

$$\begin{array}{rclclcl} \min \quad z & = & 100x_1 & + & 150x_2 & + & \theta \\ \text{s.t.} & & x_1 & + & x_2 & & \leq 120 \\ & & 83.52x_1 & + & 180.48x_2 & + & \theta \geq -520 \\ & & 211.2x_1 & & & + & \theta \geq -1584 \\ & & 115.2x_1 & + & 96x_2 & + & \theta \geq -2104 \\ & & & & x_1 & \geq & 40 \\ & & & & x_2 & \geq & 20 \end{array}$$

to get  $\mathbf{x}^{(4)} = (40, 33.75)^\top$ ,  $\theta^{(4)} = -9952$ , with  $z^{(4)} = -889.5$ .

**Step 2:** We are skipping this step, again!

**Step 3.**

$$\begin{aligned}
\mathbf{y}_1^{(4)} &= (300, 60)^\top, \quad w_1^{(4)} = -8880, \quad \text{with } \boldsymbol{\pi}_1^{(4)} = (-2.08, -1.44, 0, 0)^\top, \\
\mathbf{y}_2^{(4)} &= (300, 60)^\top, \quad w_2^{(4)} = -10320 \quad \text{with } \boldsymbol{\pi}_2^{(4)} = (-3.175, -0.05, -8.548, 0)^\top, \\
e_4 &= 0.4\boldsymbol{\pi}_1^{(4)} \mathbf{h}(\boldsymbol{\epsilon}_1) + 0.6\boldsymbol{\pi}_2^{(4)} \mathbf{h}(\boldsymbol{\epsilon}_2) = -1538.601, \\
\mathbf{E}_4 &= 0.4\boldsymbol{\pi}_1^{(4)} T(\boldsymbol{\epsilon}_1) + 0.6\boldsymbol{\pi}_2^{(4)} T(\boldsymbol{\epsilon}_2) = (164.212, 48.501)^\top, \\
J^{(4)} &= e_4 - \mathbf{E}_4 \mathbf{x}^{(4)} = -9744 > \theta^{(4)} = -9952
\end{aligned}$$

Hence, we add the cut

$$164.212x_1 + 48.501x_2 + \theta \geq -1538.601.$$

**Fifth Iteration:**

**Step 1.** We set  $\nu = \nu + 1 = 5$  and solve

$$\begin{array}{llllll}
\min \quad z & = & 100x_1 & + & 150x_2 & + & \theta \\
\text{s.t.} & & x_1 & + & x_2 & & \leq 120 \\
& & 83.52x_1 & + & 180.48x_2 & + & \theta \geq -520 \\
& & 211.2x_1 & & & + & \theta \geq -1584 \\
& & 115.2x_1 & + & 96x_2 & + & \theta \geq -2104 \\
& & 164.212x_1 & + & 48.501x_2 & + & \theta \geq -1538.601 \\
& & & & x_1 & \geq & 40 \\
& & & & x_2 & \geq & 20
\end{array}$$

to get  $\mathbf{x}^{(5)} = (46.667, 36.25)^\top, \theta^{(5)} = -10960$ , with  $z^{(5)} = -855.833$ .



Skipping **Step 2** again, we perform **Step 3**:

$$\begin{aligned}
 \mathbf{y}_1^{(5)} &= (300, 100)^\top, \quad w_1^{(5)} = -10000 \quad \text{with} \\
 \boldsymbol{\pi}_1^{(5)} &= (-0.695, -2.479, 0, -8.655)^\top \\
 \mathbf{y}_2^{(5)} &= (300, 100)^\top, \quad w_2^{(5)} = -11600 \quad \text{with} \\
 \boldsymbol{\pi}_2^{(5)} &= (-3.156, -0.088, -8.362, 0)^\top \\
 e_5 &= 0.4\boldsymbol{\pi}_1^{(5)} \mathbf{h}(\boldsymbol{\varepsilon}_1) + 0.6\boldsymbol{\pi}_2^{(5)} \mathbf{h}(\boldsymbol{\varepsilon}_2) = -1851.361 \\
 \mathbf{E}_5 &= 0.4\boldsymbol{\pi}_1^{(5)} T(\boldsymbol{\varepsilon}_1) + 0.6\boldsymbol{\pi}_2^{(5)} T(\boldsymbol{\varepsilon}_2) = (130.307, 83.522)^\top \\
 (J^{(5)} = e_5 - \mathbf{E}_5 \mathbf{x}^{(5)} = -10960) &= (\theta^{(5)} = -10960)
 \end{aligned}$$

Stop as  $\mathbf{x}^{(5)} = (46.667, 36.25)^\top$ , with  $z^{(5)} = -855.833$  is optimal.

## Chapter 3

# Markov Decision Chains

This section follows **very, very closely** the lecture notes of Richard Weber at the University of Cambridge. Very good reference books are [5] and [6] (also, followed **very** closely in parts).

We are now going to be concerned with *dynamical systems* and their optimisation *over time*. We will focus on stochastic systems, which are therefore dynamic and which are related to the *stochastic linear programs* we have seen (multi-stage recourse problems).

**Markov decision processes** are essentially models for sequential decision making when outcomes are uncertain. As we will see in this course, Markov decision process models consist of *decision epochs*, *states*, *actions*, *rewards/costs* and *transition probabilities*. Choosing an action in a state at a decision epoch generates a reward/cost, and determines the state at the next decision epoch through a transition probability function.

We will encounter the concept of *policies* or *strategies* which prescribe particular actions to choose under any eventuality at every future decision epoch. As decision makers, we seek policies that are optimal in some sense.

People make decisions every day and those decisions have both short term and long term consequences. The same is true for the actions taken within these Markov decision processes and therefore, we should take care

that decisions are not made in isolation, since any decision taken now impacts any subsequent decision, which impacts subsequent decisions...

**Defn 3.1.** *A **sequential decision model** is as follows. At a specific point in time a decision maker observes the state of a system. Based on this state, the decision maker chooses an action. This action results in two outcomes:*

1. *an immediate reward or cost and*
2. *the evolution of the process to a new state at a subsequent point in time according as some probability distribution.*

*At the subsequent point in time the decision maker faces a similar problem, but maybe with the system in a different state and with a different set of possible actions.*

We will assume that we have a random process which can, to some extent, be controlled. For example, at certain decision points, we have a choice of actions that can be taken, where different actions have the potential to affect the future behaviour of the process in different ways.

When actions can be taken to influence the evolution of a Markov process, then we call this a **Markov decision process**, or alternatively, a **Markov decision problem**. Typically, the aim of exerting some influence on the process is to minimise or maximise a reward or cost function, which in turn depends on the state trajectory of the Markov process.

### §3.1 MOTIVATION

**Example 34** (Machine repair I: A cost minimisation problem). Consider the problem of managing the operation of a machine efficiently over  $T$  time

periods, corresponding to  $T + 1$  stages, labelled  $0, 1, \dots, T$ . During each time period, the machine can be in any one of  $N$  states, denoted  $1, 2, \dots, N$ :

- State 1 corresponds to the machine being in perfect condition.

$\vdots$

- State  $N$  corresponds to the machine being in greatest disrepair.

Let  $h(i)$  be the operating cost of the machine during one time period when the machine is in state  $i$ , and assume that

$$h(1) \leq h(2) \leq h(3) \leq \dots \leq h(N).$$

During each of the  $T$  time periods, in the absence of repair, the state of the machine can either remain the same or deteriorate to one of the higher states of disrepair. Let  $p_{ij} = P(\text{next state is } j | \text{current state is } i)$ , with  $p_{ij} = 0$  if  $j < i$ , for  $i, j = 1, \dots, N$ .

At each stage of operation (the start of each time period), we know the state of the machine, and we must choose one of the following possible actions:

1. First repair the machine to perfect working order (state 1) at a cost  $R$  immediately prior to the operating period.
2. Let the machine operate one more time period in its current state.

Assume that once repaired, the machine is guaranteed to stay in state 1 until the next stage (until the beginning of the next time period). In subsequent periods, it may deteriorate to states  $j > 1$  according to the transition probabilities  $p_{1j}$ . At the  $T$ th stage (the end of the  $T$ th time period), the operation of the machine stops, and the cost incurred is zero for all states of the machine.

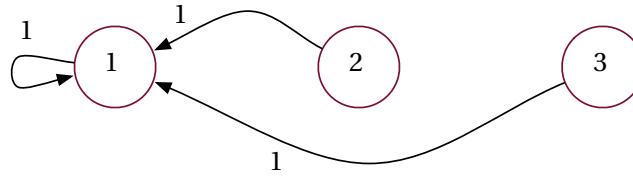
The objective is, *to determine, for each time period and for each state of the machine, whether the optimal action is to repair, or leave the machine in its current state*. In particular, we wish to obtain a set of instructions, called a

*policy*, which prescribes the sequence of actions to take when the machine starts in a given state  $i_0$ , such that the total expected cost of operation over the  $T$  time periods is minimised.

Consider a specific instance of the machine repair problem over  $T = 4$  time periods, and  $N = 3$  states. Let the operating cost in each state and cost of repair be given by  $h(1) = 2$ ,  $h(2) = 4$ ,  $h(3) = 6$ ,  $R = 6$ . The decision stages are  $n = 0, 1, 2, 3, 4$ .

If the action REPAIR is chosen, the state transition probabilities are given by

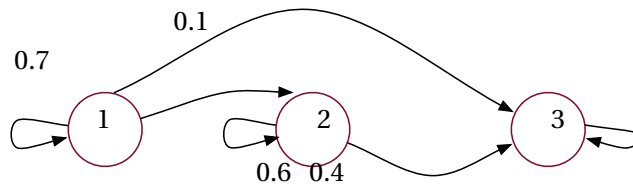
$$P(\text{repair}) = \begin{pmatrix} 1 & 0 & 0 \\ 1 & 0 & 0 \\ 1 & 0 & 0 \end{pmatrix}.$$



**Figure 3.1:** State transition diagram if action REPAIR is chosen.

Alternatively, if the action DON'T REPAIR is chosen, the state transition probabilities are given by

$$P(\text{don't repair}) = \begin{pmatrix} 0.7 & 0.2 & 0.1 \\ 0 & 0.6 & 0.4 \\ 0 & 0 & 1 \end{pmatrix}.$$



**Figure 3.2:** State transition diagram if the action DON'T REPAIR is chosen.

The transition probabilities of the Markov chain at stage (time)  $n$  therefore depend on the action that is chosen at stage  $n$ . Let  $X_n$  represent the state of the machine at stage  $n$ , that is, at the beginning of the  $n$ th operation period. If we were to employ the same action at each decision stage, then the Markov process  $X_n, n = 0, \dots, 4$ , on the states  $\{1, 2, 3\}$ , would be time-homogeneous. On the other hand, if we were to alternate the action at each stage, then the Markov process would be time-inhomogeneous.

This machine repair example is a *finite horizon problem*, because  $T$  is finite. Furthermore,  $T$  is a deterministic quantity which is given in advance.

### §3.2 ELEMENTS OF A MARKOV DECISION CHAIN

We now focus on stochastic systems, in particular where the evolution of the process is governed by a DTMC. A **Markov decision chain** (MDC)  $\{X_n\}_{n \in \mathbb{N}}$  consists of:

1. **States**, in a **state space**  $\mathcal{S} \subseteq \mathbb{N}$ . For example,  $\mathcal{S} = \{0, 1, \dots, N\}$ .  $X_n \in \mathcal{S}$  for all  $n = 0, 1, \dots$ . We will often assume  $\mathcal{S} \subset \mathbb{N}$ .
2. **Actions**,  $u_n \in \mathcal{U}_n(x_n)$ , for all  $n = 0, 1, \dots$ . We will often have  $\mathcal{U}(x_n)$ , where the set of actions available depends upon the current state of the chain only.
3. **Costs**,  $c(x_n, u_n, n)$ . Assumed separable (decomposable) — only incurred at current stage and dependent upon action. We will often assume  $c(x_n, u_n)$ , that is, the cost depends only upon current state and/or action, and not on time.

We also have a set of **terminal costs**  $C(x_T)$ , which are the costs incurred in the final stage  $T$  (if  $T < \infty$ ) and are dependent on the state  $x_T$ . This is often assumed to be 0, but not always.

4. **Transition probability distributions**. Let

$$\mathbf{x}_n = (x_1, x_2, \dots, x_n) \quad \text{and} \quad \mathbf{u}_n = (u_0, u_1, \dots, u_n)$$

denote the state and action histories at time  $n$ , respectively. We must specify

$$\begin{aligned}\Pr(X_{n+1} = x_{n+1} | \mathbf{x}_n, \mathbf{u}_n, n) &=^* \Pr(X_{n+1} = x_{n+1} | x_n, u_n, n) \\ &= p(x_{n+1} | x_n, u_n),\end{aligned}$$

where  $=^*$  is a consequence of Markov property.

#### 5. Optimisation criterion:

- a) **Finite horizon.** In this case, the behaviour of the system is considered for epochs  $n = 0$  to  $n = T$ , for a fixed positive integer  $T$ .
- b) **Infinite horizon.** In this case,  $n = 0, 1, \dots$ . We will consider this situation with and without *discounting*. *Average cost* also rests in this category.

### POLICIES

A **policy** (or **strategy**) is a rule for choosing the action under all possible circumstances as a function of the circumstances. A controller's actions under a policy can be based on the previous states visited, the actions chosen in those states, and when those visits occurred.

There are several important types of policies, these are classified by how much of the history may be utilised by the controller:

- A **stationary policy**, denoted by  $f$ , operates as follows: Associated with each state  $i$  is a distinguished action  $u(i) \in \mathcal{U}_i$ . If at any time the *controller* finds the system in state  $i$ , then the controller always chooses the action  $u(i)$ .

Thus, a stationary policy depends on the history of the process only through the current state. To implement a stationary policy, the controller needs to know the current state of the system only.

- A slightly less restrictive type of policy is the **randomised stationary policy**, denoted by  $f_R$ . Associated with each state  $i$  is a pmf  $P_{f_R}^i$  on  $\mathcal{U}_i$ . If at any time the controller finds the system in state  $i$ , then the controller always chooses action  $u$  with prob.  $P_{f_R}^i(u)$ .

As is the case for a stationary policy, a randomised stationary policy depends on the history of the process only through the current state. A stationary policy can be viewed as a degenerate randomised stationary policy.

- A **deterministic Markov policy** is a sequence  $\theta = (f_0, f_1, f_2, \dots)$  of stationary policies. It operates as follows: If the chain is in state  $i$  at time  $t = n$ , then the controller chooses action  $f_n(i)$ . Thus a deterministic Markov policy depends on the history of the chain only through the current state and the time index.
- A **randomised Markov policy** is a sequence  $\theta = (f_{R_0}, f_{R_1}, f_{R_2}, \dots)$  of randomised stationary policies. It operates as follows: If the chain is in state  $i$  at time  $t = n$ , then the controller chooses action  $u \in \mathcal{U}_i$  with probability  $P_{R_n}^i(u)$ . Thus a randomised Markov policy depends on the history of the chain only through the current state and the time index.



## §3.3 FINITE HORIZON PROBLEMS

**Defn 3.2.** A *basic finite horizon problem* is a discrete-time Markov chain with finite time horizon  $T$ , where  $T$  is given in advance. It takes values on a finite state space  $\mathcal{S} = \{1, \dots, N\}$ .

For each state  $i \in \mathcal{S}$  and each  $n \in \{0, 1, \dots, T-1\}$ , there is a finite set of actions  $\mathcal{U}_n(i)$  that can be selected. No action needs to be selected at the terminating time  $T$ .

Transitions between states  $(i, j)$  are governed by the probabilities  $p_{ij}(u_n)$ , where  $u_n \in \mathcal{U}_n(i)$ . So the state transition probabilities in state  $i$  depend on the particular action  $u_n \in \mathcal{U}_n(i)$  that is selected. We assume that the transition probabilities are a function of  $n$  only through the actions  $u_n$ .

Let  $c(i, u_n)$  be the expected single-stage cost incurred by selecting action  $u_n$  when in state  $i$  at time  $n$ , which is the expected cost incurred during the time period  $n$  to  $n+1$ , when the process is in state  $i$  at time  $n$  and the action  $u_n$  is selected.

If we assume that the expected single stage costs  $c(i, u_n)$  are stationary, then this means that the cost  $c(i, u_n)$  depends on  $n$  only through the action  $u_n$ .

Let  $C(i)$ ,  $i \in \mathcal{S}$ , be the **terminal cost** associated with ending in state  $i$ ; often this is set to zero.

We will consider a class of **policies**  $\Psi$ , each consisting of a sequence of functions

$$\Psi = \{f_0, f_1, \dots, f_{T-1}\},$$

where each function  $f_n$  maps states  $i \in \mathcal{S}$  into actions:

$$u_n = f_n(i).$$

We assume that  $f_n(i) \in \mathcal{U}_n(i)$  for all  $i \in \mathcal{S}$ ; such a policy is called **admissible**. Without this constraint, the policy does not make sense.

Given a policy  $\Psi$ , the state evolution is a Markov process,  $X_n^\Psi$ , on the state space  $\mathcal{S}$ , for  $n = 0, \dots, T$ , with transition probabilities

$$P(X_{n+1}^\Psi = j \mid X_n^\Psi = i) = p_{ij}(f_n(i)) \quad \text{for } i, j \in \mathcal{S} \text{ and } n = 0, \dots, T-1.$$

The expected cost of a policy  $\Psi$ , starting from state  $i$  at the beginning of stage  $k \in \{0, 1, \dots, T\}$ , is defined as

$$J_k^\Psi(i) = \mathbb{E} \left[ C(X_T^\Psi) + \sum_{n=k}^{T-1} c(X_n^\Psi, f_n(X_n^\Psi)) \right], \quad \text{with } X_k^\Psi = i. \quad (3.1)$$

**The objective** in solving the finite horizon MDP is to obtain a policy  $\Psi$  that minimises the expected cost for  $k = 0$

$$J_0^\Psi(i_0) = \mathbb{E} \left[ C(X_T^\Psi) + \sum_{n=0}^{T-1} c(X_n^\Psi, f_n(X_n^\Psi)) \right], \quad (3.2)$$

for a given starting state  $X_0^\Psi = i_0$ . That is, we wish to minimise the expected total cost over all stages (recall that decisions are only made at stages  $0, 1, \dots, T-1$ , and none at stage  $T$ ).

**Defn 3.3.** For a given initial state  $i_0$ , an **optimal policy**

$$\Psi^* = \{f_0^*, f_1^*, \dots, f_{T-1}^*\}$$

is a policy that satisfies

$$J^{\Psi^*}(i_0) = \min_{\Psi \in \Gamma} J_0^\Psi(i_0),$$

where  $\Gamma$  is the set of all admissible policies.

The **optimal cost** depends on the starting state  $i_0$ , and is denoted  $J^*(i_0)$ :

$$J^*(i_0) = J^{\Psi^*}(i_0) = \min_{\Psi \in \Gamma} J_0^\Psi(i_0).$$

We may view  $J^*$  as a function that assigns to each initial state  $i_0 \in \mathcal{S}$  the optimal cost  $J^*(i_0)$ .

Finite horizon problems can be solved using the technique known as **dynamic programming**, which rests on the Principle of Optimality, which is stated as follows.

**Theorem 3.1** (Principle of Optimality). *Let*

$$\Psi^* = \{f_0^*, f_1^*, \dots, f_{T-1}^*\}$$

*be an optimal policy starting from a given state  $i_0$ . Assume that when using  $\Psi^*$ , the system visits state  $i$  at stage  $k$  with positive probability; that is,*

$$P(X_k^{\Psi^*} = i \mid X_0 = i_0) > 0.$$

*Consider the sub-problem in which the process is in state  $i$  at stage  $k$ ,  $X_k^{\Psi^*} = i$ , and we wish to minimise the expected cost incurred over the “time-to-go” or remaining  $T - k$  stages, given by*

$$\mathbb{E} \left[ C(X_T^\Psi) + \sum_{n=k}^{T-1} c(X_n^\Psi, f_n(X_n^\Psi)) \right].$$

*Then, the truncated policy  $\{f_k^*, f_{k+1}^*, \dots, f_{T-1}^*\}$  is optimal for this sub-problem.*

A sketch of the proof is, if the truncated policy were not optimal, then the entire policy  $\Psi^*$  would not be optimal. An analogy: Suppose that the shortest route to travel between Adelaide and Paris is via Frankfurt. It follows that the

part of that route between Frankfurt and Paris is also the shortest route for the “sub-problem” of travelling from Frankfurt to Paris in minimum time.

The Principle of Optimality suggests that an optimal policy can be constructed piece-by-piece, starting from the final stage of the decision process and working backwards, until an optimal policy for the entire problem is constructed.

This piece-by-piece technique is precisely what is done by the **dynamic programming algorithm**. We introduce the dynamic programming algorithm by applying it to the machine repair problem in Example 34, which is a basic finite horizon problem and so fits within this general scheme.

**Example 35 (Machine repair II).** In the machine repair problem, the terminal costs are  $C(1) = C(2) = C(3) = 0$  and for every state  $i = 1, 2, 3$ , and all stages  $n = 0, \dots, 4$ , we have two actions

$$\mathcal{U}_n(i) = \{\text{REPAIR}, \text{DON'T REPAIR}\} = \{R, D\},$$

at each decision stage. For each pair of states  $(i, j)$ , we have the transition probabilities  $p_{ij}(1)$  and  $p_{ij}(2)$  corresponding to actions 1 and 2, respectively. The complete set of transition probabilities are contained in the matrices

$$P(R) = \begin{pmatrix} 1 & 0 & 0 \\ 1 & 0 & 0 \\ 1 & 0 & 0 \end{pmatrix} \text{ and } P(D) = \begin{pmatrix} 0.7 & 0.2 & 0.1 \\ 0 & 0.6 & 0.4 \\ 0 & 0 & 1 \end{pmatrix}.$$

The dynamic programming algorithm proceeds as follows: we work out costs by starting with the last time step and working backwards.

**Stage T = 4:** No decision needs to be made because this is the end of the operation, and so  $J_4(i) = 0$ , for  $i = 1, 2, 3$ .

**Stage T – 1 = 3:** At the beginning of stage 3, the machine can be in any of the states 1, 2 or 3, depending on previous decisions.

- Suppose at the beginning of stage 3, the machine is in state 1:

$$u \in \{\text{REPAIR}, \text{DON'T REPAIR}\} = \{R, D\}.$$

Because the process terminates at time 4, we consider the expected single-stage cost only. So the single-stage cost that takes us from stage 3 to 4 is the only cost that needs to be considered for the first iteration of the dynamic programming algorithm.

Hence<sup>1</sup>,

$$\begin{aligned} J_3(1) &= \min_{u \in \{R, D\}} \mathbb{E}[c(1, u)] = \min \{\mathbb{E}[c(1, R)], \mathbb{E}[c(1, D)]\} \\ &= \min \{R + h(1), p_{11}(D)h(1) + p_{12}(D)h(2) + p_{13}(D)h(3)\} \\ &= \min [6 + 2, 0.7 \times 2 + 0.2 \times 4 + 0.1 \times 6] \\ &= \min [6 + 2, 2.8] \\ &= 2.8. \end{aligned}$$

This tells us that if the machine is in state 1 at time 3, then the optimal action is not to repair the machine. Taking this action minimises the expected remaining cost (to reach the termination stage), which is equal to 2.8. This optimal cost is sometimes referred to as the **cost-to-go**.

We now have to repeat this calculation for the other states.

- Suppose the machine is in state 2 at time 3 (we don't care how it got there, for now). Here, the dynamic programming algorithm requires us

---

<sup>1</sup>Note: we assume that the transition from state  $i$  to state  $j$  happens immediately once an action is made, rather than over a unit of time. So If we are in state  $i$  and time  $n$  and take an action, then we immediately transition into state  $j$  according to probability  $p_{ij}$  (said action), and we will be paying the maintenance cost for being in state  $j$  over time interval  $[n, n + 1]$ .

to set

$$\begin{aligned}
 J_3(2) &= \min_{u \in \{R, D\}} \mathbb{E}[c(2, u)] = \min \{\mathbb{E}[c(2, R)], \mathbb{E}[c(2, D)]\} \\
 &= \min \{R + h(1), p_{21}(D)h(1) + p_{22}(D)h(2) + p_{23}(D)h(3)\} \\
 &= \min \{6 + 2, 0.6 \times 4 + 0.4 \times 6\} \\
 &= \min \{8, 4.8\} = 4.8.
 \end{aligned}$$

The action  $u = D$  corresponding to DON'T REPAIR, is thus the action which should be selected if the machine is in state 2 at time 3, which incurs a cost-to-go of 4.8.

- Similarly,

$$\begin{aligned}
 J_3(3) &= \min_{u \in \{R, D\}} \mathbb{E}[c(3, u)] = \min \{\mathbb{E}[c(3, R)], \mathbb{E}[c(3, D)]\} \\
 &= \min \{R + h(1), p_{31}(D)h(1) + p_{32}(D)h(2) + p_{33}(D)h(3)\} \\
 &= \min [6 + 2, 6] \\
 &= 6;
 \end{aligned}$$

thus the optimal action is DON'T REPAIR, which incurs a cost-to-go of 6.

We have therefore identified the optimal actions and their associated optimal (minimum) costs for all states that the machine may find itself in at time 3. That is, we have identified the optimal costs

$$J_3(1) = 2.8, \quad J_3(2) = 4.8, \quad J_3(3) = 6,$$

which result from taking the optimal actions at stage 3.

**Stage T – 2 = 2:** At the beginning of stage 2, the machine can also be in any of the states 1, 2 or 3, depending on the previous decision (we don't care how it got there, for now).

- Suppose the machine is in state 1. Since there are two stages to go until termination, we now wish to select the action that minimises

the expected cost over both stages, rather than just minimising the expected single-stage cost from stage 2 to 3.

The quantity that we therefore need to minimise is

$$\begin{aligned} & \mathbb{E}[\text{single-stage cost from stage 2 to 3}] \\ & + \mathbb{E}[\text{cost from stage 3 to 4} \mid \text{the optimal action is taken at stage 3}]. \end{aligned}$$

The dynamic programming algorithm requires us to set

$$\begin{aligned} J_2(1) &= \min_{u \in \{R, D\}} \mathbb{E} \left[ c(1, u) + \sum_{j=1}^3 p_{1j}(u) J_3(j) \right] \\ &= \min \left\{ R + h(1) + \sum_{j=1}^3 p_{1j}(R) J_3(j), \quad p_{11}(D)h(1) + p_{12}(D)h(2) + p_{13}(D)h(3) + \sum_{j=1}^3 p_{1j}(D) J_3(j) \right\} \\ &= \min \{6 + 2 + 1 \times 2.8, \quad (0.7 \times 2 + 0.2 \times 4 + 0.1 \times 6) + (0.7 \times 2.8 + 0.2 \times 4.8 + 0.1 \times 6)\} \\ &= \min \{10.8, \quad 5.78\} = \mathbf{5.78}, \end{aligned}$$

which is obtained by selecting the action DON'T REPAIR.

- Suppose the machine is in state 2 at time 2, then

$$\begin{aligned} J_2(2) &= \min_{u \in \{R, D\}} \mathbb{E} \left[ c(2, u) + \sum_{j=1}^3 p_{2j}(u) J_3(j) \right] \\ &= \min \left\{ R + h(1) + \sum_{j=1}^3 p_{2j}(R) J_3(j), \quad p_{21}(D)h(1) + p_{22}(D)h(2) + p_{23}(D)h(3) + \sum_{j=1}^3 p_{2j}(D) J_3(j) \right\} \\ &= \min \{10, (0.6 \times 4 + 0.4 \times 6) + (0.6 \times 4.8 + 0.4 \times 6)\} = \min \{10, 10.08\} = \mathbf{10}, \end{aligned}$$

which is obtained by selecting the action REPAIR.

- Similarly, suppose the machine is in state 3 at time 2, then

$$\begin{aligned} J_2(3) &= \min_{u \in \{R, D\}} \mathbb{E} \left[ c(3, u) + \sum_{j=1}^3 p_{3j}(u) J_3(j) \right] \\ &= \min \{10, \quad h(3) + J_3(3)\} = \min \{10, 6 + 6\} = \mathbf{10}, \end{aligned}$$

so that in this case, the optimal action is also REPAIR.

We have now identified the optimal costs

$$J_2(1) = 5.78, \quad J_2(2) = 10, \quad J_2(3) = 10,$$

which result from taking the optimal actions at stages 2 and 3.

**Stage T – 3 = 1:** At time 1, there are three stages to go until termination, and we now select the action that minimises the expected cost over all three stages. At the beginning of stage 1, the machine can again be in any of the states 1, 2 or 3.

- Suppose the machine is in state 1. The dynamic programming algorithm requires us to set

$$\begin{aligned} J_1(1) &= \min_{u \in \{R, D\}} \mathbb{E} \left[ c(1, u) + \sum_{j=1}^3 p_{1,j}(u) J_2(j) \right] \\ &= \min \left\{ R + h(1) + \sum_{j=1}^3 p_{1j}(R) J_2(j), \quad p_{11}(D)h(1) + p_{12}(D)h(2) + p_{13}(D)h(3) + \sum_{j=1}^3 p_{1j}(D) J_2(j) \right\} \\ &= \min \{6 + 2 + 1 \times 5.78, \quad (0.7 \times 2 + 0.2 \times 4 + 0.1 \times 6) + (0.7 \times 5.78 + 0.2 \times 10 + 0.1 \times 10)\} \\ &= \min \{13.78, \quad 9.846\} = \mathbf{9.8460}. \end{aligned}$$

Hence, the optimal action is DON'T REPAIR, with associated cost-to-go of 9.8460.

- Similarly,

$$\begin{aligned} J_1(2) &= \min_{u \in \{R, D\}} \mathbb{E} \left[ c(2, u) + \sum_{j=1}^3 p_{2,j}(u) J_2(j) \right] \\ &= \min \left\{ R + h(1) + \sum_{j=1}^3 p_{2j}(R) J_2(j), \quad p_{21}(D)h(1) + p_{22}(D)h(2) + p_{23}(D)h(3) + \sum_{j=1}^3 p_{2j}(D) J_2(j) \right\} \\ &= \min \{6 + 2 + 1 \times 5.78, \quad (0.6 \times 4 + 0.4 \times 6) + (0.6 \times 10 + 0.4 \times 10)\} \\ &= \min \{13.78, \quad 14.8\} = \mathbf{13.78}, \end{aligned}$$



and

$$\begin{aligned}
 J_1(3) &= \min_{u \in \{R, D\}} \mathbb{E} \left[ c(3, u) + \sum_{j=1}^3 p_{3j}(u) J_2(j) \right] \\
 &= \min \{13.78, \quad h(3) + J_2(3)\} \\
 &= \min \{13.78, \quad 6 + 10\} = \mathbf{13.78},
 \end{aligned}$$

which are achieved by selecting REPAIR, for both states 2 and 3.

**Stage T – 4 = 0:** We are now in a position to calculate the optimal policy and corresponding costs  $J_0(i)$  for all possible initial states  $i$  of the machine.

The dynamic programming algorithm requires that

$$\begin{aligned}
 J_0(1) &= \min_{u \in \{R, D\}} \mathbb{E} \left[ c(1, u) + \sum_{j=1}^3 p_{1j}(u) J_1(j) \right] \\
 &= \min \left\{ R + h(1) + \sum_{j=1}^3 p_{1j}(R) J_1(j), \quad p_{11}(D)h(1) + p_{12}(D)h(2) + p_{13}(D)h(3) + \sum_{j=1}^3 p_{1j}(D) J_1(j) \right\} \\
 &= \min \{6 + 2 + 1 \times 9.846, \quad (0.7 \times 2 + 0.2 \times 4 + 0.1 \times 6) + (0.7 \times 9.846 + 0.2 \times 13.78 + 0.1 \times 13.78)\} \\
 &= \min \{17.846, \quad 13.8262\} = \mathbf{13.8262}.
 \end{aligned}$$

Similarly,

$$\begin{aligned}
 J_0(2) &= \min_{u \in \{R, D\}} \mathbb{E} \left[ c(2, u) + \sum_{j=1}^3 p_{2,j}(u) J_1(j) \right] \\
 &= \min \left\{ R + h(1) + \sum_{j=1}^3 p_{2j}(R) J_1(j), \quad p_{21}(D)h(1) + p_{22}(D)h(2) + p_{23}(D)h(3) + \sum_{j=1}^3 p_{2j}(D) J_1(j) \right\} \\
 &= \min \{6 + 2 + 1 \times 9.846, \quad (0.6 \times 4 + 0.4 \times 6) + (0.6 \times 13.78 + 0.4 \times 13.78)\} \\
 &= \min \{17.846, \quad 18.58\} = \mathbf{17.8460},
 \end{aligned}$$

and

$$\begin{aligned}
 J_0(3) &= \min_{u \in \{R,D\}} \mathbb{E} \left[ c(3, u) + \sum_{j=1}^3 p_{3j}(u) J_1(j) \right] \\
 &= \min \{17.846, \quad h(3) + J_1(3)\} \\
 &= \min \{17.846, \quad 6 + 13.78\} = \mathbf{17.846},
 \end{aligned}$$

So, we find that

$$J_0(1) = 13.8262, \quad J_0(2) = 17.846, \quad J_0(3) = 17.846,$$

which are achieved by selecting the action DON'T REPAIR if the machine starts in state 1, and selecting the action REPAIR if the machine starts in either state 2 or 3.

The above algorithm has allowed us to identify an optimal policy

$$\Psi^* = \{f_0^*, f_1^*, f_2^*, f_3^*\},$$

where the functions  $f_0^*, \dots, f_3^*$  are given in the Table 3.1, summarising optimal actions and costs for every possible combination of stage and machine state:

$i$	$J_0(i)$	$f_0^*(i)$	$J_1(i)$	$f_1^*(i)$	$J_2(i)$	$f_2^*(i)$	$J_3(i)$	$f_3^*(i)$	$J_4(i)$
1	13.8262	don't repair	9.846	don't repair	5.78	don't repair	2.8	don't repair	0
2	17.846	repair	13.78	repair	10	don't repair	4.8	don't repair	0
3	17.846	repair	13.78	repair	10	repair	6.00	don't repair	0

Table 3.1: Optimal Policy

Suppose you had to translate the optimal policy  $\Psi^*$  into words, for a manager who is not mathematically inclined. What would you say?

- Don't repair the machine if it is in perfect working order (state 1).

- If the machine is slightly damaged (state 2), only repair the machine if the **time-to-go** is greater than 2.
- If the machine is badly damaged (state 3), only repair the machine if the **time-to-go** is greater than 1.

The above policy makes sense, because as a general rule, if there is not much operation time remaining, the expected cost of investing in machine repair is greater than the expected cost of just letting the machine continue in a damaged state.

Of course, the optimal policy and costs are usually different if we change the problem parameters, such as lowering the cost of fixing the machine from  $R = 6$  to  $R = 1$ . The optimal policy is then to repair the machine whenever it is not in perfect working condition (state 1). (Check!)

**Theorem 3.2** (Bellman Optimisation Principle). *For every initial state  $i_0$ , the optimal cost  $J^*(i_0)$  is equal to  $J_0(i_0)$ , where the function  $J_0$  is given by the last step of the following algorithm, which proceeds backward in time from stage  $T - 1$  to stage 0:*

$$J_T(i) = C(i),$$

and then for  $k = 0, \dots, T - 1$ ,

$$J_k(i) = \min_{u_k \in \mathcal{U}_k(i)} \mathbb{E} \left[ c(i, u_k) + \sum_{j \in \mathcal{S}} p_{ij}(u_k) J_{k+1}(j) \right], \quad (3.3)$$

where  $C(i)$  is the terminal cost associated with state  $i$ , for  $i \in \mathcal{S}$ .

Furthermore, if  $u_k^* = f_k^*(i)$  minimises the RHS of (3.3) for each stage  $k = 0, \dots, T - 1$  and each state  $i \in \mathcal{S}$ , then the policy

$$\Psi^* = \{f_0^*, f_1^*, \dots, f_{T-1}^*\} \quad \text{is optimal.}$$

Theorem 3.2 is a consequence of the Markov property and is the basis of dynamic computing. The optimality equation (3.3) is a backward recursion in time, such that the later policy is decided first. This process of backward recursion calculating expected cost (reward), is also known as **Value Iteration**.

Let us now look at a manufacturing process where we will maximise rewards  $r$  over time. Of course we could minimise the costs  $c$  (negative of the rewards) and use the same notation as before, but it is more convenient to scribe it in terms of a maximisation problem.

**Example 36 (Designer furniture: A value maximisation problem).** A small company makes reproduction designer furniture and at the end of the year, the business is in one of two states

$$\{\text{THRIVING}, \text{FLOUNDERING}\} = \{1, 2\}.$$

At the beginning of each year, the company can make a decision as to whether they will invest money on design work or not.

If the company decides not to do design work, then the state transitions from the beginning of the year to the end of that same year form a Markov chain with transition matrix on  $\{1, 2\}$

$$P(\text{NO DESIGN}) = \begin{pmatrix} 0.5 & 0.5 \\ 0.3 & 0.7 \end{pmatrix}.$$

There is also a reward matrix that is associated with the decision NO DESIGN. If the company continues to THRIVE, sales provide a net income of 9. If it moves from THRIVING to FLOUNDERING, the reduced sales during the year provide a net income of only 3. If the business moves from FLOUNDERING to THRIVING, the sales provide an income of 3. If the business remains in the floundering state, the sales revenue does not cover the costs of operation, resulting in a reward of  $-4$ .

The resultant reward matrix under NO DESIGN is thus given by

$$R(\text{NO DESIGN}) = \begin{pmatrix} 9 & 3 \\ 3 & -4 \end{pmatrix}.$$

Notice that the rewards take no account of future earnings, as they are realised only over the current stage (year) and the transition probabilities are independent of the past history of the process.

It is possible to increase the chances of moving to (or staying in) the thriving state by spending time on design, which has an associated cost which reduces the rewards by 4. Hence, under DESIGN we get

$$P(\text{DESIGN}) = \begin{pmatrix} 0.75 & 0.25 \\ 0.7 & 0.3 \end{pmatrix},$$

$$R(\text{DESIGN}) = R(\text{NO DESIGN}) - \begin{pmatrix} 4 & 4 \\ 4 & 4 \end{pmatrix} = \begin{pmatrix} 5 & -1 \\ -1 & -8 \end{pmatrix}.$$

We label the two actions NO DESIGN and DESIGN that can be taken at the beginning on the year (stage) for convenience as NO and DES respectively, so that the set of actions are

$$U_n(i) = \{\text{NO}, \text{DES}\} \quad \text{for each state } i = 1, 2.$$

Suppose now that the lease for the premises of the company runs out at the end of  $T = 4$  years and that the company members have agreed to split up and go their separate ways. Hence, they must decide whether to spend time and money on design and their decision will be based on expected returns (rewards) or **expected monetary value (EMV)**.

In this example, there is no fixed reward for starting in any particular state at the beginning of a stage, because the reward depends on what state the company is in at the beginning and what state it is in at the end of that stage.

Let  $v(i, u_k)$  be the single stage expected reward given that the company starts in state  $i$ , at the beginning of stage  $k$  under decision (action)  $u_k$ . Also, let  $p_{ij}(u_k)$  and  $r_{ij}(u_k)$  be the corresponding entries in the transition probability matrix and reward matrix, respectively, under decision  $u_k$ .

The EMV  $V_k^\Psi(i)$  of a policy  $\Psi$  starting from state  $i \in \mathcal{S}$  at the beginning of stage  $k$  is defined by a **Bellman equation** given by

$$V_k^\Psi(i) = \mathbb{E} \left[ V_T(X_T^\Psi) + \sum_{n=k}^{T-1} v(X_n^\Psi, f(X_n^\Psi)) \right],$$

for  $k = 0, 1, \dots, T - 1$  and where  $X_k^\Psi = i$  and  $V_T(X_T^\Psi)$  is the value received for finishing in state  $X_T^\Psi$ .

This problem is similar to before, except now we want to find a policy  $\Psi$  that maximises EMV. Specifically, for some initial state  $i_0$ , we want an optimal policy

$$\Psi^* = \{f_0^*, f_1^*, \dots, f_{T-1}^*\}$$

that satisfies

$$V^*(i_0) = \max_{\Psi \in \Gamma} V_0^\Psi(i_0),$$

where  $\Gamma$  is the set of all admissible policies.

#### Stage $T = 4$ :

If we assume there is no advantage to be gained by finishing in any particular state, then  $V_4^*(1) = V_4^*(2) = 0$ .

#### Stage $T - 1 = 3$ :

$$\begin{aligned} \nu(1, \text{NO}) &= (0.5 \times 9 + 0.5 \times 3) + (0.5 \times V_4^*(1) + 0.5 \times V_4^*(2)) \\ &= (0.5 \times 9 + 0.5 \times 3) + 0 = 6, \\ \nu(1, \text{DES}) &= 0.75 \times 5 + 0.25 \times (-1) = 3.5, \\ \nu(2, \text{NO}) &= 0.3 \times 3 + 0.7 \times (-4) = -1.9, \\ \nu(2, \text{DES}) &= 0.7 \times (-1) + 0.3 \times (-8) = -3.1. \end{aligned}$$

Hence,

$$\begin{aligned} V_3^*(1) &= \max\{\nu(1, \text{DES}), \nu(1, \text{NO})\} = 6, \quad \text{under NO DESIGN,} \\ V_3^*(2) &= \max\{\nu(2, \text{DES}), \nu(2, \text{NO})\} = -1.9, \quad \text{under NO DESIGN.} \end{aligned}$$

#### Stage $T - 2 = 2$ :

$$\begin{aligned} \nu(1, \text{NO}) &= (0.5 \times 9 + 0.5 \times 3) + (0.5 \times 6 + 0.5 \times (-1.9)) = 8.05, \\ \nu(1, \text{DES}) &= (0.75 \times 5 + 0.25 \times (-1)) + (0.75 \times 6 + 0.25 \times (-1.9)) = 7.53, \\ \nu(2, \text{NO}) &= (0.3 \times 3 + 0.7 \times (-4)) + (0.3 \times 6 + 0.7 \times (-1.9)) = -1.43, \\ \nu(2, \text{DES}) &= (0.7 \times (-1) + 0.3 \times (-8)) + (0.7 \times 6 + 0.3 \times (-1.9)) = 0.53. \end{aligned}$$

Hence,

$$\begin{aligned}
 V_2^*(1) &= 8.05, & \text{under NO DESIGN} \\
 V_2^*(2) &= 0.53, & \text{under DESIGN,} \\
 V_1^*(1) &= 10.29, & \text{under NO DESIGN} \\
 V_1^*(2) &= 2.69, & \text{under DESIGN.} \\
 V_0^*(1) &= 12.49, & \text{under NO DESIGN} \\
 V_0^*(2) &= 4.91, & \text{under DESIGN.}
 \end{aligned}$$

In summary, we have:

$i$	$V_0^*(i)$	$f_0^*(i)$	$V_1^*(i)$	$f_1^*(i)$	$V_2^*(i)$	$f_2^*(i)$	$V_3^*(i)$	$f_3^*(i)$	$V_4^*(i)$
1	12.49	NO	10.29	NO	8.05	NO	6.00	NO	0
2	4.91	DES	2.69	DES	0.53	DES	-1.90	NO	0

Table 3.2: The Optimal Policy

This indicates that

- Never engage in design work if the company is thriving.
- If the company is floundering, do design work if in stages 0, 1 or 2. (Equivalently, if the company is floundering, do design work if there is more than 1 stage-to-go.)

## §3.4 INDEFINITE HORIZON PROBLEMS

## STOCHASTIC SHORTEST PATH MDP (SPMDP)

We now consider one type of problems where the horizon  $T$  is not deterministically known in advance, but is a random variable. In particular, we investigate a type of MDPs where the underlying Markov chain has an **absorbing state** and where every transition incurs some cost. This implies there are associated costs that are added at each decision stage. **Our aim** will be to find a policy that minimises the expected total costs, from all possible starting states, until absorption occurs.

In the two previous examples (one minimisation and one maximisation), we analysed policies of the form  $\Psi = \{f_0, f_1, \dots, f_{T-1}\}$ , which allowed us the flexibility of specifying a different action at each stage. These are **deterministic Markov policies**, where the qualifier **deterministic** refers to the fact that we choose precisely **one** action to take at each state  $i$  and each time  $t$ . A **randomised Markov policy** specifies one or more actions  $u_k$  to take at each state  $i$  and each time  $t$ , with associated probabilities that add up to 1.

In this study of SPMDPs, we restrict ourselves to the consideration of policies of the form

$$\Psi = \{f, f, \dots, f\},$$

which are **stationary** and usually denoted by  $\Psi = f \in \Gamma$ .



**Defn 3.4.** A **Stochastic Shortest Path MDP** with a finite state space and finite set of actions is defined as follows. There are a finite number of states  $i = 1, \dots, N$ , and an extra absorbing state  $d$ . For each state  $i$ , there is a finite set of actions  $\mathcal{U}(i)$  that can be selected. At the absorbing state  $d$ , no action needs to be chosen. Transitions from state  $i$  to state  $j$  are governed by the probabilities  $p_{ij}(u)$ , where  $u \in \mathcal{U}(i)$ . So the state transition probabilities at state  $i$  depend on the particular action  $u \in \mathcal{U}(i)$  that is selected. Transitions from state  $i$  to state  $j$  incur a cost  $r_{ij}$ .

The expected single-step cost incurred by selecting action  $u$  when in state  $i$  is therefore given by

$$\mathbb{E}[c(i, u)] = \sum_{j=1}^N p_{ij}(u) r_{ij}. \quad (3.4)$$

**Example 37** (A simple stochastic shortest path problem). A guide, Maizura, is standing on top of a mountain, looking down at a hiker, Ross, who is trying to get HOME, as quickly as possible. Consider Figure 3.3, depicting three locations (states), 1, 2 and 3.

States 1 and 2 are on opposite sides of a lake, and are connected by a bridge, which Ross takes 1 minute to traverse. State 3 is HOME, the destination/termination state/absorbing state. There is a trail directly joining state 1 to the destination state 3, but it passes through a swamp, and takes 2 minutes to traverse. There is a paved road directly joining state 2 to state 3, which only takes 0.7 minutes to traverse.

The actions that Maizura can take when Ross is at 1 or 2, are to yell CROSS THE BRIDGE, or to yell FOLLOW THE LAKE.

Suppose also that Ross has probability 1/4 of not hearing Maizura's instruction correctly, and in fact doing the opposite of what he is told. Assume Ross never remains in the same location after receiving an instruction, which

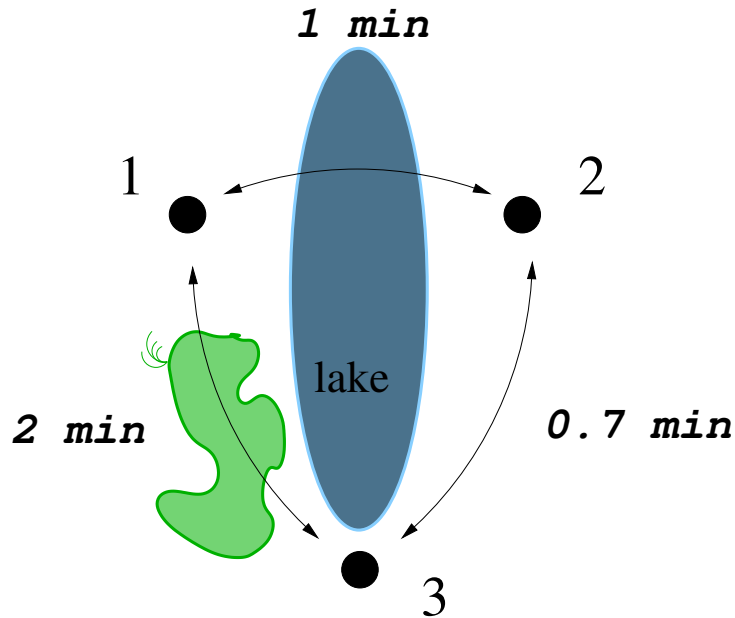


Figure 3.3: Hiker Problem.

means he always moves to another location. We are given the transition costs (time)

$$r_{12} = r_{21} = 1, \quad r_{13} = 2, \quad r_{23} = 0.7.$$

**Which instruction should Maizura yell out to Ross at each location?**

There are four possible stationary policies for this problem:

$$\begin{aligned} f_1 &= (u_1(1), u_1(2)) = (\text{FOLLOW LAKE}, \text{FOLLOW LAKE}), \\ f_2 &= (u_2(1), u_2(2)) = (\text{FOLLOW LAKE}, \text{CROSS BRIDGE}), \\ f_3 &= (u_3(1), u_3(2)) = (\text{CROSS BRIDGE}, \text{FOLLOW LAKE}), \\ f_4 &= (u_4(1), u_4(2)) = (\text{CROSS BRIDGE}, \text{CROSS BRIDGE}). \end{aligned}$$

The general form for the transition probability matrix corresponding to

the policy  $f_k$  is

$$P_{f_k} = \begin{pmatrix} p_{11}(f_k(1)) & p_{12}(f_k(1)) & p_{13}(f_k(1)) \\ p_{21}(f_k(2)) & p_{22}(f_k(2)) & p_{23}(f_k(2)) \\ 0 & 0 & 1 \end{pmatrix}.$$

Corresponding to each of the four possible policies are

$$P_{f_1} = \begin{pmatrix} 0 & \frac{1}{4} & \frac{3}{4} \\ \frac{1}{4} & 0 & \frac{3}{4} \\ 0 & 0 & 1 \end{pmatrix}, \quad P_{f_2} = \begin{pmatrix} 0 & \frac{1}{4} & \frac{3}{4} \\ \frac{3}{4} & 0 & \frac{1}{4} \\ 0 & 0 & 1 \end{pmatrix},$$

$$P_{f_3} = \begin{pmatrix} 0 & \frac{3}{4} & \frac{1}{4} \\ \frac{1}{4} & 0 & \frac{3}{4} \\ 0 & 0 & 1 \end{pmatrix}, \quad P_{f_4} = \begin{pmatrix} 0 & \frac{3}{4} & \frac{1}{4} \\ \frac{3}{4} & 0 & \frac{1}{4} \\ 0 & 0 & 1 \end{pmatrix}.$$

We will solve the problem in Example 37 by evaluating the costs associated with all possible policies, and identifying the policy that simultaneously minimises the expected cost from both of the possible starting states 1 and 2. Note that in this case, the cost measure is expressed in units of time, so that the total cost corresponds to the total amount of time taken for the hiker to reach the destination state. Of course, the cost could correspond to some other quantity, such as the amount of energy expended.

Let  $J^{f_k}(i)$  be the expected cost incurred by Ross in reaching the destination state 3 from state  $i$ , given that Maizura follows the policy  $f_k$ . Employing a first step analysis, we obtain the optimality equations

$$J^{f_k}(i) = \sum_{j=1}^3 p_{ij}(f_k(i)) (r_{ij} + J^{f_k}(j)) \quad \text{for } i = 1, 2, \quad (3.5)$$

with the boundary condition  $J^{f_k}(3) = 0$  for all policies  $f_k, k = 1, \dots, 4$ . This boundary condition makes sense physically, because the destination state is an absorbing state at which the Ross incurs no further cost.

The above system of equations can also be written

$$J^{f_k}(i) = c(i, f_k(i)) + \sum_{j=1}^3 p_{ij}(f_k(i)) J^{f_k}(j) \quad \text{for } i = 1, 2,$$

where the expected single-stage transition cost at state  $i$  under the stationary policy  $f_k$  is

$$c(i, f_k(i)) = \sum_{j=1}^3 p_{ij}(f_k(i)) r_{ij}.$$

We can then substitute the appropriate transition probabilities and state transition costs into (3.5) for each of the four policies, to obtain a system of equations.

For example, if  $k = 1$ ,

$$\begin{aligned} J^{f_1}(1) &= \frac{1}{4}(1) + \frac{3}{4}(2) + \frac{1}{4}J^{f_1}(2) + \frac{3}{4}J^{f_1}(3) = \frac{7}{4} + \frac{1}{4}J^{f_1}(2), \\ J^{f_1}(2) &= \frac{1}{4}(1) + \frac{3}{4}(0.7) + \frac{1}{4}J^{f_1}(1) + \frac{3}{4}J^{f_1}(3) = \frac{3.1}{4} + \frac{1}{4}J^{f_1}(1). \end{aligned}$$

Solving gives us

$$J^{f_1}(1) = 2.07333, \quad J^{f_1}(2) = 1.29333.$$

Similarly solving for all policies and states, we get the solutions shown in Table 3.3. The optimal policy, which minimises cost, is for the guide to always yell the instruction FOLLOW LAKE, which corresponds to  $f_1$ .

Note that if the Ross could hear the guide perfectly, thus always following the exact instruction, the optimal policy would be  $f_3$ . This can be deduced by inspection of the diagram, since the optimal path from state 2 to the destination is directly along the lake, and the optimal path from state 1 to the destination is to first cross the bridge to state 2 (and then to follow the lake to the destination).

$k$	$f_k(1)$	$f_k(2)$	$J^{f_k}(1)$	$J^{f_k}(2)$
1	FOLLOW LAKE	FOLLOW LAKE	2.07	1.29
2	FOLLOW LAKE	CROSS BRIDGE	2.44	2.75
3	CROSS BRIDGE	FOLLOW LAKE	2.25	1.34
4	CROSS BRIDGE	CROSS BRIDGE	4.44	4.26

Table 3.3: The stationary policies and corresponding costs

We have just solved the above stochastic shortest path MDP by enumerating all possible policies, evaluating their associated expected costs and identifying the policy that yields the minimum expected cost. This **exhaustive search** approach is clearly not feasible if the number of possible policies is large and so other techniques are needed for such problems (see later).

### §3.5 INFINITE HORIZON PROBLEMS

There are two types of Infinite Horizon MDPs: Average Value (Cost) MDPs, and Discounted MDPs.

#### INFINITE HORIZON AVERAGE VALUE (COST) MDPs

We now reconsider Example 36, where the company intends to operate indefinitely, and hence we do not have a predetermined value for  $T$ .

Observe that in the **Value Iteration**, the increase in  $V_t^*(i)$  between subsequent years appeared to be converging to a constant value, which we will assume exists and call it the **gain**  $g$ . This gain represents the average value of remaining in business for a year after some **relatively long time** and must be the same for all initial states  $i \in \mathcal{S}$ .

In Example 36:

$t$	1	2	3	4
$V_{t-1}^*(1) - V_t^*(1)$	2.20	2.24	2.05	6.00
average $\uparrow$	2.21	2.20	2.24	2.05
$V_{t-1}^*(2) - V_t^*(2)$	2.22	2.16	2.43	-1.90

In general, for a problem with  $n$  states, we can find the exact value of  $g$  by considering the transition matrix corresponding to the optimal policy.

First, note that finite-state Markov chains always have a stationary distribution. Let the stationary distribution of the transition matrix corresponding to the optimal policy be given by

$$\boldsymbol{\pi} = (\pi_1, \pi_2, \dots, \pi_n).$$

Then,  $\pi_j$  is the probability that after a relatively long time and under the optimal policy the system will be in state  $j$ . The gain is therefore the sum, over all possible states, of the products of the probability of being in each state and the expected reward over a single stage commencing in that state.

As the expected reward for starting in state  $i \in \{1, 2, \dots, n\}$  is given by

$$q_i = \sum_{j=1}^n p_{ij} r_{ij},$$

the gain is therefore

$$g = \sum_{i=1}^n \pi_i q_i.$$

When the company intends to operate indefinitely, it can be established that the optimal steady state policy is to DESIGN if the state of the company is floundering and NOT TO DESIGN otherwise. The steady state policy has the following transition and reward matrices

$$P_{opt} = \begin{pmatrix} 0.5 & 0.5 \\ 0.7 & 0.3 \end{pmatrix}, \quad R_{opt} = \begin{pmatrix} 9 & 3 \\ -1 & -8 \end{pmatrix}.$$

Solving the equilibrium equations  $\pi P_{opt} = \pi$ , we obtain

$$\pi = \left( \frac{7}{12}, \frac{5}{12} \right),$$

which leads to

$$g = \frac{7}{12}(0.5 \times 9 + 0.5 \times 3) + \frac{5}{12}(0.7 \times (-1) + 0.3 \times (-8)) = 2.208.$$

### POLICY ITERATION

The existence of a constant in steady state is the essence of something known as **Policy Iteration**, where we have problems that do not have a pre-determined value for  $T$ , but have a time-to-go that is infinite.

We have been performing backward recursions and calculating values at each stage, a process known as **Value iteration**. In the infinite horizon case, these values may become positively or negatively infinite. However, we observed that even though

$$\lim_{T \rightarrow \infty} \mathbb{E} \left[ \sum_{n=0}^{T-1} v(X_n, f_n(X_n) | X_0 = i) \right] \rightarrow \infty,$$

it may be true that

$$\lim_{T \rightarrow \infty} \frac{1}{T} \mathbb{E} \left[ \sum_{n=0}^{T-1} v(X_n, f_n(X_n) | X_0 = i) \right] \rightarrow g, \quad (3.6)$$

where  $-\infty < g < \infty$ .

### Policy Iteration Procedure

Consider that for a chosen policy, we have some finite horizon  $T$  with

$$V_k(i) = \sum_{j=1}^n p_{ij} (r_{ij} + V_{k+1}(j)) \quad (3.7)$$

and so for  $T \gg k$ , we have

$$V_k(i) = g(T - k) + \varphi(i), \quad (3.8)$$

where  $\varphi(i)$  is the value of starting in state  $i$  at time  $k$ . One of the consequences of (3.8) is that

$$V_k(i) = g + V_{k+1}(i), \quad (3.9)$$

which implies that the specific value of starting in state  $i$  given by  $\varphi(i)$  becomes negligible over time. Substituting (3.8) into (3.7) yields

$$\begin{aligned} (T - k)g + \varphi(i) &= \sum_{j=1}^n p_{ij} (r_{ij} + (T - k - 1)g + \varphi(j)) \\ \Leftrightarrow (T - k)g + \varphi(i) - (T - k - 1)g &= \sum_{j=1}^n p_{ij} (r_{ij} + \varphi(j)) \\ \Leftrightarrow \varphi(i) + g &= \left( \sum_{j=1}^n p_{ij} r_{ij} \right) + \left( \sum_{j=1}^n p_{ij} \varphi(j) \right). \end{aligned} \quad (3.10)$$

Equations (3.10) represent **the value determination operation**. The **Policy Improvement Routine** for each state  $i$ , is to find an alternative policy that maximises the expression on the RHS of (3.10).

#### Policy Improvement Routine (PIR)

Given the  $\varphi(j)$  for the current policy, we want to find an alternative policy  $f'$  that maximises

$$\left( \sum_{j=1}^n p'_{ij} r'_{ij} \right) + \left( \sum_{j=1}^n p'_{ij} \varphi(j) \right), \quad (3.11)$$

where  $p'_{ij}$  and  $r'_{ij}$  are, respectively, the probability transitions and rewards under the alternative policy  $\varphi'$ .



The policy iteration procedure revolves around (3.10) and (3.11). It is convenient to start with (3.10) using any policy (for example, a policy maximize immediate reward). Note that Equations (3.10) are a set of  $n$  linear equations in  $n + 1$  unknowns  $\varphi(1), \varphi(2), \dots, \varphi(n)$  and  $g$ . A unique solution may be found by setting  $\varphi(1) = 0$ , so that the other unknowns can be found and which represent relative values of starting in states  $i \neq 1$ . The next step is to apply the PIR criteria using (3.11).

The PIR will generally not immediately give the optimal policy unless our first guess (initial policy) was an exceptionally good choice, so that there is no alternative policy such that (3.11) yields an improved value. We will show later that the PIR can only yield an alternative policy that is an improvement. The procedure itself, which iterates between possible policies, stops or terminates when two policies are identical in subsequent iterations. Note that the **PIR** may also be used for minimisation problems.

### Policy Improvement Routine

- Step 1. Start with any policy.
- Step 2. Solve

$$\phi(i) + g = \left( \sum_{j=1}^n p_{ij} r_{ij} \right) + \left( \sum_{j=1}^n p_{ij} \phi(j) \right) \quad \text{to find } \phi(i) \text{ and } g.$$

- Step 3. Find an alternative policy  $f'$  and see if it improves

$$\left( \sum_{j=1}^n p'_{ij} r'_{ij} \right) + \left( \sum_{j=1}^n p'_{ij} \phi(j) \right)$$

- Step 4. If it does not, STOP. Otherwise, set  $p_{ij} = p'_{ij}$ ,  $r_{ij} = r'_{ij}$  and go back to Step 2.

**Example 38 (Designer Furniture Manufacturer).** We revisit the company and consider NO DESIGN as our initial policy.

**Iteration 1:** Using Equation (3.10) gives us the system

$$\begin{aligned} g + \varphi(1) &= 0.5 \times 9 + 0.5 \times 3 + 0.5\varphi(1) + 0.5\varphi(2), \\ g + \varphi(2) &= 0.3 \times 3 + 0.7 \times (-4) + 0.3\varphi(1) + 0.7\varphi(2). \end{aligned}$$

Setting  $\varphi(1) = 0$ , we obtain, for  $i = 1$  and for  $i = 2$ , respectively,

$$g = 0.5 \times 9 + 0.5 \times 3 + 0.5\varphi(2), \quad (3.12)$$

$$g + \varphi(2) = 0.3 \times 3 + 0.7 \times (-4) + 0.7\varphi(2). \quad (3.13)$$

The difference (3.12) – (3.13) removes the gain  $g$  and gives

$$\varphi(2) = 0.9 - 2.8 - 4.5 - 1.5 + 0.2\varphi(2) \Rightarrow \varphi(2) = -9.875.$$

Under DESIGN, using  $\varphi(1) = 0$ ,  $\varphi(2) = -9.875$  and (3.10) leads to

$$\begin{aligned} 3.75 - 0.25 + 0 + (0.25)(-9.875) &= 1.0312 \quad \text{for } i = 1, \\ -0.7 - 2.4 + 0 + (0.3)(-9.875) &= -6.0625 \quad \text{for } i = 2. \end{aligned}$$

As  $1.0625 > 1.0312$  and  $-6.0625 > -8.8125$ , our new policy is NO DESIGN if in state  $i = 1$ , and DESIGN if in state  $i = 2$ .

**Iteration 2:** Under this new policy, setting  $\varphi(1) = 0$ , Equation (3.10) gives

$$g = (0.5)(9) + (0.5)(3) + 0.5\varphi(2), \quad (3.14)$$

$$g + \varphi(2) = (0.7)(-1) + (0.3)(-8) + 0.3\varphi(2). \quad (3.15)$$

The difference (3.14) – (3.15) yields

$$\varphi(2) = -0.7 - 2.4 - 4.5 - 1.5 - 0.2\varphi(2) \Rightarrow \varphi(2) = -7.5833.$$

With  $\varphi(1) = 0$ ,  $\varphi(2) = -7.5833$ , under NO DESIGN

$$\begin{aligned} 4.5 + 1.5 + 0 + (0.5)(-7.5833) &= 2.2083 \quad \text{for } i = 1, \\ 0.9 - 2.8 + 0 + (0.7)(-7.5833) &= -7.2083 \quad \text{for } i = 2. \end{aligned}$$

and under DESIGN

$$\begin{aligned} 3.75 - 0.25 + 0 + (0.25)(-7.5833) &= 1.6042 \quad \text{for } i = 1, \\ -0.7 - 2.4 + 0 + (0.3)(-7.5833) &= -5.3750 \quad \text{for } i = 2. \end{aligned}$$

As  $2.2083 > 1.6042$  and  $-5.3750 > -7.2083$ , our new policy is NO DESIGN if in state  $i = 1$ , and DESIGN if in state  $i = 2$ . This is identical to the previous policy, and so we stop the iterations as the optimal policy has been found.

This optimal policy is what we might intuitively conclude from the previous investigation for the problem with  $T = 4$ , where the company only did design work if it was FLOUNDERING and there was more than 1 stage-to-go. (There is always more than 1 stage-to-go in the infinite horizon case!)

Note that the policy just discovered is a stationary policy, and gives us the action to be taken when in each of the states of the process *independent* of the time that we are in each of these states.

**Defn 3.5** (Average value/cost MDP). *Consider an MDP with a finite number of states, for which the underlying Markov chain is positive recurrent, and the termination time  $T$  is effectively infinite. In this case, each state is visited an infinite number of times.*

*If the expected value (cost) of all single-stage transitions is strictly positive, it is clear that the total value (cost) of any policy  $\Psi$ , is infinite, which tells us nothing about the different policies.*

*However, the **average value (cost) per stage**  $g$  is in many cases well defined and given by*

$$g = \lim_{T \rightarrow \infty} \frac{1}{T} \mathbb{E} \left[ \sum_{n=0}^T v(X_n, f_n(X_n)) \right].$$

*An MDP that finds an optimal policy that maximises (minimises)  $g$  is known as an **average value (cost) MDP**.*

### DISCOUNTED MDPs

Let  $\beta \in (0, 1]$  be a *discount factor*, which acts as a measure of the importance that we place on costs incurred into the future, as opposed to costs incurred at the present time. As a practical example, consider an asset whose value depreciates over time. In a decision process involving such an asset, we must factor in this depreciation when we evaluate the costs of different policies. We can define the cost criterion factoring in  $\beta \in (0, 1]$  over a finite horizon  $T$  as

$$\sum_{n=0}^{T-1} \beta^n c(X_n, u_n) + \beta^T C(X_T).$$

Suppose the costs are uniformly bounded, such that

$$|c(X_n, u_n)| < B \quad \text{for all } X_n \in \mathcal{S} \text{ and } n \geq 0.$$

Then, if discounting is strict so that  $0 < \beta < 1$ , the infinite horizon cost is bounded by  $B/(1 - \beta)$ , which we can show by using geometric series.

**The minimal present value**  $J_k^\Psi(i)$  at time  $k$  of expected cost from time  $k$  up to time  $T$ , starting in state  $i$  at time  $k$  under policy  $\Psi$ , is given by

$$\inf_{\Psi} \mathbb{E} \left[ \sum_{n=k}^{T-1} \beta^{n-k} c(X_n^\Psi, f_n(X_n^\Psi)) + \beta^{T-k} C(X_T) \mid X_k = i \right]. \quad (3.16)$$

We then obtain the optimality equations

$$J_T(i) = C(X_T), \quad (3.17)$$

$$J_k(i) = \inf_{u_k \in \mathcal{U}_k(i)} \{ \mathbb{E}[c(i, u_k)] + \beta \mathbb{E}[J_{k+1}(X_{k+1})] \} \quad \text{for } k \leq T. \quad (3.18)$$

### The $s$ -horizon cost $F_s^\Psi(i)$ under policy $\Psi$

In the finite horizon case, the value function can simply be obtained from (3.17) and (3.18) using backward recursions from the terminal time  $T$ .

Consider a stationary policy  $\Psi$ , where  $c(i, u_k) = c(i, u)$  for all  $k \geq 0$ , and assume that  $C(X_T) = 0$  for all final states, so that we can define the  **$s$ -horizon cost  $F_s^\Psi(i)$  under policy  $\Psi$** :

$$F_s^\Psi(i) = \mathbb{E} \left[ \sum_{n=0}^{s-1} \beta^n c(X_n, u) \mid X_0 = i \right]. \quad (3.19)$$

The subscript  $s$  means that we are looking at the cost in terms of  $s$  steps-to-go. If we take the infimum with respect to  $\Psi$ , we obtain the optimal  $s$ -horizon cost

$$F_s(i) = \inf_{\Psi \in \Gamma} F_s^\Psi(i).$$

Note that this infimum always exists and satisfies

$$F_s(i) = \inf_u \{ \mathbb{E}[c(i, u)] + \beta \mathbb{E}[F_{s-1}(X_1) \mid X_0 = i] \},$$

with terminal condition  $F_0(i) = 0$  for all  $i$ .

#### The infinite-horizon cost $F^\Psi(i)$ under policy $\Psi$

We now consider the case where the discounted MDP has an infinite horizon, or where  $s \rightarrow \infty$ . The **infinite horizon cost  $F^\Psi(i)$  under policy  $\Psi$**  is given by

$$F^\Psi(i) = \lim_{s \rightarrow \infty} F_s^\Psi(i). \quad (3.20)$$

The limit in Equation (3.20) need not always exist, so we will now investigate conditions under which it does. To this end, we will formulate everything in terms of the minimisation problem.

**Defn 3.6.** *Minimisation problems can be classified into three categories.*

**D:** *Discounted Programming.*

$$0 < \beta < 1$$

and  $|c(i, u)| < B$  for all  $i \in \mathcal{S}$  and  $u \in \mathcal{U}(i)$ .

**N:** *Negative Programming.*

$$0 < \beta \leq 1$$

and  $c(i, u) \geq 0$  for all  $i \in \mathcal{S}$  and  $u \in \mathcal{U}(i)$ .

**P:** *Positive Programming.*

$$0 < \beta \leq 1$$

and  $c(i, u) \leq 0$  for all  $i \in \mathcal{S}$  and  $u \in \mathcal{U}(i)$ .

The terms **positive** and **negative** arises from considering that

- maximising positive rewards (P) is equivalent to minimising negative costs;
- maximising negative rewards (N) is equivalent to minimising positive costs.

In cases of  $N$  and  $P$ , wlog, we will usually take  $\beta = 1$ .

The existence of the limit in (3.20) (possibly  $\pm\infty$ ) is assured in cases  $N$  and  $P$  by monotone convergence and in the case  $D$  because the total cost occurring after the  $s$ th step is bounded by

$$\beta^s \frac{B}{1-\beta} \quad \text{and} \quad \lim_{s \rightarrow \infty} \beta^s \frac{B}{1-\beta} \rightarrow 0.$$

The infimal infinite-horizon cost is defined by

$$F(i) = \inf_{\Psi \in \Gamma} F^\Psi(i) = \inf_{\Psi \in \Gamma} \lim_{s \rightarrow \infty} F_s^\Psi(i). \quad (3.21)$$

**Theorem 3.3.** Suppose  $D$ ,  $N$  or  $P$  hold, then

$$F(i) = \inf_u \{ \mathbb{E}[c(i, u)] + \beta \mathbb{E}[F(X_1) \mid X_0 = i, u_0 = u] \}. \quad (3.22)$$

*Proof.* First, we show that  $F(i) \geq \text{RHS of (3.22)}$ , and then  $F(i) \leq \text{RHS of (3.22)}$ .

(a) **To show**  $F(i) \geq \text{RHS of (3.22)}$ :

Suppose  $\Psi$  is a policy that chooses  $u_0 = u$  when  $X_0 = i$ , then

$$F_s^\Psi(i) = \mathbb{E}[c(i, u)] + \beta \mathbb{E}[F_{s-1}^\Psi(X_1) \mid X_0 = i, u_0 = u] \quad (3.23)$$

Now, either  $D$ ,  $N$  or  $P$  is sufficient to allow us to take limits on both sides of (3.23) and interchange the order of limit and expectation. In cases  $N$  and  $P$ , this is because of monotone convergence, where infinity is allowed as a possible limiting value.

We obtain the following

$$\begin{aligned} F^\Psi(i) &= \mathbb{E}[c(i, u)] + \beta \mathbb{E}[F^\Psi(X_1) \mid X_0 = i, u_0 = u] \\ &\geq \mathbb{E}[c(i, u)] + \beta \mathbb{E}[F(X_1) \mid X_0 = i, u_0 = u] \\ &\geq \inf_u \{ \mathbb{E}[c(i, u)] + \beta \mathbb{E}[F(X_1) \mid X_0 = i, u_0 = u] \}. \end{aligned}$$

Taking the infimum the LHS over  $\Psi \in \Gamma$  gives us  $F(i) \geq \text{RHS of (3.22)}$ .

(b) **To show**  $F(i) \leq \text{RHS of (3.22)}$ : We consider another policy  $\Psi^*$  that after choosing some  $u_0$  when  $X_0 = i$ , reaches state  $X_1$  and chooses a different policy  $\Psi_1$ , which is suboptimal by less than some  $\varepsilon$  from then on; that is,

$$F^{\Psi_1}(X_1) \leq F(X_1) + \varepsilon.$$

Such a policy must exist by definition of  $F$ , although  $\Psi_1$  will depend on  $X_1$ . Thus, we have

$$\begin{aligned}
 F(i) &\leq F^{\Psi^*}(i) \\
 &= \mathbb{E}[c(i, u)] + \beta \mathbb{E}[F^{\Psi_1}(X_1) \mid X_0 = i, u_0 = u] \\
 &\leq \mathbb{E}[c(i, u)] + \beta \mathbb{E}\left[F(X_1) + \varepsilon \mid X_0 = i, u_0 = u\right] \\
 &\leq \mathbb{E}[c(i, u)] + \beta \mathbb{E}[F(X_1) \mid X_0 = i, u_0 = u] + \beta \varepsilon
 \end{aligned}$$

which implies

$$F(i) \leq \inf_u \{\mathbb{E}[c(i, u)] + \beta \mathbb{E}[F(X_1) \mid X_0 = i, u_0 = u]\} + \beta \varepsilon.$$

Recalling that  $\varepsilon$  is arbitrary yields the result.  $\square$

**Example 39 (Factory).** At a factory, a machine can be operated in either AUTOMATIC MODE or MANUAL MODE. If it is operated in AUTOMATIC MODE on a particular day, there is a probability  $q$  that it will fail to start the next day. Similarly, if it is operated in MANUAL MODE, there is a probability  $p < q$  (otherwise why would you operate manually?) that it will fail to start the next day.

Running the machine in AUTOMATIC MODE costs \$0 per day, while running in MANUAL MODE costs \$ $c$  per day for the attendant. If the machine fails to start on a particular day, it must be repaired, a process that takes the entire day at a cost of \$ $8c$ , due to lost production, labour, parts etc. Costs are discounted by a factor  $\beta \in (0, 1)$ .

We would like to minimise the expected discounted cost over an essentially infinite horizon. Let  $F(0)$  and  $F(1)$  denote the minimal costs from NOT WORKING in the morning and WORKING in the morning respectively. Then, the state space is  $\mathcal{S} = \{\text{NOT WORKING}, \text{WORKING}\} \equiv \{0, 1\}$ .

**If we want to show that it is optimal to run the machine continuously in automatic mode rather than continuously in manual mode,** we need to determine some conditions based on the parameterisations that we have. That



is, using the two possible modes (actions) {MANUAL, AUTOMATIC}, Theorem 3.3 says that we can write down the following

$$\begin{aligned} F(0) &= \inf \{8c + \beta [pF(0) + (1-p)F(1)], 8c + \beta [qF(0) + (1-q)F(1)]\}, \\ F(1) &= \inf \{c + \beta [pF(0) + (1-p)F(1)], \beta [qF(0) + (1-q)F(1)]\}. \end{aligned}$$

Then operating in MANUAL MODE only we have

$$F^M(0) = 8c + \beta [pF^M(0) + (1-p)F^M(1)], \quad (3.24)$$

$$F^M(1) = c + \beta [pF^M(0) + (1-p)F^M(1)]. \quad (3.25)$$

The difference (3.24)-(3.25) gives us

$$F^M(0) - F^M(1) = 7c. \quad (3.26)$$

Subsequently using (3.26) in (3.24) and (3.25) leads to

$$F^M(0) = \frac{c(8-7\beta(1-p))}{1-\beta}, \quad F^M(1) = \frac{c(1+7\beta p)}{1-\beta}.$$

Similarly, operating in AUTOMATIC MODE only it can be shown that

$$F^A(0) = \frac{8c(1-\beta(1-q))}{1-\beta}, \quad F^A(1) = \frac{8c\beta q}{1-\beta}.$$

If running continuously in AUTOMATIC MODE is better, we must have

$$\begin{aligned} F^A(1) &\leq F^M(1) \\ \Leftrightarrow \frac{8c\beta q}{(1-\beta)} &< \frac{c(1+7\beta p)}{(1-\beta)} \\ \Leftrightarrow \beta &\leq \frac{1}{(8q-7p)}, \end{aligned}$$

since  $8q - 7p > 0$ . Again, if running continuously in AUTOMATIC MODE is better, we must also have that

$$\begin{aligned}
 F^A(0) &\leq F^M(0) \\
 \Leftrightarrow \frac{8c(1 - \beta(1 - q))}{(1 - \beta)} &\leq \frac{c(8 - 7\beta(1 - p))}{(1 - \beta)} \\
 \Leftrightarrow -8(1 - q) &\leq -7(1 - p) \\
 \Leftrightarrow (8q - 7p) &\leq 1,
 \end{aligned}$$

which means  $0 \leq 8q - 7p \leq 1$ .

We should of course check other stationary policies to find the optimal policy. That is, cross over policies, where we run in one mode if in one state and the other mode if in the other state (see later).

### §3.6 POSITIVE PROGRAMMING & THE VALUE ITERATION ALGORITHM

#### POSITIVE PROGRAMMING

Positive programming concerns maximising non-negative rewards or minimising non-positive costs and in some cases, there is no optimal policy.

**Example 40 (No optimal policy).** Suppose that  $\mathcal{S} = \mathbb{Z}^+$ , the set of non-negative integers. When in state  $i$ , we have a choice of either

- moving to state  $i + 1$  and receive no reward, or
- moving to state 0, receiving a reward  $(1 - 1/i)$ , and then remaining in state 0 thereafter with no further reward.

The optimality equation may be written

$$F(i) = \sup \left\{ \left(1 - \frac{1}{i}\right), F(i + 1) \right\} \quad \text{for } i > 0.$$

Clearly,  $F(i) = 1, i > 0$ , but the policy that chooses the maximising action always moves up one state and hence receives zero reward. Hence, there is no policy that yields the reward of 1.

## CHARACTERISATION OF THE OPTIMAL POLICY

The following theorem provides a necessary and sufficient condition for a policy to be optimal.

**Theorem 3.4.** *Suppose that  $D$  or  $P$  holds, and that  $\Psi$  is a policy whose value function  $F^\Psi(i)$  satisfies the optimality equation:*

$$F^\Psi(i) = \sup_u \{ \mathbb{E}[r(i, u)] + \beta \mathbb{E}[F^\Psi(X_1) | X_0 = i, u_0 = u] \}.$$

*Then  $\Psi$  is the optimal policy. (Note: This is a maximisation problem.)*

*Proof.* Let  $\Psi_1$  be any other policy and suppose it takes  $u_t(i) = f_t(i)$ . Since  $F^\Psi(i)$  satisfies the optimality equation,

$$\begin{aligned} F^\Psi(i) &= \sup_u \{ \mathbb{E}[r(i, u)] + \beta \mathbb{E}[F^\Psi(X_1) | X_0 = i, u_0 = u] \} \\ &\geq \underbrace{\mathbb{E}[r(i, f_0(i))]}_{\text{expected single-stage reward under policy } \Psi_1} + \beta \mathbb{E}[F^\Psi(X_1) | X_0 = i, u_0 = f_0(i)]. \end{aligned}$$

By repeated substitution of this into itself, we find

$$F^\Psi(i) \geq \mathbb{E}_{\Psi_1} \left[ \sum_{t=0}^{s-1} \beta^t r(X_t, u_t) \middle| X_0 = i \right] + \beta^s \mathbb{E}[F^\Psi(X_s) | X_0 = i], \quad (3.27)$$

where the subscript  $\Psi_1$  in the first expectation on the RHS indicates that the actions  $u_t$  at each time  $t = 0, \dots, s-1$  are determined by the policy  $\Psi_1$ .

**In the case  $P$ ,** we know that  $0 < \beta \leq 1$  and that  $c(i, u) \leq 0$ , the latter implying  $r(i, u) \geq 0$  (as  $r(i, u) := -c(i, u)$ ). Consequently,

$$\beta^s \mathbb{E}[F^\Psi(X_s) | X_0 = i] \geq 0.$$

Hence, we have

$$F^\Psi(i) \geq \mathbb{E}_{\Psi_1} \left[ \sum_{t=0}^{s-1} \beta^t r(X_t, u_t) | X_0 = i \right].$$

Letting  $s \rightarrow \infty$  and interchanging the limit and the expectation operators on the RHS (which we are allowed to do because of monotone convergence!), we have

$$F^\Psi(i) \geq F^{\Psi_1}(i).$$

**In the case D**, we have  $0 < \beta < 1$  and  $c(i, u)$  being uniformly bounded. Thus, we can directly let  $s \rightarrow \infty$  and observe that

$$\lim_{s \rightarrow \infty} \beta^s \mathbb{E} [F^\Psi(X_s) | X_0 = i] = 0,$$

which implies that

$$F^\Psi(i) \geq F^{\Psi_1}(i).$$

□

**Example 41 (Optimal gambling).** A gambler has  $\$i$  and wants to increase this to  $\$N$ . At each stage they can bet any whole number of  $\$$  not exceeding their capital, say  $j \leq i$ . If they win, with probability  $p$ , they now have  $\$(i + j)$ , or if they lose, with probability  $q = (1 - p) > 0$ , they now have  $\$(i - j)$ . There are no draws, only wins or losses.

Suppose  $p \geq 1/2$ , and we have to prove that the **timid strategy**  $\Psi_1$ , of always betting only  $\$1$ , is optimal for the gambler to attain  $\$N$ . The optimality equation for maximising the probability of attaining  $\$N$  is

$$F(i) = \max_{j, j \leq i} \{pF(i + j) + qF(i - j)\},$$

with  $F(0) = 0$  and  $F(N) = 1$ .

To show that the timid strategy,  $\Psi_1$ , is optimal we need to find its value function, say  $G(i) = F^{\Psi_1}(i)$ , and then show that it is a solution to the optimality

equation. We have that  $G(i) = pG(i+1) + qG(i-1)$ , with  $G(0) = 0$  and  $G(N) = 1$  and trying a solution of the form  $G(i) = w^i$  for some  $w \in R$  yields the quadratic

$$\frac{w^i}{w^{i-1}} = p \frac{w^{i+1}}{w^{i-1}} + q \frac{w^{i-1}}{w^{i-1}} \Rightarrow (w-1) \left( w - \frac{q}{p} \right) = 0,$$

so that  $w = 1$  or  $w = \frac{q}{p}$ , which for  $p > \frac{1}{2}$  (by assumption) gives us

$$G(i) = A \left( \frac{q}{p} \right)^i + B, \quad \text{for } A, B \in R.$$

Using the boundary conditions

$$\begin{aligned} G(0) = 0 &\Rightarrow B = -A, \\ G(N) = 1 &\Rightarrow A = \frac{1}{\left( \frac{q}{p} \right)^N - 1}, \end{aligned}$$

we obtain

$$G(i) = \frac{1 - (q/p)^i}{1 - (q/p)^N}, \quad \text{for } p > \frac{1}{2}.$$

When  $p = 1/2$ , we use the boundary conditions in  $G(i) = Ai + B$ , to give us the final part of the problem, which yields

$$G(i) = \begin{cases} \frac{1 - (q/p)^i}{1 - (q/p)^N}, & p > \frac{1}{2}, \\ \frac{i}{N}, & p = \frac{1}{2}. \end{cases}$$

If  $p = 1/2 = q$ , then  $G(i) = i/N$  clearly satisfies the optimality equation:

$$\max_{j, j \leq i} \left\{ \frac{1}{2} \left[ \frac{i+j}{N} \right] + \frac{1}{2} \left[ \frac{i-j}{N} \right] \right\} = \max_{j, j \leq i} \left\{ \frac{1}{2} \left[ \frac{i+i}{N} \right] + \frac{1}{2} \left[ \frac{j-j}{N} \right] \right\} = \frac{i}{N}.$$

If  $p > 1/2$ , we have to verify that  $G(i) = \frac{1 - (q/p)^i}{1 - (q/p)^N}$  satisfies

$$G(i) = \max_{j, j \leq i} \left\{ p \left[ \frac{1 - (q/p)^{i+j}}{1 - (q/p)^N} \right] + q \left[ \frac{1 - (q/p)^{i-j}}{1 - (q/p)^N} \right] \right\}. \quad (3.28)$$

Let  $W_j$  be the expression inside the braces  $\{ \}$  on the RHS of (3.28), then we can easily show that  $W_{j+1} < W_j$ , for  $1 \leq j \leq i$ . In particular, since  $p > q$  and  $p, q, N > 0$ , this is equivalent to establishing that for all  $1 \leq k \leq i$ ,

$$\begin{aligned} p \left( \frac{q}{p} \right)^{i+k} + q \left( \frac{q}{p} \right)^{i-k} &< p \left( \frac{q}{p} \right)^{i+k+1} + q \left( \frac{q}{p} \right)^{i-k-1} \\ \Rightarrow p \left( \frac{q}{p} \right)^{i+k} \left( 1 - \left( \frac{q}{p} \right) \right) &< q \left( \frac{q}{p} \right)^{i-k-1} \left( 1 - \left( \frac{q}{p} \right) \right) \\ \Rightarrow p \left( \frac{q}{p} \right)^{i+k} &< q \left( \frac{q}{p} \right)^{i-k-1} \end{aligned}$$

as all terms are positive because  $1 > p > 1/2$ .

Then,

$$\begin{aligned} p \left( \frac{q}{p} \right)^{i+k} &< q \left( \frac{q}{p} \right)^{i-k-1} \\ \Leftrightarrow \left( \frac{q}{p} \right)^{i-k-1} \left( p \left( \frac{q}{p} \right)^{2k+1} - q \right) &< 0 \\ \Leftrightarrow q \left( \left( \frac{q}{p} \right)^{2k} - 1 \right) &< 0, \end{aligned}$$

which directly follows because  $\left( \frac{q}{p} \right)^{2k} < 1$  for all  $k > 0$ .

Hence,  $j = 1$  maximises the RHS of (3.28) as required and thus the timid strategy  $\Psi_1$  maximises the probability of making  $\$N$ .

### THE VALUE ITERATION ALGORITHM

Recall now the **Value Iteration method** that we used on finite horizon problems. Starting with  $F_0(i) = 0$ , we worked backwards iteratively for  $s = 1, 2, \dots$  to calculate

$$F_s(i) = \inf_{\Psi} \left\{ \mathbb{E}[c(i, u)] + \beta \mathbb{E} \left[ F_{s-1}^{\Psi}(X_1) \mid X_0 = i, u_0 = u \right] \right\}.$$

So  $F_s(i)$  is the infimal cost over  $s$  steps. Now let

$$F_{\infty}(i) = \lim_{s \rightarrow \infty} F_s(i) = \lim_{s \rightarrow \infty} \inf_{\Psi} F_s^{\Psi}(i) \quad (3.29)$$

The limit exists, either by monotone convergence under  $N$  or  $P$ , or by the fact that  $\lim_{s \rightarrow \infty}$  of cost incurred  $\rightarrow 0$  under  $D$ . Note that (3.29) reverses the order of  $\lim_{s \rightarrow \infty}$  and  $\inf_{\Psi}$  in (3.21), which is the infimal infinite-horizon cost

$$F(i) = \inf_{\Psi \in \Gamma} F^{\Psi}(i) = \inf_{\Psi \in \Gamma} \lim_{s \rightarrow \infty} F_s^{\Psi}(i).$$

The following theorem states that we can interchange the order of the operations and therefore  $F_s(i) \rightarrow F(i)$ . In the case  $N$ , we will need an additional assumption:

**Defn 3.7** (Assumption A: Finite actions). *Under assumption A, there are only finitely many possible values of  $u$  (finitely many actions) in each state.*

**Theorem 3.5.** *Suppose  $D$  or  $P$  holds, or  $N$  and Assumption A hold. Then,  $F_{\infty}(i) = F(i)$ .*

*Proof.* We prove first ‘ $\leq$ ’, and then ‘ $\geq$ ’ to get the result.

(‘ $\leq$ ’) Given any policy  $\bar{\Psi}$ ,

$$F_{\infty}(i) = \lim_{s \rightarrow \infty} F_s(i) = \lim_{s \rightarrow \infty} \inf_{\Psi} F_s^{\Psi}(i) \leq \lim_{s \rightarrow \infty} F_s^{\bar{\Psi}}(i) = F^{\bar{\Psi}}(i).$$

Taking the infimum over  $\bar{\Psi}$  gives  $F_{\infty}(i) \leq F(i)$ .

(‘ $\geq$ ’) **In the case  $P$ ,** we have  $c(i, u) \leq 0$ , so  $F_s(i) \geq F(i)$ . Now, let  $s \rightarrow \infty$ , then  $F_{\infty}(i) \geq F(i)$ .

**In the case  $D$ ,** with  $|c(i, u)| < B$ , imagine subtracting  $B > 0$  from every cost, which reduces the infinite-horizon cost under any policy by exactly  $B/(1 - \beta)$ , and  $F(i)$  and  $F_{\infty}(i)$  also decrease by this amount. All costs are now negative, so the result above in the case  $P$ , again applies.

Alternatively, note that

$$F_s(i) - \beta^s B/(1 - \beta) \leq F(i) \leq F_s(i) + \beta^s B/(1 - \beta)$$

and hence  $\lim_{s \rightarrow \infty} F_s(i) = F(i)$ .

**In the case  $N$  with  $A$ ,**

$$\begin{aligned} F_{\infty}(i) &= \lim_{s \rightarrow \infty} \min_u \{c(i, u) + \mathbb{E}[F_{s-1}(X_1)|X_0 = i, u_0 = u]\} \\ &= \min_u \{c(i, u) + \lim_{s \rightarrow \infty} \mathbb{E}[F_{s-1}(X_1)|X_0 = i, u_0 = u]\} \end{aligned} \quad (3.30)$$

$$= \min_u \{c(i, u) + \mathbb{E}[F_{\infty}(X_1)|X_0 = i, u_0 = u]\} \quad (3.31)$$

where Equality (3.30) follows because the minimum is over a finite number of terms (condition  $A$ ), and Equality (3.31) follows by monotone convergence (since  $F_s(i)$  increases in  $s$ ).

Now let  $\Psi^*$  be the policy that chooses the minimising action on the RHS of (3.31). This implies, by substitution of (3.31) into itself, and using the fact that  $N$  implies  $F_{\infty} \geq 0$ ,

$$\begin{aligned} F_{\infty}(i) &= \mathbb{E}_{\Psi^*} \left[ \sum_{t=0}^{s-1} c(X_t, u_t) + F_{\infty}(X_s) | X_0 = i \right] \\ &\geq \mathbb{E}_{\Psi^*} \left[ \sum_{t=0}^{s-1} c(X_t, u_t) | X_0 = i \right]. \end{aligned}$$



Letting  $s \rightarrow \infty$  gives  $F_\infty(i) \geq F^{\Psi^*}(i) \geq F(i)$ , which leads to the required result  $F_\infty(i) = F(i)$  if either  $D$  or  $P$  holds, or  $N$  and  $A$  hold.  $\square$

**Example 42 (Factory revisited).** The  $s$ -horizon values are given by

$$\begin{aligned} F_s(0) &= \inf\{F_s^M(0), F_s^A(0)\} \\ &= \inf \left\{ \begin{array}{l} 8c + \beta [pF_{s-1}^M(0) + (1-p)F_{s-1}^M(1)], \\ 8c + \beta [qF_{s-1}^A(0) + (1-q)F_{s-1}^A(1)] \end{array} \right\} \\ F_s(1) &= \inf\{F_s^M(1), F_s^A(1)\} \\ &= \inf \left\{ \begin{array}{l} c + \beta [pF_{s-1}^M(0) + (1-p)F_{s-1}^M(1)], \\ \beta [qF_{s-1}^A(0) + (1-q)F_{s-1}^A(1)] \end{array} \right\} \end{aligned}$$

with  $F_0(0) = \inf[0, 0]$  and  $F_0(1) = \inf[0, 0]$ .

We will use the MATLAB file **FactoryExample.m** on MyUni to investigate this algorithm. Effectively we have a function in the optimality equation which we can use to set up some recurrence relations that we can iterate on computer and investigate solutions for various policies. The function is monotonically increasing and because of the discounting parameter  $\beta$ , we would hope that the function will converge to the same constant values we calculated analytically. We can thus check our iterations using the analytic solutions for the fixed policies of running in either automatic mode or manual mode, as well as investigating other mixed policies.

### §3.7 NEGATIVE PROGRAMMING & OPTIMAL STOPPING

Recall that  $N$ , denoting Negative Programming, is about maximising non-positive rewards,  $r(i, u) \leq 0$ , or minimising non-negative costs,  $c(i, u) \geq 0$ , and that

1. a **Markov policy** is a policy that specifies the control at time  $t$  to be simply a function of the state and time. Earlier we used the notation

$$u_t(i) = f_t(i)$$

to specify the action at time  $t$ ; convenient notation for a Markov policy.

2. If in addition the policy does not depend on time in its choice of action, then it is said to be a **stationary Markov policy**. We write  $\Psi = (f, f, \dots) = f$ , where the stationary action  $u = f(i)$ .

### CHARACTERISATION OF THE OPTIMAL POLICY

The following theorem gives a necessary and sufficient condition for a stationary policy to be optimal: namely, it must choose the optimal  $u$  on the RHS of the optimality equation. (Note that in this theorem we are requiring that the infimum over  $u$  is attained as the minimum over  $u$ , as would be the case if we make the finite action assumption, A.)

**Theorem 3.6.** *Suppose D or N holds. Suppose  $\Psi = f$  is the stationary Markov policy such that*

$$f(i) = \operatorname{argmin}_u \{ \mathbb{E}[c(i, u)] + \beta \mathbb{E}[F(X_1) | X_0 = i, u_0 = u] \}.$$

*Then  $F^\Psi(i) = F(i)$ , and  $\Psi$  is optimal. That is,  $u = f(i)$  is the stationary value of  $u$  that minimises the RHS of the optimality equation.*

*Proof.* With  $\Psi = f$ , by substituting the optimality equation into itself and using the fact that  $\Psi$  specifies the minimising action at each stage,

$$F(i) = \mathbb{E}_\Psi \left[ \sum_{t=0}^{s-1} \beta^t c(X_t, u_t) | X_0 = i \right] + \beta^s \mathbb{E}_\Psi [F(X_s) | X_0 = i]. \quad (3.32)$$

We now show that  $F(i) \geq F^\Psi(i)$ :

- In the case  $N$ , we can drop the final term on the RHS of (3.32) (because it is non-negative) and then let  $s \rightarrow \infty$ .

- In the case  $D$ , we can let  $s \rightarrow \infty$  directly, observing that this term tends to 0.

On the other hand, by definition,  $F(i) \leq F^\Psi(i)$ , so we arrive at equality.  $\square$

**Corollary 3.1.** *If  $D$  holds or  $N$  holds under Assumption A, then an optimal policy exists.*

Neither Theorem 3.6 nor Corollary 3.1 are true for  $P$  (see Example 40).

### OPTIMAL STOPPING OVER A FINITE HORIZON

One way that the total expected cost can be finite is if it is possible to enter a state from which no further costs are incurred.

Suppose  $u$  can take on two possible values:  $u = 0$  (stop) and  $u = 1$  (continue). Suppose there is a termination (absorbing) state, say 0, that is entered upon choosing the stopping action. Once this state is entered the system stays in that state and no further cost is incurred. We let  $c(x, 0) = k(x)$  (stopping cost) and  $c(x, 1) = c(x)$  (continuation cost).

Suppose that  $F_s(x)$  denotes the minimum total cost when we are constrained to stop within the next  $s$  steps. This gives a finite-horizon problem with dynamic programming equation

$$F_s(x) = \min \{k(x), c(x) + \mathbb{E}[F_{s-1}(X_1)|X_0 = x, u_0 = 1]\} \quad (3.33)$$

with  $F_0(x) = k(x)$ ,  $c(0) = 0$ .

Consider now the set of states in which it is at least as good to stop now as to continue one more step, and then stop:

$$\mathcal{L} = \{x : k(x) \leq c(x) + \mathbb{E}[k(X_1)|X_0 = x, u_0 = 1]\}.$$

It is obviously not optimal to stop if  $x \notin \mathcal{L}$ . If  $\mathcal{L}$  is closed then the following theorem gives us the form of the optimal policies for all finite-horizons.

**Theorem 3.7** (OSLA). *Suppose  $\mathcal{L}$  is closed (so that once the system enters  $\mathcal{L}$  it remains in  $\mathcal{L}$ ). Then an optimal policy for all finite horizons is: stop iff  $x \in \mathcal{L}$ .*

*Proof.* This proof is by induction. If the horizon is  $s = 1$ , then obviously it is optimal to stop only if  $x \in \mathcal{S}$ . Suppose the theorem is true for a horizon of  $s - 1$ . As above, if  $x \notin \mathcal{S}$ , then it is better to continue for one more step and stop rather than stop in state  $x$ . If  $x \in \mathcal{S}$ , then the fact that  $\mathcal{S}$  is closed implies  $x_1 \in \mathcal{S}$  and so  $F_{s-1}(x_1) = k(x_1)$ . But then (3.33) gives  $F_s(x) = k(x)$ . So we should stop if  $x \in \mathcal{S}$ .  $\square$

The optimal policy is known as the one-step look-ahead rule (OSLA rule).

**Example 43.** Optimal parking. Recall Problem 2 from Lecture 1. Let  $x = 1$  denote the car park is available  $x$  spaces from the entrance, and  $x = 0$  denote that it is full. We have

$$\begin{aligned} F_s(0) &= (1 - p)F_{s-1}(0) + pF_{s-1}(1), \\ F_s(1) &= \min\{s, (1 - p)F_{s-1}(0) + pF_{s-1}(1)\} \end{aligned}$$

where  $F_0(0) = D$  and  $F_0(1) = 0$ .

Now, consider the set  $\mathcal{L}$  for this problem. We have

$$\mathcal{L} = \{s : s \leq e(s - 1)\}$$

where  $e(s)$  is the expected cost if we take the first available parking space that is  $s$  spaces or closer to the entrance. (Note that “taking the first available parking space is equivalent to stopping at the first opportunity.”) Then,

$$e(s) = ps + (1 - p)e(s - 1),$$

with  $e(0) = (1 - p)D$ . Solving, we have

$$e(s) = \frac{-(1 - p)}{p} + s + \left(D + \frac{1}{p}\right)(1 - p)^{s+1}, \quad s = 0, 1, \dots$$

Simplifying, we have

$$\mathcal{L} = \{s : (Dp + 1)(1 - p)^s \geq 1\}.$$

This set is closed, as  $s$  is decreasing, and hence Theorem 3.7 holds and  $\mathcal{L}$  characterises the optimal policy.

### OPTIMAL STOPPING OVER AN INFINITE HORIZON

Let  $F_s(x)$  be the infimal cost given that we are required to stop by the  $s$ th step. Let  $F(x)$  be the equivalent where the only requirement is that we eventually stop. Since less cost can be incurred if we are allowed more time in which to stop, we have

$$F_s(x) \geq F_{s+1}(x) \geq F(x).$$

Hence, by monotone convergence  $F_s(x)$  tends to a limit, say  $F_\infty(x)$ , and  $F_\infty(x) \geq F(x)$ .

The following Lemma gives conditions under which the infimal finite-horizon cost does converge to the infimal infinite-horizon cost.

**Lemma 3.1.** *Suppose all costs are bounded as follows.*

$$K = \sup_x k(x) < \infty, \text{ and } C = \inf_x c(x) > 0. \quad (3.34)$$

*Then  $F_s(x) \rightarrow F(x)$  as  $s \rightarrow \infty$ .*

*Proof.* Suppose  $\pi$  is an optimal policy for the infinite horizon problem and stops at the random time  $\tau$ . Then its cost is at least  $(s + 1)C \Pr(\tau > s)$ . However, since it would be possible to stop at time 0 the cost is also no more than  $K$ , so

$$(s + 1)C \Pr(\tau > s) \leq F(x) \leq K.$$

In the  $s$ -horizon problem we could follow  $\pi$ , but stop at time  $s$  if  $\tau > s$ . This implies

$$F(x) \leq F_s(x) \leq_* F(x) + K \Pr(\tau > s) \leq F(x) + \frac{K^2}{(s+1)C}.$$

By letting  $s \rightarrow \infty$ , we have  $F_\infty(x) = F(x)$ .

Note that the inequality  $\leq_*$  comes considering  $F_s(x)$  conditioning on whether  $\tau \leq s$  or  $\tau > s$ .  $\square$

Note that the problem posed here is identical to one in which we pay  $K$  at the start and receive a terminal reward  $r(x) = K - k(x)$ .

**Theorem 3.8.** *Suppose  $\mathcal{L}$  is closed and 3.34 hold. Then an optimal policy for the infinite horizon is: stop iff  $x \in \mathcal{L}$ .*

*Proof.* By Theorem 3.7 we have for all finite  $s$ ,

$$F_s(x) \begin{cases} = k(x), & x \in \mathcal{L} \\ < k(x), & x \notin \mathcal{L} \end{cases}$$

Lemma 3.1 gives  $F(x) = F_\infty(x)$ .  $\square$

# Bibliography

- [GS01][1] GRIMMETT, G., AND STIRZAKER, D. (2001) Probability and Random Processes. Oxford.
- [KM11][2] KALL, P., AND MAYER, J. (2011) Stochastic Linear Programming. Springer.
- [KW94][3] KALL, P., AND WALLACE, S. (1994) Stochastic Programming. John Wiley & Sons.
- [KC13][4] KROESE, D., AND CHAN, J. (2013) Statistical Modelling and Computation: An Inclusive Approach to Statistics. Springer.
- [Put][5] PUTERMAN, M. *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. John Wiley & Sons, 2005.
- [Sen][6] SENNOTT, L. *Stochastic Dynamics Programming and the Control of Queueing Systems*. John Wiley & Sons, 1998.
- [T16][7] TAHA, H. (2016) Operations Research: An Introduction. Pearson.