# Large Deviation Theory

J.M. Swart

December 21, 2016

## Preface

The earliest origins of large deviation theory lie in the work of Boltzmann on entropy in the 1870ies and Cramér's theorem from 1938 [Cra38]. A unifying mathematical formalism was only developed starting with Varadhan's definition of a 'large deviation principle' (LDP) in 1966 [Var66].

Basically, large deviation theory centers around the observation that suitable functions F of large numbers of i.i.d. random variables  $(X_1, \ldots, X_n)$  often have the property that

$$\mathbb{P}[F(X_1, \dots, X_n) \in \mathrm{d}x] \sim e^{-s_n I(x)} \quad \text{as } n \to \infty, \qquad (\mathrm{LDP})$$

where  $s_n$  are real contants such that  $\lim_{n\to\infty} s_n = \infty$  (in most cases simply  $s_n = n$ ). In words, (LDP) says that the probability that  $F(X_1, \ldots, X_n)$  takes values near a point x decays exponentially fast, with speed  $s_n$ , and rate function I.

Large deviation theory has two different aspects. On the one hand, there is the question of how to formalize the intuitive formula (LDP). This leads to the already mentioned definition of 'large deviation principles' and involves quite a bit of measure theory and real analysis. The most important basic results of the abstract theory were proved more or less between 1966 and 1991, when O'Brian en Verwaat [OV91] and Puhalskii [Puk91] proved that exponential tightness implies a subsequential LDP. The abstract theory of large deviation principles plays more or less the same role as measure theory in (usual) probability theory.

On the other hand, there is a much richer and much more important side of large deviation theory, which tries to identify rate functions I for various functions F of independent random variables, and study their properties. This part of the theory is as rich as the branch of probability theory that tries to prove limit theorems for functions of large numbers of random variables, and has many relations to the latter.

There exist a number of good books on large deviation theory. The oldest book that I am aware of is the one by Ellis [Ell85], which is still useful for applications of large deviation theory in statistical mechanics and gives a good intuitive feeling for the theory, but lacks some of the standard results.

The classical books on the topic are the ones of Deuschel and Stroock [DS89] and especially Dembo and Zeitouni [DZ98], the latter originally published in 1993. While these are very thorough introductions to the field, they can at places be a bit hard to read due to the technicalities involved. Also, both books came a bit

too early to pick the full fruit of the development of the abstract theory.

A very pleasant book to read as a first introduction to the field is the book by Den Hollander [Hol00], which avoids many of the technicalities in favour of a clear exposition of the intuitive ideas and a rich choice of applications. A disadvantage of this book is that it gives little attention to the abstract theory, which means many results are not proved in their strongest form.

Two modern books on the topic, which each try to stress certain aspects of the theory, are the books by Dupuis and Ellis [DE97] and Puhalskii [Puh01]. These books are very strong on the abstract theory, but, unfortunately, they indulge rather heavily in the introduction of their own terminology and formalism (for example, in [DE97], replacing the large deviation principle by the almost equivalent 'Laplace principle') which makes them somewhat inaccessible, unless read from the beginning to the end.

A difficulty encountered by everyone who tries to teach large deviation theory is that in order to do it properly, one first needs quite a bit of abstract theory, which however is intuitively hard to grasp unless one has seen at least a few examples. I have tried to remedy this by first stating, without proof, a number of motivating examples. In the proofs, I have tried to make optimal use of some of the more modern abstract theory, while sticking with the classical terminology and formulations as much as possible.

# Contents

0	Som	e motivating examples	<b>7</b>	
	0.1	Cramér's theorem	7	
	0.2	Moderate deviations	0	
	0.3	Relative entropy	.1	
	0.4	Non-exit probabilities	4	
	0.5	Outlook 1	.6	
1	Abstract theory 17			
	1.1	Weak convergence on Polish spaces	7	
	1.2	Large deviation principles	22	
	1.3	Varadhan's lemma	28	
	1.4	The contraction principle	30	
	1.5	Exponential tilts	31	
	1.6	Robustness 3	33	
	1.7	Tightness	35	
	1.8	LDP's on compact spaces	<b>37</b>	
	1.9	Exponential tightness	12	
	1.10	Applications of exponential tightness	17	
<b>2</b>	Sun	ns of i.i.d. random variables 5	5	
	2.1	The Legendre transform	55	
	2.2	Cramér's theorem	52	
	2.3	The multi-dimensional Legendre transform	55	
	2.4	Relative entropy	'2	
	2.5	Cramér's theorem in more dimensions	32	
	2.6	Sanov's theorem	38	
3	Markov chains 91			
	3.1	Basic notions	)1	
	3.2	A LDP for Markov chains	)3	
	3.3	The empirical process	)4	
	3.4	Perron-Frobenius eigenvalues	1	
	3.5	Continuous time	7	
	3.6	Excercises	28	

CONTENTS

## Chapter 0

## Some motivating examples

## 0.1 Cramér's theorem

Let  $(X_k)_{k\geq 1}$  be a sequence of i.i.d. absolutely integrable (i.e.,  $\mathbb{E}[|X_1|] < \infty$ ) real random variables with mean  $\rho := \mathbb{E}[X_1]$ , and let

$$T_n := \frac{1}{n} \sum_{k=1}^n X_k \qquad (n \ge 1).$$

be their empirical avarages. Then the weak law of large numbers states that

$$\mathbb{P}[|T_n - \rho| \ge \varepsilon] \underset{n \to \infty}{\longrightarrow} 0 \qquad (\varepsilon > 0).$$

In 1938, the Swedish statistician and probabilist Harald Cramér [Cra38] studied the question how fast this probability tends to zero. For laws with sufficiently light tails (as stated in the condition (0.1) below), he arrived at the following conclusion.

### Theorem 0.1 (Cramér's theorem) Assume that

$$Z(\lambda) := \mathbb{E}[e^{\lambda X_1}] < \infty \qquad (\lambda \in \mathbb{R}).$$

$$(0.1)$$

Then

(i) 
$$\lim_{n \to \infty} \frac{1}{n} \log \mathbb{P}[T_n \ge y] = -I(y) \qquad (y > \rho),$$
  
(ii) 
$$\lim_{n \to \infty} \frac{1}{n} \log \mathbb{P}[T_n \le y] = -I(y) \qquad (y < \rho),$$
  
(0.2)

where I is defined by

$$I(y) := \sup_{\lambda \in \mathbb{R}} \left[ y\lambda - \log Z(\lambda) \right] \qquad (y \in \mathbb{R}).$$
(0.3)

The function Z in (0.1) is called the moment generating function or cumulant generating function, and its logarithm is consequently called the *logarithmic moment* generating function (or *logarithmic cumulant generating function* of the law of  $X_1$ . In the context of large deviation theory,  $\log Z(\lambda)$  is also called the *free energy* function, see [Ell85, Section II.4].

The function I defined in (0.3) is called the *rate function*. In order to see what Cramér's theorem tells us exactly, we need to know some elementary properties of this function. Note that (0.1) implies that  $\mathbb{E}[|X_1|^2] < \infty$ . To avoid trivial cases, we assume that the  $X_k$  are not a.s. constant, i.e.,  $\operatorname{Var}(X_1) > 0$ .

Below, int(A) denotes the interior of a set A, i.e., the largest open set contained in A. We recall that for any finite measure  $\mu$  on  $\mathbb{R}$ ,  $support(\mu)$  is the smallest closed set such that  $\mu$  is concentrated on  $support(\mu)$ .

**Lemma 0.2 (Properties of the rate function)** Let  $\mu$  be the law of  $X_1$ , let  $\rho := \langle \mu \rangle$  and  $\sigma^2 := \operatorname{Var}(\mu)$  denote its mean and variance, and assume that  $\sigma > 0$ . Let  $y_- := \inf(\operatorname{support}(\mu)), y_+ := \sup(\operatorname{support}(\mu))$ . Let I be the function defined in (0.3) and set

$$\mathcal{D}_I := \{ y \in \mathbb{R} : I(y) < \infty \}$$
 and  $\mathcal{U}_I := \operatorname{int}(\mathcal{D}_I).$ 

Then:

- (i) I is convex.
- (ii) I is lower semi-continuous.
- (iii)  $0 \leq I(y) \leq \infty$  for all  $y \in \mathbb{R}$ .
- (iv) I(y) = 0 if and only if  $y = \rho$ .
- (v)  $\mathcal{U}_I = (y_-, y_+).$
- (vi) I is infinitely differentiable on  $\mathcal{U}_I$ .
- (vii)  $\lim_{y \downarrow y_{-}} I'(y) = -\infty$  and  $\lim_{y \uparrow y_{+}} I'(y) = \infty$ .
- (viii) I'' > 0 on  $\mathcal{U}_I$  and  $I''(\rho) = 1/\sigma^2$ .



Figure 1: A typical example of a rate function.

(ix) If 
$$-\infty < y_-$$
, then  $I(y_-) = -\log \mu(\{y_-\})$ , and  
if  $y_+ < \infty$ , then  $I(y_+) = -\log \mu(\{y_+\})$ .

See Figure 1 for a picture. Here, if E is any metric space (e.g.  $E = \mathbb{R}$ ), then we say that a function  $f: E \to [-\infty, \infty]$  is *lower semi-continuous* if one (and hence both) of the following equivalent conditions are satisfied:

- (i)  $\liminf_{n\to\infty} f(x_n) \ge f(x)$  whenever  $x_n \to x$ .
- (ii) For each  $-\infty \le a \le \infty$ , the *level set*  $\{x \in E : I(x) \le a\}$  is a closed subset of E.

In view of Lemma 0.2, Theorem 0.1 tells us that the probability that the empirical average  $T_n$  deviates by any given constant from its mean decays exponentially fast in n. More precisely, formula (0.2) (i) says that

$$\mathbb{P}[T_n \ge y] = e^{-nI(y) + o(n)} \quad \text{as} \quad n \to \infty \quad (y > \rho),$$

were, as usual, o(n) denotes any function such that

$$o(n)/n \to 0$$
 as  $n \to \infty$ .

Note that formulas (0.2) (i) and (ii) only consider one-sided deviations of  $T_n$  from its mean  $\rho$ . Nevertheless, the limiting behavior of two-sided deviations can easily be derived from Theorem 0.1. Indeed, for any  $y_- < \rho < y_+$ ,

$$\mathbb{P}[T_n \le y_- \text{ or } T_n \ge y_+] = e^{-nI(y_-) + o(n)} + e^{-nI(y_+) + o(n)}$$
$$= e^{-n\min\{I(y_-), I(y_+)\} + o(n)} \quad \text{as} \quad n \to \infty.$$

In particular,

$$\lim_{n \to \infty} \frac{1}{n} \log \mathbb{P}[|T_n - \rho| \ge \varepsilon] = \min\{I(\rho - \varepsilon), I(\rho + \varepsilon)\} \qquad (\varepsilon > 0)$$

**Exercise 0.3** Use Theorem 0.1 and Lemma 0.2 to deduce that, under the assumptions of Theorem 0.1,

$$\lim_{n \to \infty} \frac{1}{n} \log \mathbb{P}[T_n > y] = -I_{\rm up}(y) \qquad (y \ge \rho),$$

where  $I_{up}$  is the upper semi-continuous modification of I, i.e.,  $I_{up}(y) = I(y)$  for  $y \neq y_{-}, y_{+}$  and  $I_{up}(y_{-}) = I_{up}(y_{+}) := \infty$ .

## 0.2 Moderate deviations

As in the previous section, let  $(X_k)_{k\geq 1}$  be a sequence of i.i.d. absolutely integrable real random variables with mean  $\rho := \mathbb{E}[|X_1|]$  and assume that (0.1) holds. Let

$$S_n := \sum_{k=1}^n X_k \qquad (n \ge 1).$$

be the *partial sums* of the first n random variables. Then Theorem 0.1 says that

$$\mathbb{P}[S_n - \rho n \ge yn] = e^{-nI(\rho + y) + o(n)} \quad \text{as } n \to \infty \qquad (y > 0).$$

On the other hand, by the central limit theorem, we know that

$$\mathbb{P}[S_n - \rho n \ge y\sqrt{n}] \xrightarrow[n \to \infty]{} \Phi(y/\sigma) \qquad (y \in \mathbb{R}),$$

where  $\Phi$  is the distribution function of the standard normal distribution and

$$\sigma^2 = \operatorname{Var}(X_1),$$

which we assume to be positive. One may wonder what happens at in-between scales, i.e., how does  $\mathbb{P}[S_n - \rho n \geq y_n]$  decay to zero if  $\sqrt{n} \ll y_n \ll n$ ? This is the question of *moderate deviations*. We will only consider the case  $y_n = yn^{\alpha}$  with  $\frac{1}{2} < \alpha < 1$ , even though other timescales (for example in connection with the law of the iterated logarithm) are also interesting.

**Theorem 0.4 (Moderate deviations)** Let  $(X_k)_{k\geq 1}$  be a sequence of i.i.d. absolutely integrable real random variables with mean  $\rho := \mathbb{E}[|X_1|]$ , variance  $\sigma^2 = \operatorname{Var}(X_1) > 0$ , and  $\mathbb{E}[e^{\lambda X_1}] < \infty$  ( $\lambda \in \mathbb{R}$ ). Then

$$\lim_{n \to \infty} \frac{1}{n^{2\alpha - 1}} \log \mathbb{P}[S_n - \rho n \ge y n^{\alpha}] = -\frac{1}{2\sigma^2} y^2 \qquad (y > 0, \ \frac{1}{2} < \alpha < 1).$$
(0.4)

**Remark** Setting  $y_n := yn^{\alpha-1}$  and naively applying Cramér's theorem, pretending that  $y_n$  is a constant, using Lemma 0.2 (viii), we obtain

$$\log \mathbb{P}[S_n - \rho n \ge y n^{\alpha}] = \log \mathbb{P}[S_n - \rho n \ge y_n n]$$
  
$$\approx -nI(y_n) \approx -n\frac{1}{2\sigma^2}y_n^2 = -\frac{1}{2\sigma^2}y^2 n^{2\alpha - 1}.$$

Dividing both sides of this equation by  $n^{2\alpha-1}$  yields formula (0.4) (although this derivation is not correct). There does not seem to be a good basic reference for moderate deviations. Some more or less helpful references are [DB81, Led92, Aco02, EL03].

### 0.3 Relative entropy

Imagine that we throw a dice n times, and keep record of how often each of the possible outcomes  $1, \ldots, 6$  comes up. Let  $N_n(x)$  be the number of times outcome x has turned up in the first n throws, let  $M_n(x) := N_n(x)/x$  be the relative frequency of x, and set

$$\Delta_n := \max_{1 \le x \le 6} M_n(x) - \min_{1 \le x \le 6} M_n(x).$$

By the strong law of large numbers, we know that  $M_n(x) \to 1/6$  a.s. as  $n \to \infty$  for each  $x \in \{1, \ldots, 6\}$ , and therefore  $\mathbb{P}[\Delta_n \ge \varepsilon] \to 0$  as  $n \to \infty$  for each  $\varepsilon > 0$ . It turns out that this convergence happens exponentially fast.

**Proposition 0.5 (Deviations from uniformity)** There exists a continuous, strictly increasing function  $I : [0, 1] \to \mathbb{R}$  with I(0) = 0 and  $I(1) = \log 6$ , such that

$$\lim_{n \to \infty} \frac{1}{n} \log \mathbb{P}[\Delta_n \ge \varepsilon] = -I(\varepsilon) \qquad (0 \le \varepsilon \le 1).$$
(0.5)

Proposition 0.5 follows from a more general result that was already discovered by the physicist Boltzmann in 1877. A much more general version of this result for random variables that do not need to take values in a finite space was proved by the Russian mathematician Sanov [San61]. We will restrict ourselves to finite state spaces for the moment. To state the theorem, we first need a few definitions.

Let S be a finite set and let  $\mathcal{M}_1(S)$  be the set of all probability measures on S. Since S is finite, we may identify  $\mathcal{M}_1(S)$  with the set

$$\mathcal{M}_1(S) := \big\{ \pi \in \mathbb{R}^S : \pi(x) \ge 0 \ \forall x \in S, \ \sum_{x \in S} \pi(1) = 1 \big\},\$$

where  $\mathbb{R}^S$  denotes the space of all functions  $\pi : S \to \mathbb{R}$ . Note that  $\mathcal{M}_1(S)$  is a compact, convex subset of the (|S|-1)-dimensional space  $\{\pi \in \mathbb{R}^S : \sum_{x \in S} \pi(1) = 1\}$ .

Let  $\mu, \nu \in \mathcal{M}_1(S)$  and assume that  $\mu(x) > 0$  for all  $x \in S$ . Then we define the *relative entropy* of  $\nu$  with respect to  $\mu$  by

$$H(\nu|\mu) := \sum_{x \in S} \nu(x) \log \frac{\nu(x)}{\mu(x)} = \sum_{x \in S} \mu(x) \frac{\nu(x)}{\mu(x)} \log \frac{\nu(x)}{\mu(x)},$$

where we use the conventions that  $\log(0) := -\infty$  and  $0 \cdot \infty := 0$ . Note that since  $\lim_{z \downarrow 0} z \log z = 0$ , the second formula shows that  $H(\nu|\mu)$  is continuous in  $\nu$ . The function  $H(\nu|\mu)$  is also known as the Kullback-Leibler distance or divergence.

**Lemma 0.6 (Properties of the relative entropy)** Assume that  $\mu \in \mathcal{M}_1(S)$ and assume that  $\mu(x) > 0$  for all  $x \in S$ . Then the function  $\nu \mapsto H(\nu|\mu)$  has the following properties.

(i)  $0 \leq H(\nu|\mu) < \infty$  for all  $\nu \in \mathcal{M}_1(S)$ .

- (ii)  $H(\mu|\mu) = 0.$
- (iii)  $H(\nu|\mu) > 0$  for all  $\nu \neq \mu$ .
- (iv)  $\nu \mapsto H(\nu|\mu)$  is convex and continuous on  $\mathcal{M}_1(S)$ .
- (v)  $\nu \mapsto H(\nu|\mu)$  is infinitely differentiable on the interior of  $\mathcal{M}_1(S)$ .

Assume that  $\mu \in \mathcal{M}_1(S)$  satisfies  $\mu(x) > 0$  for all  $x \in S$  and let  $(X_k)_{k \ge 1}$  be an i.i.d. sequence with common law  $\mathbb{P}[X_1 = x] = \mu(x)$ . As in the example of the dice

#### 0.3. RELATIVE ENTROPY

throws, we let

$$M_n(x) := \frac{1}{n} \sum_{k=1}^n \mathbb{1}_{\{X_k = x\}} \qquad (x \in S, \ n \ge 1).$$
(0.6)

Note that  $M_n$  is a  $\mathcal{M}_1(S)$ -valued random variable. We call  $M_n$  the *empirical* distribution.

**Theorem 0.7 (Boltzmann-Sanov)** Let C be a closed subset of  $\mathcal{M}_1(S)$  such that C is the closure of its interior. Then

$$\lim_{n \to \infty} \frac{1}{n} \log \mathbb{P}[M_n \in C] = -\min_{\nu \in C} H(\nu|\mu).$$
(0.7)

Note that (0.7) says that

$$\mathbb{P}[M_n \in C] = e^{-nI_C + o(n)} \text{ as } n \to \infty \quad \text{where} \quad I_C = \min_{\nu \in C} H(\nu|\mu). \tag{0.8}$$

This is similar to what we have already seen in Cramér's theorem: if I is the rate function from Theorem 0.1, then  $I(y) = \min_{y' \in [y,\infty)} I(y')$  for  $y > \rho$  and  $I(y) = \min_{y' \in (-\infty,y]} I(y')$  for  $y < \rho$ . Likewise, as we have seen in (0.1), the probability that  $T_n \in (-\infty, y_-] \cup [y_+, \infty)$  decays exponentially with rate  $\min_{y' \in (-\infty, y_-] \cup [y_+, \infty)} I(y')$ .

The proof of Theorem 0.7 will be delayed till later, but we will show here how Theorem 0.7 implies Proposition 0.5.

**Proof of Proposition 0.5** We set  $S := \{1, \ldots, 6\}$ ,  $\mu(x) := 1/6$  for all  $x \in S$ , and apply Theorem 0.7. For each  $0 \le \varepsilon < 1$ , the set

$$C_{\varepsilon} := \left\{ \nu \in \mathcal{M}_1(S) : \max_{x \in S} \nu(x) - \min_{x \in S} \nu(x) \ge \varepsilon \right\}$$

is a closed subset of  $\mathcal{M}_1(S)$  that is the closure of its interior. (Note that the last statement fails for  $\varepsilon = 1$ .) Therefore, Theorem 0.7 implies that

$$\lim_{n \to \infty} \frac{1}{n} \log \mathbb{P}[\Delta_n \ge \varepsilon] = \lim_{n \to \infty} \frac{1}{n} \log \mathbb{P}[M_n \in C_\varepsilon] = -\min_{\nu \in C_\varepsilon} H(\nu|\mu) =: -I(\varepsilon). \quad (0.9)$$

The fact that I is continuous and satisfies I(0) = 0 follows easily from the properties of  $H(\nu, \mu)$  listed in Lemma 0.6. To see that I is strictly increasing, fix  $0 \leq \varepsilon_1 < \varepsilon_2 < 1$ . Since  $H(\cdot | \mu)$  is continuous and the  $C_{\varepsilon_2}$  are compact, we can find a  $\nu_*$  (not necessarily unique) such that  $H(\cdot | \mu)$  assumes its minimum over  $C_{\varepsilon_2}$  in  $\nu_*$ . Now by the fact that  $H(\cdot | \mu)$  is convex and assumes its unique minimum in  $\mu$ , we see that  $\nu' := \frac{\varepsilon_1}{\varepsilon_2}\nu_* + (1 - \frac{\varepsilon_1}{\varepsilon_2})\mu \in C_{\varepsilon_1}$  and therefore  $I(\varepsilon_1) \leq H(\nu'|\mu) < H(\nu_*|\mu) = I(\varepsilon_2).$ 

Finally, by the continuity of  $H(\cdot | \mu)$ , we see that

$$I(\varepsilon) \uparrow \min_{\nu \in C_1} H(\nu|\mu) = H(\delta_1|\mu) = \log 6$$
 as  $\varepsilon \uparrow 1$ .

To see that (0.5) also holds for  $\varepsilon = 1$  (which does not follow directly from Theorem 0.7 since  $C_1$  is not the closure of its interior), it suffices to note that  $\mathbb{P}[\Delta_n = 1] = (\frac{1}{6})^{n-1}$ .

**Remark 1** It is quite tricky to calculate the function I from Proposition 0.5 explicitly. For  $\varepsilon$  sufficiently small, it seems that the minimizers of the entropy  $H(\cdot | \mu)$  on the set  $C_{\varepsilon}$  are (up to permutations of the coordinates) of the form  $\nu(1) = \frac{1}{6} - \frac{1}{2}\varepsilon$ ,  $\nu(2) = \frac{1}{6} + \frac{1}{2}\varepsilon$ , and  $\nu(3), \ldots, \nu(6) = \frac{1}{6}$ . For  $\varepsilon > \frac{1}{3}$ , this solution is of course no longer well-defined and the minimizer must look differently.

**Remark 2** I do not know whether the function I is convex.

### 0.4 Non-exit probabilities

In this section we move away from the i.i.d. setting and formulate a large deviation result for Markov processes. To keep the technicalities to a minimum, we restrict ourselves to Markov processes with a finite state space. We recall that a continuous-time, time-homogeneous Markov process  $X = (X_t)_{t\geq 0}$  taking value in a finite set S is uniquely characterized (in law) by its initial law  $\mu(x) := \mathbb{P}[X_0 = x]$ and its transition probabilities  $P_t(x, y)$ . Indeed, X has piecewise constant, rightcontinuous sample paths and its finite-dimensional distributions are characterized by

$$\mathbb{P}[X_{t_1} = x_1, \dots, X_{t_n} = x_n] = \sum_{x_0} \mu(x_0) P_{t_1}(x_0, x_1) P_{t_2 - t_1}(x_1, x_2) \cdots P_{t_n - t_{n-1}}(x_n, x_n)$$

 $(t_1 < \cdots < t_n, x_1, \ldots, x_n \in S)$ . The transition probabilities are continuous in t, have  $P_0(x, y) = 1_{\{x=y\}}$  and satisfy the Chapman-Kolmogorov equation

$$\sum_{y} P_s(x, y) P_t(y, z) = P_{s+t}(x, z) \qquad (s, t \ge 0, \ x, z \in S).$$

### 0.4. NON-EXIT PROBABILITIES

As a result, they define a semigroup  $(P_t)_{t>0}$  of linear operators  $P_t : \mathbb{R}^S \to \mathbb{R}^S$  by

$$P_t f(x) := \sum_y P_t(x, y) f(y) = \mathbb{E}^x [f(X_t)],$$

where  $\mathbb{E}^x$  denotes expectation with respect to the law  $\mathbb{P}^x$  of the Markov process with initial state  $X_0 = x$ . One has

$$P_t = e^{Gt} = \sum_{n=0}^{\infty} \frac{1}{n!} G^n t^n,$$

where  $G : \mathbb{R}^S \to \mathbb{R}^S$ , called the *generator* of the semigroup  $(P_t)_{t \ge 0}$ , is an operator of the form

$$Gf(x) = \sum_{y: y \neq x} r(x, y) \big( f(y) - f(x) \big) \qquad (f \in \mathbb{R}^S, \ x \in S),$$

where r(x, y)  $(x, y \in S, x \neq y)$  are nonnegative contants. We call r(x, y) the rate of jumps from x to y. Indeed, since  $P_t = 1 + tG + O(t^2)$  as  $t \to 0$ , we have that

$$\mathbb{P}^{x}[X_{t} = y] = \begin{cases} tr(x, y) + O(t^{2}) & \text{if } x \neq y, \\ 1 - t \sum_{z: z \neq x} r(x, z) + O(t^{2}) & \text{if } x = y. \end{cases}$$

Let  $U \subset S$  be some strict subset of S and assume that  $X_0 \in U$  a.s. We will be interested in the probability that  $X_t$  stays in U for a long time. Let us say that the transition rates r(x, y) are *irreducible* on U if for each  $x, z \in U$  we can find  $y_0, \ldots, y_n$  such that  $y_0 = x$ ,  $y_n = z$ , and  $r(y_{k-1}, y_k) > 0$  for each  $k = 1, \ldots, n$ . Note that this says that it is possible for the Markov process to go from any point in Uto any other point in U without leaving U.

**Theorem 0.8 (Non-exit probability)** Let X be a Markov process with finite state space S, transition rates r(x, y)  $(x, y \in S, x \neq y)$ , and generator G. Let  $U \subset S$  and assume that the transition rates are irreducible on U. Then there exists a function f, unique up to a multiplicative constant, and a constant  $\lambda \geq 0$ , such that

(i) 
$$f > 0$$
 on  $U$ ,  
(ii)  $f = 0$  on  $S \setminus U$ ,  
(iii)  $Gf(x) = -\lambda f(x)$   $(x \in U)$ .

Moreover, the process X started in any initial law such that  $X_0 \in U$  a.s. satisfies

$$\lim_{t \to \infty} \frac{1}{t} \log \mathbb{P} \big[ X_s \in U \ \forall 0 \le s \le t \big] = -\lambda.$$
(0.10)

## 0.5 Outlook

Our aim will be to prove Theorems 0.1, 0.4, 0.7 and 0.8, as well as similar and more general results in a *unified framework*. Therefore, in the next chapter, we will give a formal definition of when a sequence of probability measures satisfies a *large deviation principle* with a given *rate function*. This will allow us to formulate our theorems in a unified framework that is moreover powerful enough to deal with generalizations such as a multidimensional version of Theorem 0.1 or a generalization of Theorem 0.7 to continuous spaces. We will see that large deviation principles satisfy a number of abstract principles such as the *contraction principle* which we have already used when we derived Proposition 0.5 from Theorem 0.7. Once we have set up the general framework in Chapter 1, in the following chapters, we set out to prove Theorems 0.1, 0.7, and 0.8, as well as similar and more general results,<sup>1</sup> and show how these are related.

<sup>&</sup>lt;sup>1</sup>Unfortunately, we will not have time to prove Theorem 0.4.

## Chapter 1

## Abstract theory

### 1.1 Weak convergence on Polish spaces

Recall that a topological space is a set E equipped with a collection  $\mathcal{O}$  of subsets of E that are called *open* sets, such that

- (i) If  $(O_{\gamma})_{\gamma \in \Gamma}$  is any collection of (possibly uncountably many) sets  $O_{\gamma} \in \mathcal{O}$ , then  $\bigcup_{\gamma \in \Gamma} O_{\gamma} \in \mathcal{O}$ .
- (ii) If  $O_1, O_2 \in \mathcal{O}$ , then  $O_1 \cap O_2 \in \mathcal{O}$ .
- (iii)  $\emptyset, E \in \mathcal{O}$ .

Any such collection of sets is called a *topology*. It is fairly standard to also assume the *Hausdorff* property

(iv) For each  $x_1, x_2 \in E, x_1 \neq x_2 \exists O_1, O_2 \in \mathcal{O}$  s.t.  $O_1 \cap O_2 = \emptyset, x_1 \in O_1, x_2 \in O_2$ .

A sequence of points  $x_n \in E$  converges to a limit x in a given topology  $\mathcal{O}$  if for each  $O \in \mathcal{O}$  such that  $x \in O$  there is an n such that  $x_m \in O$  for all  $m \geq n$ . (If the topology is Hausdorff, then such a limit is unique, i.e.,  $x_n \to x$  and  $x_n \to x'$ implies x = x'.) A set  $C \subset E$  is called *closed* if its complement is open.

Because of property (i) in the definition of a topology, for each  $A \subset E$ , the union of all open sets contained in A is itself an open set. We call this the *interior* of A, denoted as  $int(A) := \bigcup \{ O : U \subset A, O \text{ open} \}$ . Then clearly int(A) is the smallest open set contained in A. Similarly, by taking complements, for each set  $A \subset E$ there exists a smallest closed set containing A. We call this the *closure* of A, denoted as  $\overline{A} := \bigcap \{C : C \supset A, C \text{ closed}\}$ . A topological space is called *separable* if there exists a countable set  $D \subset E$  such that D is dense in E, where we say that a set  $D \subset E$  is *dense* if its closure is E, or equivalently, if every nonempty open subset of E has a nonempty intersection with D.

In particular, if d is a metric on E, and  $B_{\varepsilon}(x) := \{y \in E : d(x, y) < \varepsilon\}$ , then

$$\mathcal{O} := \left\{ O \subset E : \forall x \in O \; \exists \varepsilon > 0 \; \text{s.t.} \; B_{\varepsilon}(x) \subset O \right\}$$

defines a Hausdorff topology on E such that convergence  $x_n \to x$  in this topology is equivalent to  $d(x_n, x) \to 0$ . We say that the metric d generates the topology  $\mathcal{O}$ . If for a given topology  $\mathcal{O}$  there exists a metric d that generates  $\mathcal{O}$ , then we say that the topological space  $(E, \mathcal{O})$  is metrizable.

Recall that a sequence  $x_n$  in a metric space (E, d) is a *Cauchy sequence* if for all  $\varepsilon > 0$  there is an *n* such that  $d(x_k, x_l) \leq \varepsilon$  for all  $k, l \geq n$ . A metric space is *complete* if every Cauchy sequence converges.

A Polish space is a separable topological space  $(E, \mathcal{O})$  such that there exists a metric d on E with the property that (E, d) is complete and d generates  $\mathcal{O}$ . Warning: there may be many different metrics on E that generate the same topology. It may even happen that E is not complete in some of these metrics, and complete in others (in which case E is still Polish). Example:  $\mathbb{R}$  is separable and complete in the usual metric d(x, y) = |x - y|, and therefore  $\mathbb{R}$  is a Polish space. But  $d'(x, y) := |\arctan(x) - \arctan(y)|$  is another metric that generates the same topology, while  $(\mathbb{R}, d')$  is not complete. (Indeed, the completion of  $\mathbb{R}$  w.r.t. the metric d' is  $[-\infty, \infty]$ .)

On any Polish space  $(E, \mathcal{O})$  we let  $\mathcal{B}(E)$  denote the Borel- $\sigma$ -algebra, i.e., the smallest  $\sigma$ -algebra containing the open sets  $\mathcal{O}$ . We let  $B_b(E)$  and  $\mathcal{C}_b(E)$  denote the linear spaces of all bounded Borel-measurable and bounded continuous functions  $f: E \to \mathbb{R}$ , respectively. Then  $\mathcal{C}_b(E)$  is complete in the supermumnorm  $||f||_{\infty} :=$  $\sup_{x \in E} |f(x)|$ , i.e.,  $(\mathcal{C}_b(E), || \cdot ||_{\infty})$  is a Banach space [Dud02, Theorem 2.4.9]. We let  $\mathcal{M}(E)$  denote the space of all finite measures on  $(E, \mathcal{B}(E))$  and write  $\mathcal{M}_1(E)$ for the space of all probability measures. It is possible to equip  $\mathcal{M}(E)$  with a metric  $d_M$  such that [EK86, Theorem 3.1.7]

- (i)  $(\mathcal{M}(E), d_H)$  is a separable complete metric space.
- (ii)  $d_M(\mu_n, \mu) \to 0$  if and only if  $\int f d\mu_n \to \int f d\mu$  for all  $f \in \mathcal{C}_b(E)$ .

The precise choice of  $d_M$  (there are several canonical ways to define such a metric) is not important to us. We denote convergence in  $d_M$  as  $\mu_n \Rightarrow \mu$  and call the associated topology (which is uniquely determined by the requirements above) the topology of weak convergence. By property (i), the space  $\mathcal{M}(E)$  equipped with the topology of weak convergence is a Polish space.

**Proposition 1.1 (Weak convergence)** Let *E* be a Polish space and let  $\mu_n, \mu \in \mathcal{M}(E)$ . Then one has  $\mu_n \Rightarrow \mu$  if and only if the following two conditions are satisfied.

- (i)  $\limsup_{n \to \infty} \mu_n(C) \le \mu(C) \qquad \forall C \ closed,$
- (ii)  $\liminf_{n \to \infty} \mu_n(O) \ge \mu(O) \quad \forall O \text{ open.}$

If the  $\mu_n, \mu$  are probability measures, then it suffices to check either (i) or (ii).

Before we give the proof of Proposition 1.1, we need a few preliminaries. Recall the definition of lower semi-continuity from Section 0.1. Upper semi-continuity is defined similarly: a function  $f: E \to [-\infty, \infty)$  is upper semi-continuous if and only if -f is lower semi-continuous. We set  $\mathbb{R} := [-\infty, \infty]$  and define

$$\mathcal{U}(E) := \{ f : E \to \overline{\mathbb{R}} : f \text{ upper semi-continuous} \},\$$
$$\mathcal{U}_b(E) := \{ f \in \mathcal{U}(E) : \sup_{x \in E} |f(x)| < \infty \},\$$
$$\mathcal{U}_+(E) := \{ f \in \mathcal{U}(E) : f \ge 0 \},\$$

and  $\mathcal{U}_{b,+}(E) := \mathcal{U}_b(E) \cap \mathcal{U}_+(E)$ . We define  $\mathcal{L}(E), \mathcal{L}_b(E), \mathcal{L}_+(E), \mathcal{L}_{b,+}(E)$  respectively  $\mathcal{C}(E), \mathcal{C}_b(E), \mathcal{C}_+(E), \mathcal{C}_{b,+}(E)$  similarly, with upper semi-continuity replaced by lower semi-continuity and resp. continuity. We will also sometimes use the notation  $B(E), B_b(E), B_+(E), B_{b,+}(E)$  for the space of Borel measurable functions  $f: E \to \mathbb{R}$  and its subspaces of bounded, nonnegative, and bounded nonnegative functions, respectively.

Exercise 1.2 (Topologies of semi-continuity) Let  $\mathcal{O}_{up} := \{[-\infty, a) : -\infty < a \leq \infty\} \cup \{\emptyset, \overline{\mathbb{R}}\}$ . Show that  $\mathcal{O}_{up}$  is a topology on  $\overline{\mathbb{R}}$  (albeit a non-Hausdorff one!) and that a function  $f : E \to \overline{\mathbb{R}}$  is upper semi-continuous if and only if it is continuous with respect to the topology  $\mathcal{O}_{up}$ . The topology  $\mathcal{O}_{up}$  is known as the *Scott topology*.

The following lemma lists some elementary properties of upper and lower semicontinuous functions. We set  $a \lor b := \max\{a, b\}$  and  $a \land b := \min\{a, b\}$ .

Lemma 1.3 (Upper and lower semi-continuity) (a)  $C(E) = U(E) \cap \mathcal{L}(E)$ .

(b)  $f \in \mathcal{U}(E)$  (resp.  $f \in \mathcal{L}(E)$ ) and  $\lambda \ge 0$  implies  $\lambda f \in \mathcal{U}(E)$  (resp.  $\lambda f \in \mathcal{L}(E)$ ).

(c)  $f, g \in \mathcal{U}(E)$  (resp.  $f, g \in \mathcal{L}(E)$ ) implies  $f + g \in \mathcal{U}(E)$  (resp.  $f + g \in \mathcal{L}(E)$ ).

(d)  $f, g \in \mathcal{U}(E)$  (resp.  $f, g \in \mathcal{L}(E)$ ) implies  $f \lor g \in \mathcal{U}(E)$  and  $f \land g \in \mathcal{U}(E)$  (resp.  $f \lor g \in \mathcal{L}(E)$  and  $f \land g \in \mathcal{L}(E)$ ).

(e)  $f_n \in \mathcal{U}(E)$  and  $f_n \downarrow f$  (resp.  $f_n \in \mathcal{L}(E)$  and  $f_n \uparrow f$ ) implies  $f \in \mathcal{U}(E)$  (resp.  $f \in \mathcal{L}(E)$ ).

(f) An upper (resp. lower) semi-continuous function assumes its maximum (minimum) over a compact set.

**Proof** Part (a) is obvious from the fact that if  $x_n \to x$ , then  $f(x_n) \to f(x)$  if and only if  $\limsup_n f(x_n) \leq f(x)$  and  $\liminf_n f(x_n) \geq f(x)$ . Since f is lower semicontinuous iff -f is upper semi-continuous, it suffices to prove parts (b)–(f) for upper semi-continuous functions. Parts (b) and (d) follow easily from the fact that f is upper semi-continuous if and only if  $\{x : f(x) \ge a\}$  is closed for each  $a \in \mathbb{R}$ , which is equivalent to  $\{x : f(x) < a\}$  being open for each  $a \in \mathbb{R}$ . Indeed,  $f \in \mathcal{U}(E)$ implies that  $\{x : \lambda f(x) < a\} = \{x : f(x) < \lambda^{-1}a\}$  is open for each  $a \in \mathbb{R}, \lambda > 0$ , hence  $\lambda f \in \mathcal{U}(E)$  for each  $\lambda > 0$ , while obviously also  $0 \cdot f \in \mathcal{U}(E)$ . Likewise,  $f, g \in \mathcal{U}(E)$  implies that  $\{x : f(x) \lor g(x) < a\} = \{x : f(x) < a\} \cap \{x : g(x) < a\}$  is open for each  $a \in \mathbb{R}$  hence  $f \lor g \in \mathcal{U}(E)$  and similarly  $\{x : f(x) \land g(x) < a\} = \{x : f(x) \land g(x) < a\}$  $f(x) < a \} \cup \{x : g(x) < a\}$  is open implying that  $f \land g \in \mathcal{U}(E)$ . Part (e) is proved in a similar way: since  $\{x : f_n(x) < a\} \uparrow \{x : f(x) < a\}$ , we conclude that the latter set is open for all  $a \in \mathbb{R}$  hence  $f \in \mathcal{U}(E)$ . Part (c) follows by observing that  $\limsup_{n \to \infty} (f(x_n) + g(x_n)) \le \limsup_{n \to \infty} f(x_n) + \limsup_{m \to \infty} g(x_m) \le f(x) + g(x)$ for all  $x_n \to x$ . To prove part (f), finally let f be upper semi-continuous, K compact, and choose  $a_n \uparrow \sup_{x \in K} f(x)$ . Then  $A_n := \{x \in K : f(x) \ge a_n\}$  is a decreasing sequence of nonempty compact sets, hence (by [Eng89, Corollary 3.1.5]) there exists some  $x \in \bigcap_n A_n$  and f assumes its maximum in x. 

We say that an upper or lower semi-continuous function is *simple* if it assumes only finitely many values.

Lemma 1.4 (Approximation with simple functions) For each  $f \in \mathcal{U}(E)$ there exists simple  $f_n \in \mathcal{U}(E)$  such that  $f_n \downarrow f$ . Analogue statements hold for  $\mathcal{U}_b(E)$ ,  $\mathcal{U}_+(E)$  and  $\mathcal{U}_{b,+}(E)$ . Likewise, lower semi-continuous functions can be approximated from below with simple lower semi-continuous functions.

**Proof** Let  $r_{-} := \inf_{x \in E} f(x)$  and  $r_{+} := \sup_{x \in E} f(x)$ . Let  $\mathcal{D} \subset (r_{-}, r_{+})$  be countable and dense and let  $\Delta_n$  be finite sets such that  $\Delta_n \uparrow \mathcal{D}$ . Let  $\Delta_n = \{a_0, \ldots, a_{m(n)}\}$ with  $a_0 < \cdots < a_{m(n)}$  and set

$$f_n(x) := \begin{cases} a_0 & \text{if } f(x) < a_0, \\ a_k & \text{if } a_{k-1} \le f(x) < a_k \\ r_+ & \text{if } a_{m(n)} \le f(x). \end{cases} \quad (k = 1, \dots, m(n)),$$

Then the  $f_n$  are upper semi-continuous, simple, and  $f_n \downarrow f$ . If  $f \in \mathcal{U}_b(E)$ ,  $\mathcal{U}_+(E)$  or  $\mathcal{U}_{b,+}(E)$  then also the  $f_n$  are in these spaces. The same arguments applied to -f yield the statements for lower semi-continuous functions.

For any set  $A \subset E$  and  $x \in E$ , we let

$$d(x,A) := \inf\{d(x,y) : y \in A\}$$

denote the distance from x to A. Recall that  $\overline{A}$  denotes the closure of A.

**Lemma 1.5 (Distance to a set)** For each  $A \subset E$ , the function  $x \mapsto d(x, A)$  is continuous and satisfies d(x, A) = 0 if and only if  $x \in \overline{A}$ .

**Proof** See [Eng89, Theorem 4.1.10 and Corollary 4.1.11].

**Lemma 1.6 (Approximation of indicator functions)** For each closed  $C \subset E$ there exist continuous  $f_n : E \to [0, 1]$  such that  $f_n \downarrow 1_C$ . Likewise, for each open  $O \subset E$  there exist continuous  $f_n : E \to [0, 1]$  such that  $f_n \uparrow 1_C$ .

**Proof** Set 
$$f_n(x) := (1 - nd(x, C)) \lor 0$$
 resp.  $f_n(x) := nd(x, E \setminus O) \land 1$ .

**Proof of Proposition 1.1** Let  $\mu_n, \mu \in \mathcal{M}(E)$  and define the 'good sets'

$$\mathcal{G}_{up} := \left\{ f \in \mathcal{U}_{b,+}(E) : \limsup_{n \to \infty} \int f d\mu_n \leq \int f d\mu \right\},\$$
$$\mathcal{G}_{low} := \left\{ f \in \mathcal{L}_{b,+}(E) : \liminf_{n \to \infty} \int f d\mu_n \geq \int f d\mu \right\}$$

We claim that

- (a)  $f \in \mathcal{G}_{up}$  (resp.  $f \in \mathcal{G}_{low}$ ),  $\lambda \ge 0$  implies  $\lambda f \in \mathcal{G}_{up}$  (resp.  $\lambda f \in \mathcal{G}_{low}$ ).
- (b)  $f, g \in \mathcal{G}_{up}$  (resp.  $f, g \in \mathcal{G}_{low}$ ) implies  $f + g \in \mathcal{G}_{up}$  (resp.  $f + g \in \mathcal{G}_{low}$ ).
- (c)  $f_n \in \mathcal{G}_{up}$  and  $f_n \downarrow f$  (resp.  $f_n \in \mathcal{G}_{low}$  and  $f_n \uparrow f$ ) implies  $f \in \mathcal{G}_{up}$  (resp.  $f \in \mathcal{G}_{low}$ ).

The statements (a) and (b) are easy. To prove (c), let  $f_n \in \mathcal{G}_{up}$ ,  $f_n \downarrow f$ . Then, for each k,

$$\limsup_{n \to \infty} \int f d\mu_n \le \limsup_{n \to \infty} \int f_k d\mu_n \le \int f_k d\mu.$$

Since  $\int f_k d\mu \downarrow \int f d\mu$ , the claim follows. An analogue argument works for functions in  $\mathcal{G}_{low}$ .

We now show that  $\mu_n \Rightarrow \mu$  implies the conditions (i) and (ii). Indeed, by Lemma 1.6, for each closed  $C \subset E$  we can find continuous  $f_k : E \to [0,1]$  such that  $f_k \downarrow 1_C$ . Then  $f_k \in \mathcal{G}_{up}$  by the fact that  $\mu_n \Rightarrow \mu$  and therefore, by our claim (c) above, it follows that  $1_C \in \mathcal{G}_{up}$ , which proves condition (i). The proof of condition (ii) is similar.

Conversely, if condition (i) is satisfied, then by our claims (a) and (b) above, every simple nonnegative bounded upper semi-continuous function is in  $\mathcal{G}_{up}$ , hence by Lemma 1.4 and claim (c),  $\mathcal{U}_{b,+}(E) \subset \mathcal{G}_{up}$ . Similarly, condition (ii) implies that  $\mathcal{L}_{b,+}(E) \subset \mathcal{G}_{low}$ . In particular, this implies that for every  $f \in \mathcal{C}_{b,+}(E) = \mathcal{U}_{b,+}(E) \cap \mathcal{L}_{b,+}(E)$ ,  $\lim_{n\to\infty} \int f d\mu_n = \int f d\mu$ , which by linearity implies that  $\mu_n \Rightarrow \mu$ .

If the  $\mu_n, \mu$  are probability measures, then conditions (i) and (ii) are equivalent, by taking complements.

### **1.2** Large deviation principles

A subset K of a topological space  $(E, \mathcal{O})$  is called *compact* if every open covering of K has a finite subcovering, i.e., if  $\bigcup_{\gamma \in \Gamma} O_{\gamma} \supset K$  implies that there exist finitely many  $O_{\gamma_1}, \ldots, O_{\gamma_n}$  with  $\bigcup_{k=1}^n O_{\gamma_1} \supset K$ . If  $(E, \mathcal{O})$  is metrizable, then this is equivalent to the statement that every sequence  $x_n \in K$  has a subsequence  $x_{n(m)}$  that converges to a limit  $x \in K$  [Eng89, Theorem 4.1.17]. If  $(E, \mathcal{O})$  is Hausdorff, then each compact subset of E is closed.

#### 1.2. LARGE DEVIATION PRINCIPLES

Let E be a Polish space. We say that a function  $f: E \to \overline{\mathbb{R}}$  has compact level sets if

$$\{x \in E : f(x) \le a\}$$
 is compact for all  $a \in \mathbb{R}$ .

Note that since compact sets are closed, this is (a bit) stronger than the statement that f is lower semi-continuous. We say that I is a good rate function if I has compact level sets,  $-\infty < I(x)$  for all  $x \in E$ , and  $I(x) < \infty$  for at least one  $x \in E$ . Note that by Lemma 1.3 (f), such a function is necessarily bounded from below.

Recall that  $B_b(E)$  denotes the space of all bounded Borel-measurable real functions on E. If  $\mu$  is a finite measure on  $(E, \mathcal{B}(E))$  and  $p \ge 1$  is a real constant, then we define the  $L^p$ -norm associated with  $\mu$  by

$$||f||_{p,\mu} := \left(\int \mathrm{d}\mu |f|^p\right)^{1/p} \qquad (f \in B_b(E)).$$

Likewise, if I is a good rate function, then we can define a sort of 'weighted supremumnorm' by

$$||f||_{\infty,I} := \sup_{x \in E} e^{-I(x)} |f(x)| \qquad (f \in B_b(E)).$$
(1.1)

Note that  $||f||_{\infty,I} < \infty$  by the boundedness of f and the fact that I is bounded from below. It is easy to check that  $|| \cdot ||_{\infty,I}$  is a *seminorm*, i.e.,

- $\|\lambda f\|_{\infty,I} = |\lambda| \|f\|_{\infty,I},$
- $||f + g||_{\infty,I} \le ||f||_{\infty,I} + ||g||_{\infty,I}.$

If  $I < \infty$  then  $\| \cdot \|_{\infty,I}$  is moreover a norm, i.e.,

•  $||f||_{\infty,I} = 0$  implies f = 0.

Note that what we have just called  $L^p$ -norm is in fact only a seminorm, since  $||f||_{p,\mu} = 0$  only implies that f = 0 a.e. w.r.t.  $\mu$ . (This is usually resolved by looking at equivalence classes of a.e. equal functions, but we won't need this here.)

(Large deviation principle) Let  $s_n$  be positive constants converging to  $\infty$ , let  $\mu_n$  be finite measures on E, and let I be a good rate function on E. We say that the  $\mu_n$  satisfy the large deviation principle (LDP) with speed (also called rate)  $s_n$  and rate function I if

$$\lim_{n \to \infty} \|f\|_{s_n, \mu_n} = \|f\|_{\infty, I} \qquad (f \in \mathcal{C}_{b, +}(E)).$$
(1.2)

While this definition may look a bit strange at this point, the next proposition looks already much more similar to things we have seen in Chapter 0.

**Proposition 1.7 (Large Deviation Principle)** A sequence of finite measures  $\mu_n$  satisfies the large deviation principle with speed  $s_n$  and rate function I if and only if the following two conditions are satisfied.

(i)  $\limsup_{n \to \infty} \frac{1}{s_n} \log \mu_n(C) \le -\inf_{x \in C} I(x) \quad \forall C \text{ closed,}$ (ii)  $\liminf_{n \to \infty} \frac{1}{s_n} \log \mu_n(O) \ge -\inf_{x \in O} I(x) \quad \forall O \text{ open.}$ 

**Remark 1** Recall that  $\overline{A}$  and  $\operatorname{int}(A)$  denote the closure and interior of a set  $A \subset E$ , respectively. Since for any measurable set A, one has  $\mu_n(A) \leq \mu_n(\overline{A})$  and  $\mu_n(A) \geq \mu_n(\operatorname{int}(A))$ , conditions (i) and (ii) of Proposition 1.7 are equivalent to

(i)' 
$$\limsup_{n \to \infty} \frac{1}{s_n} \log \mu_n(A) \le -\inf_{x \in \overline{A}} I(x),$$
  
(ii)' 
$$\liminf_{n \to \infty} \frac{1}{s_n} \log \mu_n(A) \ge -\inf_{x \in \operatorname{int}(A)} I(x),$$

for all  $A \in \mathcal{B}(E)$ . We say that a set  $A \in \mathcal{B}(E)$  is *I*-continuous if

$$\inf_{x \in \operatorname{int}(A)} I(x) = \inf_{x \in \overline{A}} I(x)$$

It is now easy to see that if  $\mu_n$  satisfy the large deviation principle with speed  $s_n$ and good rate function I, then

$$\lim_{n \to \infty} \frac{1}{s_n} \log \mu_n(A) = -\inf_{x \in A} I(x)$$

for each *I*-continuous set *A*. For example, if *I* is continuous and  $\overline{A} = int(A)$ , then *A* is *I*-continuous. This is the reason, for example, why in our formulation of the Boltzmann-Sanov Theorem 0.7 we looked at sets that are the closure of their interior.

**Remark 2** The two conditions of Proposition 1.7 are the traditional definition of a large deviation principle. Moreover, large deviation principles are often only defined for the special case that the speed  $s_n$  equals n. However, as the example

### 1.2. LARGE DEVIATION PRINCIPLES

of moderate deviations (Theorem 0.4) showed, it is sometimes convenient to allow more general speeds. Also parts of the abstract theory (in particular, connected to the concept of exponential tightness) are more easy to formulate if one allows general speeds. As we will see, allowing more general speeds will not cause any technical complications so this generality comes basically 'for free'.

To prepare for the proof of Proposition 1.7, we start with some preliminary lemmas.

Lemma 1.8 (Properties of the generalized supremumnorm) Let I be a good rate function and let  $\|\cdot\|_{\infty,I}$  be defined as in (1.1). Then

- (a)  $||f \vee g||_{\infty,I} = ||f||_{\infty,I} \vee ||g||_{\infty,I} \quad \forall f, g \in B_{b,+}(E).$
- (b)  $||f_n||_{\infty,I} \uparrow ||f||_{\infty,I} \forall f_n \in B_{b,+}(E), f_n \uparrow f.$
- (c)  $||f_n||_{\infty,I} \downarrow ||f||_{\infty,I} \forall f_n \in \mathcal{U}_{b,+}(E), f_n \downarrow f.$

**Proof** Property (a) follows by writing

$$\|f \vee g\|_{\infty,I} = \sup_{\substack{x \in E \\ x \in E}} e^{-I(x)} (f(x) \vee g(x))$$
  
=  $\left(\sup_{x \in E} e^{-I(x)} f(x)\right) \vee \left(\sup_{y \in E} e^{-I(x)} g(y)\right) = \|f\|_{\infty,I} \vee \|g\|_{\infty,I}$ 

To prove (b), we start by observing that the  $||f_n||_{\infty,I}$  form an increasing sequence and  $||f_n||_{\infty,I} \leq ||f||_{\infty,I}$  for each n. Moreover, for any  $\varepsilon > 0$  we can find  $y \in E$  such that  $e^{-I(y)}f(y) \geq \sup_{x \in E} e^{-I(x)}f(x) - \varepsilon$ , hence  $\liminf_n ||f_n||_{\infty,I} \geq \lim_n e^{-I(y)}f_n(y) = e^{-I(y)}f(y) \geq ||f||_{\infty,I} - \varepsilon$ . Since  $\varepsilon > 0$  is arbitrary, this proves the claim.

To prove also (c), we start by observing that the  $||f_n||_{\infty,I}$  form a decreasing sequence and  $||f_n||_{\infty,I} \ge ||f||_{\infty,I}$  for each n. Since the  $f_n$  are upper semi-continuous and Iis lower semi-continuous, the functions  $e^{-I}f_n$  are upper semi-continuous. Since the  $f_n$  are bounded and I has compact level sets, the sets  $\{x : e^{-I(x)}f_n(x) \ge a\}$ are compact for each a > 0. In particular, for each  $a > \sup_{x \in E} e^{-I(x)}f(x)$ , the sets  $\{x : e^{-I(x)}f_n(x) \ge a\}$  are compact and decrease to the empty set, hence  $\{x : e^{-I(x)}f_n(x) \ge a\} = \emptyset$  for n sufficiently large, which shows that  $\limsup_n ||f_n||_{\infty,I} \le a$ .

**Lemma 1.9 (Good sets)** Let  $\mu_n \in \mathcal{M}(E)$ ,  $s_n \to \infty$ , and let I be a good rate function. Define the 'good sets'

$$\mathcal{G}_{up} := \left\{ f \in \mathcal{U}_{b,+}(E) : \limsup_{n \to \infty} \|f\|_{s_n,\mu_n} \le \|f\|_{\infty,I} \right\},$$
  
$$\mathcal{G}_{low} := \left\{ f \in \mathcal{L}_{b,+}(E) : \liminf_{n \to \infty} \|f\|_{s_n,\mu_n} \ge \|f\|_{\infty,I} \right\}.$$

Then

- (a)  $f \in \mathcal{G}_{up}$  (resp.  $f \in \mathcal{G}_{low}$ ),  $\lambda \ge 0$  implies  $\lambda f \in \mathcal{G}_{up}$  (resp.  $\lambda f \in \mathcal{G}_{low}$ ).
- (b)  $f, g \in \mathcal{G}_{up}$  (resp.  $f, g \in \mathcal{G}_{low}$ ) implies  $f \lor g \in \mathcal{G}_{up}$  (resp.  $f \lor g \in \mathcal{G}_{low}$ ).
- (c)  $f_n \in \mathcal{G}_{up}$  and  $f_n \downarrow f$  (resp.  $f_n \in \mathcal{G}_{low}$  and  $f_n \uparrow f$ ) implies  $f \in \mathcal{G}_{up}$  (resp.  $f \in \mathcal{G}_{low}$ ).

The proof of Lemma 1.9 makes use of the following elementary lemma.

**Lemma 1.10 (The strongest growth wins)** For any  $0 \le a_n, b_n \le \infty$  and  $s_n \to \infty$ , one has

$$\limsup_{n \to \infty} \left( a_n^{s_n} + b_n^{s_n} \right)^{1/s_n} = \left( \limsup_{n \to \infty} a_n \right) \vee \left( \limsup_{n \to \infty} b_n \right). \tag{1.3}$$

Moreover, for any  $0 \leq c_n, d_n \leq \infty$  and  $s_n \to \infty$ ,

$$\limsup_{n \to \infty} \frac{1}{s_n} \log(c_n + d_n) = \left(\limsup_{n \to \infty} \frac{1}{s_n} \log c_n\right) \vee \left(\limsup_{n \to \infty} \frac{1}{s_n} \log d_n\right).$$
(1.4)

**Proof** To see this, set  $a_{\infty} := \limsup_{n \to \infty} a_n$  and  $b_{\infty} := \limsup_{n \to \infty} b_n$ . Then, for each  $\varepsilon > 0$ , we can find an m such that  $a_n \leq a_{\infty} + \varepsilon$  and  $b_n \leq b_{\infty} + \varepsilon$  for all  $n \geq m$ . It follows that

$$\limsup_{n \to \infty} \left( a_n^{s_n} + b_n^{s_n} \right)^{1/s_n} \le \lim_{n \to \infty} \left( (a_\infty + \varepsilon)^{s_n} + (b_\infty + \varepsilon)^{s_n} \right)^{1/s_n} = (a_\infty + \varepsilon) \lor (b_\infty + \varepsilon).$$

Since  $\varepsilon > 0$  is arbitrary, this shows that  $\limsup_{n\to\infty} (a_n^{s_n} + b_n^{s_n})^{1/s_n} \leq a_\infty \vee b_\infty$ . Since  $a_n, b_n \leq (a_n^{s_n} + b_n^{s_n})^{1/s_n}$ , the other inequality is trivial. This completes the proof of (1.3).

We claim that (1.4) is just (1.3) in another guise. Indeed, setting  $a_n := c_n^{1/s_n}$  and  $b_n := d_n^{1/s_n}$  we see, using (1.3), that

$$e^{\limsup_{n \to \infty} \frac{1}{s_n} \log(c_n + d_n)} = \limsup_{n \to \infty} (a_n^{s_n} + d_n^{s_n})^{1/s_n}$$
$$= (\limsup_{n \to \infty} a_n) \lor (\limsup_{n \to \infty} b_n)$$
$$= e^{(\limsup_{n \to \infty} \frac{1}{s_n} \log(c_n))} \lor (\limsup_{n \to \infty} \frac{1}{s_n} \log(d_n)).$$

**Proof of Lemma 1.9** Part (a) follows from the fact that for any seminorm  $\|\lambda f\| = \lambda \|f\|$  ( $\lambda > 0$ ). To prove part (b), we first make a simple observation. Now

### 1.2. LARGE DEVIATION PRINCIPLES

assume that  $f, g \in \mathcal{G}_{up}$ . Then, by (1.3),

$$\lim_{n \to \infty} \sup_{n \to \infty} \|f \lor g\|_{s_n,\mu_n} \\
= \limsup_{n \to \infty} \left( \int_{\{x: f(x) \ge g(x)\}} f(x)^{s_n} \mu_n(\mathrm{d}x) + \int_{\{x: f(x) < g(x)\}} g(x)^{s_n} \mu_n(\mathrm{d}x) \right)^{1/s_n} \\
\leq \limsup_{n \to \infty} \left( \|f\|_{s_n,\mu_n}^{s_n} + \|g\|_{s_n,\mu_n}^{s_n} \right)^{1/s_n} \le \|f\|_{\infty,I} \lor \|g\|_{\infty,I} = \|f \lor g\|_{\infty,I},$$
(1.5)

proving that  $f \lor g \in \mathcal{G}_{up}$ . Similarly, but easier, if  $f, g \in \mathcal{G}_{low}$ , then

$$\liminf_{n \to \infty} \|f \vee g\|_{s_n,\mu_n} \ge \left(\liminf_{n \to \infty} \|f\|_{s_n,\mu_n}\right) \vee \left(\liminf_{n \to \infty} \|g\|_{s_n,\mu_n}\right)$$
$$\ge \|f\|_{\infty,I} \vee \|g\|_{\infty,I} = \|f \vee g\|_{\infty,I},$$

which proves that  $f \lor g \in \mathcal{G}_{low}$ .

To prove part (c), finally, assume that  $f_k \in \mathcal{G}_{up}$  satisfy  $f_k \downarrow f$ . Then f is upper semi-continuous and

$$\limsup_{n \to \infty} \|f\|_{s_n, \mu_n} \le \limsup_{n \to \infty} \|f_k\|_{s_n, \mu_n} \le \|f_k\|_{\infty, I}$$

for each k. Since  $||f_k||_{\infty,I} \downarrow ||f||_{\infty,I}$ , by Lemma 1.8 (c), we conclude that  $f \in \mathcal{G}_{up}$ . The proof for  $f_k \in \mathcal{G}_{low}$  is similar, using Lemma 1.8 (b).

**Proof of Proposition 1.7** If the  $\mu_n$  satisfy the large deviation principe with speed  $s_n$  and rate function I, then by Lemmas 1.6 and 1.9 (c),  $1_C \in \mathcal{G}_{up}$  for each closed  $C \subset E$  and  $1_O \in \mathcal{G}_{up}$  for each open  $O \subset E$ , which shows that conditions (i) and (ii) are satisfied. Conversely, if conditions (i) and (ii) are satisfied, then by Lemma 1.9 (a) and (b),

$$\mathcal{G}_{up} \supset \{f \in \mathcal{U}_{b,+}(E) : f \text{ simple}\} \text{ and } \mathcal{G}_{low} \supset \{f \in \mathcal{L}_{b,+}(E) : f \text{ simple}\}.$$

By Lemmas 1.4 and 1.9 (c), it follows that  $\mathcal{G}_{up} = \mathcal{U}_{b,+}(E)$  and  $\mathcal{G}_{low} = \mathcal{L}_{b,+}(E)$ . In particular, this proves that

$$\lim_{n \to \infty} \|f\|_{s_n, \mu_n} = \|f\|_{\infty, I} \qquad \forall f \in \mathcal{C}_{b, +}(E),$$

which shows that the  $\mu_n$  satisfy the large deviation principe with speed  $s_n$  and rate function I.

**Exercise 1.11 (Robustness of LDP)** Let  $(X_k)_{k\geq 1}$  be i.i.d. random variables with  $\mathbb{P}[X_k = 0] = \mathbb{P}[X_k = 1] = \frac{1}{2}$ , let  $Z(\lambda) := \mathbb{E}[e^{\lambda X_1}]$  ( $\lambda \in \mathbb{R}$ ) and let  $I : \mathbb{R} \to [0, \infty]$  be defined as in (0.3). Let  $\varepsilon_n \downarrow 0$  and set

$$T_n := \frac{1}{n} \sum_{k=1}^n X_k$$
 and  $T'_n := (1 - \varepsilon_n) \frac{1}{n} \sum_{k=1}^n X_k$ 

In Theorem 2.17 below, we will prove that the laws  $\mathbb{P}[T_n \in \cdot]$  satisfy the large deviation principle with speed n and rate function I. Using this fact, prove that also the laws  $\mathbb{P}[T'_n \in \cdot]$  satisfy the large deviation principle with speed n and rate function I. Use Lemma 0.2 to conclude that

$$\lim_{n \to \infty} \frac{1}{n} \log \mathbb{P}[T'_n \ge y] = -I(y) \quad (\frac{1}{2} \le y < 1),$$

but this formula does *not* hold for y = 1.

### 1.3 Varadhan's lemma

The two conditions of Proposition 1.7 are the traditional definition of the large deviation principle, which is due to Varadhan [Var66]. Our alternative, equivalent definition in terms of convergence of  $L_p$ -norms is very similar to the road followed in Puhalskii's book [Puh01]. A very similar definition is also given in [DE97], where this is called a 'Laplace principle' instead of a large deviation principle.

From a purely abstract point of view, our definition is frequently a bit easier to work with. On the other hand, the two conditions of Proposition 1.7 are closer to the usual interpretation of large deviations in terms of exponentially small probabilities. Also, when in some practical situation one wishes to prove a large deviation principle, the two conditions of Proposition 1.7 are often a very natural way to do so. Here, condition (ii) is usually easier to check than condition (i). Condition (ii) says that certain rare events occur wih at least a certain probability. To prove this, one needs to find one strategy by which a stochastic system can make the desired event happen, with a certain small probability. Condition (i) says that there are no other strategies that yield a higher probability for the same event, which requires one to prove something about all possible ways in which a certain event can happen.

In practically all applications, we will only be interested in the case that the measures  $\mu_n$  are probability measures and the rate function satisfies  $\inf_{x \in E} I(x) = 0$ , but being slightly more general comes at virtually no cost.

Varadhan [Var66] was not only the first one who formulated large deviation principles in the generality that is now standard, he also first proved the lemma that is called after him, and that reads as follows.

**Lemma 1.12 (Varadhan's lemma)** Let E be a Polish space and let  $\mu_n \in \mathcal{M}(E)$ satisfy the large deviation principle with speed  $s_n$  and good rate function I. Let

### 1.3. VARADHAN'S LEMMA

 $F: E \to \mathbb{R}$  be continuous and assume that  $\sup_{x \in E} F(x) < \infty$ . Then

$$\lim_{n \to \infty} \frac{1}{s_n} \log \int e^{s_n F} \mathrm{d}\mu_n = \sup_{x \in E} [F(x) - I(x)].$$

**Proof** Applying the exponential function to both sides of our equation, this says that

$$\lim_{n \to \infty} \left( \int e^{s_n F} \mathrm{d}\mu_n \right)^{1/s_n} = \sup_{x \in E} e^{F(x) - I(x)}$$

Setting  $f := e^F$ , this is equivalent to

$$\lim_{n \to \infty} \|f\|_{s_n, \mu_n} = \|f\|_{\infty, I},$$

where our asumptions on F translate into  $f \in C_{b,+}(E)$ . Thus, Varadhan's lemma is just a trivial reformulation of our definition of a large deviation principle. If we take the traditional definition of a large deviation principle as our starting point, then Varadhan's lemma corresponds to the 'if' part of Proposition 1.7.

As we have just seen, Varadhan's lemma is just the statement that the two conditions of Proposition 1.7 are sufficient for (1.2). The fact that these conditions are also necessary was only proved 24 years later, by Bryc [Bry90].

We conclude this section with a little lemma that says that a sequence of measures satisfying a large deviation principle determines its rate function uniquely.

**Lemma 1.13 (Uniqueness of the rate function)** Let E be a Polish space,  $\mu_n \in \mathcal{M}(E)$ , and let  $s_n$  be real constants converging to infinity. Assume that the  $\mu_n$  satisfy the large deviation principle with speed  $s_n$  and good rate function I and also that the  $\mu_n$  satisfy the large deviation principle with speed  $s_n$  and good rate function I'. Then I = I'.

**Proof** It follows immediately from our definition of the large deviation principle that  $||f||_{\infty,I} = ||f||_{\infty,I'}$  for all  $f \in \mathcal{C}_{b,+}(E)$ . By Lemma 1.6, for each  $x \in E$ , we can find continuous  $f_n : E \to [0,1]$  such that  $f_n \downarrow 1_{\{x\}}$ . By Lemma 1.8 (c), it follows that

$$e^{-I(x)} = \|1_{\{x\}}\|_{\infty,I} = \lim_{n \to \infty} \|f_n\|_{\infty,I} = \lim_{n \to \infty} \|f_n\|_{\infty,I'} = \|1_{\{x\}}\|_{\infty,I'} = e^{-I'(x)}$$

for each  $x \in E$ .

### **1.4** The contraction principle

As we have seen in Propositions 1.1 and 1.7, there is a lot of similarity between weak convergence and the large deviation principle. Elaborating on this analogy, we recall that if  $X_n$  is a sequence of random variables, taking values in some Polish space E, whose laws converge weakly to the law of a random variable X, and  $\psi : E \to F$  is a continuous function from E into some other Polish space, then the laws of the random variables  $\psi(X_n)$  converge weakly to the law of  $\psi(X)$ . As we will see, an analogue statement holds for sequences of measures satisfying a large deviation principle.

Recall that if X is a random variable taking values in some measurable space  $(E, \mathcal{E})$ , with law  $\mathbb{P}[X \in \cdot] = \mu$ , and  $\psi : E \to F$  is a measurable function from E into some other measurable space  $(F, \mathcal{F})$ , then the law of  $\psi(X)$  is the *image measure* 

 $\mu \circ \psi^{-1}(A) \quad (A \in \mathcal{F}), \quad \text{where} \quad \psi^{-1}(A) := \{x \in E : \psi(x) \in A\}$ 

is the *inverse image* (or *pre-image*) of A under  $\psi$ .

The next result shows that if  $X_n$  are random variables whose laws satisfy a large deviation principle, and  $\psi$  is a continuous function, then also the laws of the  $\psi(X_n)$  satify a large deviation principle. This fact is known a the *contraction principle*. Note that we have already seen this principle at work when we derived Proposition 0.5 from Theorem 0.7. As is clear from this example, it is in practice not always easy to explicitly calculate the 'image' of a rate function under a continuous map, as defined formally in (1.6) below.

**Proposition 1.14 (Contraction principle)** Let E, F be Polish spaces and let  $\psi: E \to F$  be continuous. Let  $\mu_n$  be finite measures on E satisfying a large deviation principle with speed  $s_n$  and good rate function I. Then the image measures  $\mu \circ \psi^{-1}$  satisfying the large deviation principle with speed  $s_n$  and good rate function J given by

$$J(y) := \inf_{x \in \psi^{-1}(\{y\})} I(x) \qquad (y \in F),$$
(1.6)

where  $\inf_{x \in \emptyset} I(x) := \infty$ .

**Proof** Recall that a function  $\psi$  from one topological space E into another topological space F is continuous if and only if the inverse image under  $\psi$  of any open set is open, or equivalently, the inverse image of any closed set is closed (see, e.g.,

#### 1.5. EXPONENTIAL TILTS

[Eng89, Proposition 1.4.1] or [Kel75, Theorem 3.1]). As a result, condition (i) of Proposition 1.7 implies that

$$\limsup_{n \to \infty} \frac{1}{s_n} \log \mu_n \circ \psi^{-1}(C) \le -\inf_{x \in \psi^{-1}(C)} I(x) = -\inf_{y \in C} \inf_{x \in \psi^{-1}(\{y\})} I(x) = -\inf_{y \in C} J(y),$$
(1.7)

where we have used that  $\psi^{-1}(C) = \bigcup_{y \in C} \psi^{-1}(y)$ . Condition (ii) of Proposition 1.7 carries over in the same way. We are left with the task of showing that J is a good rate function. Indeed, for each  $a \in \mathbb{R}$  the level set

$$\{y \in F : J(y) \le a\} = \{y \in F : \inf_{x \in \psi^{-1}(\{y\})} I(x) \le a\}$$
  
=  $\{y \in F : \exists x \in E \text{ s.t. } \psi(x) = y, \ I(x) \le a\}$   
=  $\{\psi(x) : x \in E, \ I(x) \le a\} = \psi(\{x : I(x) \le a\})$ 

is the image under  $\psi$  of the level set  $\{x : I(x) \leq a\}$ . Since the continuous image of a compact set is compact [Eng89, Theorem 3.1.10],<sup>1</sup> this proves that Jhas compact level sets. Finally, we observe (compare (1.7)) that  $\inf_{y \in F} J(y) = \inf_{x \in \psi^{-1}(F)} I(x) = \inf_{x \in E} I(x) < \infty$ , proving that J is a good rate function.

### **1.5** Exponential tilts

It is not hard to see that if  $\mu_n$  are measures satisfying a large deviation principle, then we can transform these measures by weighting them with an exponential density, in such a way that the new measures also satisfy a large deviation principle. Recall that if  $\mu$  is a measure and f is a nonnegative measurable function, then setting

$$f\mu(A) := \int_A f \mathrm{d}\mu$$

defines a new measure  $f\mu$  which is  $\mu$  weighted with the density f.

**Lemma 1.15 (Exponential weighting)** Let E be a Polish space and let  $\mu_n \in \mathcal{M}(E)$  satisfy the large deviation principle with speed  $s_n$  and good rate function I. Let  $F : E \to \mathbb{R}$  be continuous and assume that  $-\infty < \sup_{x \in E} F(x) < \infty$ . Then the measures

$$\tilde{\mu}_n := e^{s_n F} \mu_n$$

<sup>&</sup>lt;sup>1</sup>This is a well-known fact that can be found in any book on general topology. It is easy to show by counterexample that the continuous image of a *closed* set needs in general not be closed!

satisfy the large deviation principle with speed  $s_n$  and good rate function  $\tilde{I} := I - F$ .

**Proof** Note that  $e^F \in \mathcal{C}_{b,+}(E)$ . Therefore, for any  $f \in \mathcal{C}_{b,+}(E)$ ,

$$\|f\|_{s_{n},\tilde{\mu}_{n}} = \int f^{s_{n}} e^{s_{n}F} d\mu_{n} = \|fe^{F}\|_{s_{n},\mu_{n}}$$
$$\xrightarrow[n \to \infty]{} \|fe^{F}\|_{\infty,I} = \sup_{x \in E} f(x)e^{F(x)}e^{-I(x)} = \|f\|_{\infty,\tilde{I}}$$

Since F is continuous, I - F is lower semi-continuous. Since F is bounded from above, any level set of I - F is contained in some level set of I, and therefore compact. Since F is not identically  $-\infty$ , finally,  $\inf_{x \in I}(I(x) - F(x)) < \infty$ , proving that I - F is a good rate function.

Lemma 1.15 is not so useful yet, since in practice we are usually interested in probability measures, while exponential weighting may spoil the normalization. Likewise, we are usually interested in rate functions that are properly 'normalized'. Let us say that a function I is a normalized rate function if I is a good rate function and  $\inf_{x \in E} I(x) = 0$ . Note that if  $\mu_n$  are probability measures satisfying a large deviation principle with speed  $s_n$  and rate function I, then I must be normalized, since E is both open and closed, and therefore by conditions (i) and (ii) of Proposition 1.7

$$-\inf_{x\in E} I(x) = \lim_{n\to\infty} \frac{1}{s_n} \log \mu_n(E) = 0.$$

**Lemma 1.16 (Exponential tilting)** Let E be a Polish space and let  $\mu_n$  be probability measures on E satisfy the large deviation principle with speed  $s_n$  and normalized rate function I. Let  $F : E \to \mathbb{R}$  be continuous and assume that  $-\infty < \sup_{x \in E} F(x) < \infty$ . Then the measures

$$\tilde{\mu}_n := \frac{1}{\int e^{s_n F} \mathrm{d}\mu_n} e^{s_n F} \mu_n$$

satisfy the large deviation principle with speed  $s_n$  and normalized rate function  $\tilde{I}(x) := I(x) - F(x) - \inf_{y \in E} (I(y) - F(y)).$ 

**Proof** Since  $e^F \in \mathcal{C}_{b,+}(E)$ , much in the same way as in the proof of the previous lemma, we see that

$$\begin{split} \|f\|_{s_{n},\tilde{\mu}_{n}} &= \left(\frac{1}{\int e^{s_{n}F} d\mu_{n}} \int f^{s_{n}} e^{s_{n}F} d\mu_{n}\right)^{1/s_{n}} = \frac{\|fe^{F}\|_{s_{n},\mu_{n}}}{\|e^{F}\|_{\infty,I}} \\ & \xrightarrow[n \to \infty]{} \frac{\|fe^{F}\|_{\infty,I}}{\|e^{F}\|_{\infty,I}} = \frac{\sup_{x \in E} f(x)e^{F(x)}e^{-I(x)}}{\sup_{x \in E} e^{F(x)}e^{-I(x)}} \\ &= e^{-\inf_{y \in E}(I(y)-F(y))} \sup_{x \in E} f(x)e^{-(I(x)-F(x))} = \|f\|_{\infty,\tilde{I}}. \end{split}$$

The fact that  $\tilde{I}$  is a good rate function follows from the same arguments as in the proof of the previous lemma, and  $\tilde{I}$  is obviously normalized.

### 1.6 Robustness

Often, when one wishes to prove that the laws  $\mathbb{P}[X_n \in \cdot]$  of some random variables  $X_n$  satisfy a large deviation principle with a given speed and rate function, it is convenient to replace the random variables  $X_n$  by some other random variables  $Y_n$  that are 'sufficiently close', so that the large deviation principle for the laws  $\mathbb{P}[Y_n \in \cdot]$  implies the LDP for  $\mathbb{P}[X_n \in \cdot]$ . The next result (which we copy from [DE97, Thm 1.3.3]) gives sufficient conditions for this to be allowed.

**Proposition 1.17 (Superexponential approximation)** Let  $(X_n)_{n\geq 1}$ ,  $(Y_n)_{n\geq 1}$ be random variables taking values in a Polish space E and assume that the laws  $\mathbb{P}[Y_n \in \cdot]$  satisfy a large deviation principle with speed  $s_n$  and rate function I. Let d be any metric generating the topology on E, and assume that

$$\lim_{n \to \infty} \frac{1}{s_n} \log \mathbb{P}[d(X_n, Y_n) \ge \varepsilon] = -\infty \qquad (\varepsilon > 0).$$
(1.8)

Then the laws  $\mathbb{P}[X_n \in \cdot]$  satisfy the large deviation principle with speed  $s_n$  and rate function I.

**Remark** If (1.8) holds, then we say that the random variables  $X_n$  and  $Y_n$  are exponentially close. Note that condition (1.8) is in particular satisfied if for each  $\varepsilon > 0$  there is an N such that  $d(X_n, Y_n) < \varepsilon$  a.s. for all  $n \ge N$ . We can even allow for  $d(X_n, Y_n) \ge \varepsilon$  with a small probability, but in this case these probabilities must tend to zero faster than any exponential.

**Proof of Proposition 1.17** Let  $C \subset E$  be closed and let  $C_{\varepsilon} := \{x \in E : d(x, C) \leq \varepsilon\}$ . Then

$$\begin{split} \limsup_{n \to \infty} \frac{1}{s_n} \log \mathbb{P}[X_n \in C] \\ &\leq \limsup_{n \to \infty} \frac{1}{s_n} \log \left( \mathbb{P}[Y_n \in C_{\varepsilon}, \ d(X_n, Y_n) \leq \varepsilon] + \mathbb{P}[d(X_n, Y_n) > \varepsilon] \right) \\ &\leq \limsup_{n \to \infty} \frac{1}{s_n} \log \mathbb{P}[Y_n \in C_{\varepsilon}] = -\inf_{x \in C_{\varepsilon}} I(x) \xrightarrow[\varepsilon \downarrow 0]{} - \inf_{x \in C} I(x), \end{split}$$

where we have used (1.4) and in the last step we have applied (the logarithmic version of) Lemma 1.8 (c). Similarly, if  $O \subset E$  is open and  $O_{\varepsilon} := \{x \in E : d(x, E \setminus O) > \varepsilon\}$ , then

$$\liminf_{n \to \infty} \frac{1}{s_n} \log \mathbb{P}[X_n \in O] \ge \liminf_{n \to \infty} \frac{1}{s_n} \log \mathbb{P}[Y_n \in O_{\varepsilon}, \ d(X_n, Y_n) \le \varepsilon]$$

The large deviations lower bound is trivial if  $\inf_{x\in O} I(x) = \infty$ , so without loss of generality we may assume that  $\inf_{x\in O} I(x) < \infty$ . Since  $\inf_{x\in O_{\varepsilon}} I(x) \downarrow \inf_{x\in O} I(x)$ , it follows that for  $\varepsilon$  sufficiently small, also  $\inf_{x\in O_{\varepsilon}} I(x) < \infty$ . By the fact that the  $Y_n$  satisfy the large deviation lower bound and by (1.8),

$$\mathbb{P}[Y_n \in O_{\varepsilon}, \ d(X_n, Y_n) \le \varepsilon] \ge \mathbb{P}[Y_n \in O_{\varepsilon}] - \mathbb{P}[d(X_n, Y_n) > \varepsilon]$$
  
>  $e^{-s_n \inf_{x \in O_{\varepsilon}} I(x) + o(s_n)} - e^{-s_n/o(s_n)}$ 

as  $n \to \infty$ , where  $o(s_n)$  is the usual small 'o' notation, i.e.,  $o(s_n)$  denotes any term such that  $o(s_n)/s_n \to 0$ . It follows that

$$\liminf_{n \to \infty} \frac{1}{s_n} \log \mathbb{P}[Y_n \in O_{\varepsilon}, \ d(X_n, Y_n) \le \varepsilon] \ge -\inf_{x \in O_{\varepsilon}} I(x) \xrightarrow[\varepsilon \downarrow 0]{} - \inf_{x \in O} I(x),$$

which proves the the large deviation lower bound for the  $X_n$ .

Proposition 1.17 shows that large deviation principles are 'robust', in a certain sense, with repect to small perturbations. The next result is of a similar nature: we will prove that weighting measures with densities does not affect a large deviation principle, as long as these densities do not grow exponentially fast. This complements the case of exponentially growing densities which has been treated in Section 1.5.

**Lemma 1.18 (Subexponential weighting)** Let E be a Polish space and let  $\mu_n \in \mathcal{M}(E)$  satisfy the large deviation principle with speed  $s_n$  and good rate function I. Let  $F_n : E \to \mathbb{R}$  be measurable and assume that  $\lim_{n\to\infty} ||F_n||_{\infty} = 0$ , where

 $||F_n||_{\infty} := \sup_{x \in E} |F_n(x)|$ . Then the measures

$$\tilde{\mu}_n := e^{s_n F_n} \mu_n$$

satisfy the large deviation principle with speed  $s_n$  and rate function I.

**Proof** We check the large deviations upper and lower bound from Proposition 1.7. For any closed set  $C \subset E$ , by the fact that the  $\mu_n$  satisfy the large deviation principle, we have

$$\limsup_{n \to \infty} \frac{1}{s_n} \log \tilde{\mu}_n(C) = \limsup_{n \to \infty} \frac{1}{s_n} \log \int_C \mu_n(\mathrm{d}x) e^{s_n F_n(x)}$$
  
$$\leq \limsup_{n \to \infty} \frac{1}{s_n} \log \left( e^{s_n \|F_n\|} \mu_n(C) \right) = \limsup_{n \to \infty} \left( \|F_n\| + \frac{1}{s_n} \log \mu_n(C) \right),$$

which equals  $-\inf_{x\in C} I(x)$ . Similarly, for any open  $O \subset E$ , we have

$$\liminf_{n \to \infty} \frac{1}{s_n} \log \tilde{\mu}_n(O) = \liminf_{n \to \infty} \frac{1}{s_n} \log \int_O \mu_n(\mathrm{d}x) e^{s_n F_n(x)}$$
$$\geq \liminf_{n \to \infty} \frac{1}{s_n} \log \left( e^{-s_n \|F_n\|} \mu_n(O) \right) = \liminf_{n \to \infty} \left( -\|F_n\| + \frac{1}{s_n} \log \mu_n(O) \right),$$

which yields  $-\inf_{x \in O} I(x)$ , as required.

### 1.7 Tightness

In Sections 1.1 and 1.2, we have stressed the similarity between weak convergence of measures and large deviation principles. In the remainder of this chapter, we will pursue this idea further. In the present section, we recall the concept of tightness and Prohorov's theorem. In particular, we will see that any tight sequence of probability measures on a Polish space has a weakly convergent subsequence. In the next sections (to be precise, in Theorem 1.24), we will prove an analogue of this result, which says that every exponentially tight sequence of probability measures on a Polish space has a subsequence that satisfies a large deviation principle.

A set A is called *relatively compact* if its closure A is compact. The next result is known as Prohorov's theorem (see, e.g., [Ste87, Theorems III.3.3 and III.3.4] or [Bil99, Theorems 5.1 and 5.2]).

**Proposition 1.19 (Prohorov)** Let E be a Polish space and let  $\mathcal{M}_1(E)$  be the space of probability measures on  $(E, \mathcal{B}(E))$ , equipped with the topology of weak convergence. Then a subset  $\mathcal{C} \subset \mathcal{M}_1(E)$  is relatively compact if and only if  $\mathcal{C}$  is tight, *i.e.*,

 $\forall \varepsilon > 0 \ \exists K \subset E \ compact, \ s.t. \ \sup_{\mu \in \mathcal{C}} \mu(E \backslash K) \leq \varepsilon.$ 

Note that since sets consisting of a single point are always compact, Proposition 1.19 implies that every probability measure (and therefore also every finite measure) on a Polish space E has the property that for all  $\varepsilon > 0$  there exists a compact K such that  $\mu(E \setminus K) \leq \varepsilon$ . This fact in itself is already nontrivial, since Polish spaces need in general not be locally compact.

By definition, a set of functions  $\mathcal{D} \subset \mathcal{C}_b(E)$  is called *distribution determining* if for any  $\mu, \nu \in \mathcal{M}_1(E)$ ,

$$\int f d\mu = \int f d\nu \quad \forall f \in \mathcal{D} \quad \text{implies} \quad \mu = \nu.$$

We say that a sequence of probability measures  $(\mu_n)_{n\geq 1}$  is *tight* if the set  $\{\mu_n : n \geq 1\}$  is tight, i.e.,  $\forall \varepsilon > 0$  there exists a compact K such that  $\sup_n \mu_n(E \setminus K) \leq \varepsilon$ . By Prohorov's theorem, each tight sequence of probability measures has a convergent subsequence. This fact is often applied as in the following lemma.

**Lemma 1.20 (Tight sequences)** Let E be a Polish space and let  $\mu_n, \mu$  be probability measures on E. Assume that  $\mathcal{D} \subset \mathcal{C}_b(E)$  is distribution determining. Then one has  $\mu_n \Rightarrow \mu$  if and only if the following two conditions are satisfied:

- (i) The sequence  $(\mu_n)_{n>1}$  is tight.
- (ii)  $\int f d\mu_n \to \int f d\mu$  for all  $f \in \mathcal{D}$ .

**Proof** In any metrizable space, if  $(x_n)_{n\geq 1}$  is a convergent sequence, then  $\{x_n : n \geq 1\}$  is relatively compact. Thus, by Prohorov's theorem, conditions (i) and (ii) are clearly necessary.

Now assume that (i) and (ii) are satisfied but  $\mu_n \neq \mu$ . Then we can find some  $f \in \mathcal{C}_b(E)$  such that  $\int f d\mu_n \neq \int f d\mu$ . It follows that we can find some  $\varepsilon > 0$  and  $n(m) \to \infty$  such that

$$\left|\int f \mathrm{d}\mu_{n(m)} - \int f \mathrm{d}\mu\right| \ge \varepsilon \qquad \forall m \ge 1.$$
(1.9)
Since the  $(\mu_{n(m)})_{m\geq 1}$  are tight, we can select a further subsequence  $\tilde{n}(m)$  and probability measure  $\mu'$  such that  $\mu_{\tilde{n}(m)} \Rightarrow \mu'$ . By condition (ii), we have  $\mu' = \mu$ . It follows that

$$\int f \mathrm{d}\mu_{\tilde{n}(m)} \xrightarrow[m \to \infty]{} \int f \mathrm{d}\mu_{\tilde{n}(m)}$$

contradicting (1.9).

## **1.8** LDP's on compact spaces

Our aim is to prove an analogue of Lemma 1.20 for large deviation principles. To prepare for this, in the present section, we will study large deviation principles on compact spaces. The results in this section will also shed some light on some elements of the theory that have up to now not been very well motivated, such as why rate functions are lower semi-continuous.

It is well-known that a compact metrizable space is separable, and complete in any metric that generates the topology. In particular, all compact metrizable spaces are Polish. Note that if E is a compact metrizable space, then  $\mathcal{C}(E) = \mathcal{C}_b(E)$ , i.e., continuous functions are automatically bounded. We equip  $\mathcal{C}(E)$  with the supremumnorm  $\|\cdot\|_{\infty}$ , under which it is a separable Banach space.<sup>2</sup> Below, |f|denotes the absolute value of a function, i.e., the function  $x \mapsto |f(x)|$ .

**Proposition 1.21 (Generalized supremumnorms)** Let *E* be a compact metrizable space and let  $\Lambda : \mathcal{C}(E) \to [0, \infty)$  be a function such that

- (i)  $\Lambda$  is a seminorm.
- (ii)  $\Lambda(f) = \Lambda(|f|)$  for all  $f \in \mathcal{C}(E)$ .
- (iii)  $\Lambda(f) \leq \Lambda(g)$  for all  $f, g \in \mathcal{C}_+(E), f \leq g$ .
- (iv)  $\Lambda(f \lor g) = \Lambda(f) \lor \Lambda(g)$  for all  $f, g \in \mathcal{C}_+(E)$ .

<sup>&</sup>lt;sup>2</sup>The separability of  $\mathcal{C}(E)$  is an easy consequence of the Stone-Weierstrass theorem [Dud02, Thm 2.4.11]. Let  $\mathcal{D} \subset E$  be dense and let  $\mathcal{A} := \{\phi_{n,x} : x \in \mathcal{D}, n \geq 1\}$ , where  $\phi_{\delta,x}(y) := 0 \vee (1 - nd(x, y))$ . Let  $\mathcal{B}$  be the set containing the function that is identically 1 and all functions of the form  $f_1 \cdots f_m$  with  $m \geq 1$  and  $f_1, \ldots, f_m \in \mathcal{A}$ . Let  $\mathcal{C}$  be the linear span of  $\mathcal{B}$  and let  $\mathcal{C}'$  be the set of functions of the form  $a_1f_1 + \cdots + a_mf_m$  with  $m \geq 1, a_1, \ldots, a_m \in \mathbb{Q}$  and  $f_1, \ldots, f_m \in \mathcal{B}$ . Then  $\mathcal{C}$  is an algebra that separates points, hence by the Stone-Weierstrass theorem,  $\mathcal{C}$  is dense in  $\mathcal{C}(E)$ . Since  $\mathcal{C}'$  is dense in  $\mathcal{C}'$  and  $\mathcal{C}'$  is countable, it follows that  $\mathcal{C}(E)$  is separable.

Then

(a)  $\Lambda : \mathcal{C}(E) \to [0, \infty)$  is continuous w.r.t. the supremumnorm.

Moreover, there exits a function  $I: E \to (-\infty, \infty]$  such that

- (b)  $\Lambda(f_n) \downarrow e^{-I(x)}$  for any  $f_n \in \mathcal{C}_+(E)$  s.t.  $f_n \downarrow \mathbb{1}_{\{x\}}$ .
- (c) I is lower semi-continuous.
- (d)  $\Lambda(f) = \sup_{x \in E} e^{-I(x)} |f(x)| \qquad (f \in \mathcal{C}(E)).$

**Proof** To prove part (a), we observe that by (ii), (iii) and (i)

$$\Lambda(f) = \Lambda(|f|) \le \Lambda(||f||_{\infty} \cdot 1) = ||f||_{\infty}\Lambda(1),$$

where  $1 \in \mathcal{C}(E)$  denotes the function that is identically one. Using again that  $\Lambda$  is a seminorm, we see that

$$|\Lambda(f) - \Lambda(g)| \le \Lambda(f - g) \le \Lambda(1) ||f - g||_{\infty}.$$

This shows that  $\Lambda$  is continuous w.r.t. the supremumnorm.

Next, define  $I: E \to (-\infty, \infty]$  (or equivalently  $e^{-I}: E \to [0, \infty)$ ) by

$$e^{-I(x)} := \inf\{\Lambda(f) : f \in \mathcal{C}_+(E), \ f(x) = 1\}$$
  $(x \in E)$ 

We claim that this function satisfies the properties (b)–(d). Indeed, if  $f_n \in \mathcal{C}_+(E)$  satisfy  $f_n \downarrow 1_{\{x\}}$  for some  $x \in E$ , then the  $\Lambda(f_n)$  decrease to a limit by the monotonicity of  $\Lambda$ . Since

$$\Lambda(f_n) \ge \Lambda(f_n/f_n(x)) \ge \inf\{\Lambda(f) : f \in \mathcal{C}_+(E), \ f(x) = 1\} = e^{-I(x)}$$

we see that this limit is larger or equal than  $e^{-I(x)}$ . To prove the other inequality, we note that by the definition of I, for each  $\varepsilon > 0$  we can choose  $f \in \mathcal{C}_+(E)$ with f(x) = 1 and  $\Lambda(f) \leq e^{-I(x)} + \varepsilon$ . We claim that there exists an n such that  $f_n < (1 + \varepsilon)f$ . Indeed, this follows from the fact that the the sets  $C_n :=$  $\{y \in E : f_n(y) \geq (1 + \varepsilon)f(y)\}$  are compact sets decreasing to the empty set, hence  $C_n = \emptyset$  for some n [Eng89, Corollary 3.1.5]. As a result, we obtain that  $\Lambda(f_n) \leq (1 + \varepsilon)\Lambda(f) \leq (1 + \varepsilon)(e^{-I(x)} + \varepsilon)$ . Since  $\varepsilon > 0$  is arbitrary, this completes the proof of property (b).

#### 1.8. LDP'S ON COMPACT SPACES

To prove part (c), consider the functions

$$\phi_{\delta,y}(x) := 0 \lor (1 - d(y, x)/\delta) \qquad (x, y \in E, \ \delta > 0).$$

Observe that  $\phi_{\delta,y}(y) = 1$  and  $\phi_{\delta,y} = 0$  on  $B_{\delta}(y)^c$ , and recall from Lemma 1.5 that  $\phi_{\delta,y}: E \to [0,1]$  is continuous. Since

$$\|\phi_{\delta,y} - \phi_{\delta,z}\|_{\infty} \le \delta^{-1} \sup_{x \in E} |d(x,y) - d(x,z)| \le \delta^{-1} d(y,z),$$

we see that the map  $x \mapsto \phi_{\delta,x}$  is continuous w.r.t. the supremumnorm. By part (a), it follows that for each  $\delta > 0$ , the functions

$$x \mapsto \Lambda(\phi_{\delta,x})$$

are continuous. Since by part (b) these functions decrease to  $e^{-I}$  as  $\delta \downarrow 0$ , we conclude that  $e^{-I}$  is upper semi-continuous or equivalently I is lower semi-continuous.

To prove part (d), by assumption (ii), it suffices to consider the case that  $f \in C_+(E)$ . We start by observing that

$$e^{-I(x)} \le \Lambda(f)$$
  $\forall x \in E, f \in \mathcal{C}_+(E), f(x) = 1,$ 

hence, more generally, for any  $x \in E$  and  $f \in \mathcal{C}_+(E)$  such that f(x) > 0,

- / >

$$e^{-I(x)} \le \Lambda(f/f(x)) = \Lambda(f)/f(x),$$

which implies that

$$e^{-I(x)}f(x) \le \Lambda(f) \qquad \forall x \in E, \ f \in \mathcal{C}_+(E),$$

and therefore

$$\Lambda(f) \ge \sup_{x \in E} e^{-I(x)} f(x) \qquad (f \in \mathcal{C}_+(E)).$$

To prove the other inequality, we claim that for each  $f \in \mathcal{C}_+(E)$  and  $\delta > 0$  we can find some  $x \in E$  and  $g \in \mathcal{C}_+(E)$  supported on  $B_{2\delta}(x)$  such that  $f \geq g$  and  $\Lambda(f) = \Lambda(g)$ . To see this, consider the functions

$$\psi_{\delta,y}(x) := 0 \lor (1 - d(B_{\delta}(y), x)/\delta) \qquad (x, y \in E, \ \delta > 0).$$

Note that  $\psi_{\delta,y} : E \to [0,1]$  is continuous and equals one on  $B_{\delta}(y)$  and zero on  $B_{2\delta}(y)^c$ . Since E is compact, for each  $\delta > 0$  we can find a finite set  $\Delta \subset E$  such that  $\bigcup_{x \in \Delta} B_{\delta}(x) = E$ . By property (iv), it follows that

$$\Lambda(f) = \Lambda\big(\bigvee_{x \in \Delta} \psi_{\delta,x} f\big) = \bigvee_{x \in \Delta} \Lambda(\psi_{\delta,x} f).$$

In particular, we may choose some x such that  $\Lambda(f) = \Lambda(\psi_{\delta,x}f)$ . Continuing this process, we can find  $x_k \in E$  and  $f_k \in \mathcal{C}_+(E)$  supported on  $B_{1/k}(x_k)$  such that  $f \geq f_1 \geq f_2$  and  $\Lambda(f) = \Lambda(f_1) = \Lambda(f_2) = \cdots$ . It is not hard to see that the  $f_n$  decrease to zero except possibly in one point x, i.e.,

$$f_n \downarrow c1_{\{x\}}$$

for some  $0 \le c \le f(x)$  and  $x \in E$ . By part (b), it follows that  $\Lambda(f) = \Lambda(f_n) \downarrow ce^{-I(x)} \le f(x)e^{-I(x)}$ . This completes the proof of part (d).

Recall the definition of a normalized rate function from page 32. The following proposition prepares for Theorem 1.24 below.

**Proposition 1.22 (LDP along a subsequence)** Let *E* be a compact metrizable space, let  $\mu_n$  be probability measures on *E* and let  $s_n$  be positive constants converging to infinity. Then there exists  $n(m) \to \infty$  and a normalized rate function *I* such that the  $\mu_{n(m)}$  satisfy the large deviation principle with speed  $s_{n(m)}$  and rate function *I*.

**Proof** Since  $\mathcal{C}(E)$ , the space of continuous real functions on E, equipped with the supremumnorm, is a separable Banach space, we can choose a countable dense subset  $\mathcal{D} = \{f_k : k \geq 1\} \subset \mathcal{C}(E)$ . Using the fact that the  $\mu_n$  are probability measures, we see that

$$||f||_{s_n,\mu_n} = \left(\int |f|^{s_n} d\mu_n\right)^{1/s_n} \le \left(||f||_{\infty}^{s_n}\right)^{1/s_n} = ||f||_{\infty} \qquad (f \in \mathcal{C}(\overline{E})).$$

By Tychonoff's theorem, the product space

$$X := \bigotimes_{k=1}^{\infty} \left[ 0, \|f_k\|_{\infty} \right],$$

equipped with the product topology is compact. Therefore, we can find  $n(m) \to \infty$  such that

$$(\|f\|_{s_{n(m)},\mu_{n(m)}})_{k\geq 1}$$

converges as  $m \to \infty$  to some limit in X. In other words, this says that we can find a subsequence such that

$$\lim_{m \to \infty} \|f\|_{s_{n(m)},\mu_{n(m)}} =: \Lambda(f)$$

exists for each  $f \in \mathcal{D}$ . We claim that this implies that for the same subsequence, this limit exists in fact for all  $f \in \mathcal{C}(E)$ . To prove this, we observe that for each  $f, g \in \mathcal{C}(E)$ ,

$$\left| \|f\|_{s_n,\mu_n} - \|g\|_{s_n,\mu_n} \right| \le \|f - g\|_{s_n,\mu_n} \le \|f - g\|_{\infty}.$$

Letting  $n(m) \to \infty$  we see that also

$$|\Lambda(f) - \Lambda(g)| \le ||f - g||_{\infty} \tag{1.10}$$

for all  $f, g \in \mathcal{D}$ . Since a uniformly continuous function from one metric space into another can uniquely be extended to a continuous function from the completion of one space to the completion of the other, we see from (1.10) that  $\Lambda$  can be uniquely extended to a function  $\Lambda : \mathcal{C}(E) \to [0, \infty)$  such that (1.10) holds for all  $f, g \in \mathcal{C}(E)$ . Moreover, if  $f \in \mathcal{C}(E)$  is arbitrary and  $f_i \in \mathcal{D}$  satisfy  $||f - f_i||_{\infty} \to 0$ , then

$$\begin{aligned} \left| \|f\|_{s_{n(m)},\mu_{n(m)}} - \Lambda(f) \right| \\ &\leq \left| \|f\|_{s_{n(m)},\mu_{n(m)}} - \|f_{i}\|_{s_{n(m)},\mu_{n(m)}} \right| + \left| \|f_{i}\|_{s_{n(m)},\mu_{n(m)}} - \Lambda(f_{i}) \right| + \left| \Lambda(f_{i}) - \Lambda(f) \right| \\ &\leq \left| \|f_{i}\|_{s_{n(m)},\mu_{n(m)}} - \Lambda(f_{i}) \right| + 2\|f - f_{i}\|_{\infty}, \end{aligned}$$

hence

$$\limsup_{m \to \infty} \left| \|f\|_{s_{n(m)},\mu_{n(m)}} - \Lambda(f) \right| \le 2\|f - f_i\|_{\infty}$$

for each *i*, which proves that  $||f||_{s_{n(m)},\mu_{n(m)}} \to \Lambda(f)$ .

Our next aim is to show that the function  $\Lambda : \mathcal{C}(E) \to [0, \infty)$  satisfies properties (i)–(iv) of Proposition 1.21. Properties (i)–(iii) are satisfied by the norms  $\|\cdot\|_{s_{n(m)},\mu_{n(m)}}$  for each m, so by taking the limit  $m \to \infty$  we see that also  $\Lambda$  has these properties. To prove also property (iv), we use an argument similar to the one used in the proof of Lemma 1.9 (b). Arguing as in (1.5), we obtain

$$\begin{split} \Lambda(f \lor g) &= \lim_{m \to \infty} \|f \lor g\|_{s_{n(m)}, \mu_{n(m)}} \leq \limsup_{m \to \infty} \left( \|f\|_{s_{n(m)}, \mu_{n(m)}}^{s_{n(m)}} + \|g\|_{s_{n(m)}, \mu_{n(m)}}^{s_{n(m)}} \right)^{1/s_{n(m)}} \\ &= \left(\limsup_{m \to \infty} \|f\|_{s_{n(m)}, \mu_{n(m)}}\right) \lor \left(\limsup_{m \to \infty} \|g\|_{s_{n(m)}, \mu_{n(m)}}\right) = \Lambda(f) \lor \Lambda(g), \end{split}$$

where we have used (1.3). Since  $f, g \leq f \vee g$ , it follows from property (iii) that moreover  $\Lambda(f) \vee \Lambda(g) \leq \Lambda(f \vee g)$ , completing the proof of property (iv).

By Proposition 1.21, it follows that there exists a lower semi-continuous function  $I: E \to (-\infty, \infty]$  such that

$$\Lambda(f) = \sup_{x \in E} e^{-I(x)} |f(x)| \qquad (f \in \mathcal{C}(E)).$$

Since E is compact, I has compact level sets, i.e., I is a good rate function, hence the  $\mu_{n(m)}$  satisfy the large deviation principle with speed  $s_{n(m)}$  and rate function I. Since the  $\mu_{n(m)}$  are probability measures, it follows that I is normalized.

## **1.9** Exponential tightness

We wish to generalize Proposition 1.22 to spaces that are not compact. To do this, we need a condition whose role is similar to that of tightness in the theory of weak convergence.

Let  $\mu_n$  be a sequence of finite measures on a Polish space E and let  $s_n$  be positive contants, converging to infinity. We say that the  $\mu_n$  are *exponentially tight* with speed  $s_n$  if

$$\forall M \in \mathbb{R} \; \exists K \subset E \text{ compact, s.t. } \limsup_{n \to \infty} \frac{1}{s_n} \log \mu_n(E \setminus K) \leq -M.$$

Letting  $A^c := E \setminus A$  denote the complement of a set  $A \subset E$ , it is easy to check that exponential tightness is equivalent to the statement that

$$\forall \varepsilon > 0 \; \exists K \subset E \text{ compact, s.t. } \limsup_{n \to \infty} \| \mathbf{1}_{K^c} \|_{s_n, \mu_n} \le \varepsilon.$$

The next lemma says that exponential tightness is a necessary condition for a large deviation principle.

**Lemma 1.23 (LDP implies exponential tightness)** Let E be a Polish space and let  $\mu_n$  be finite measures on E satisfying a large deviation principle with speed  $s_n$  and good rate function I. Then the  $\mu_n$  are exponentially tight with speed  $s_n$ .

**Proof** This proof of this statement is more tricky than might be expected at first sight. We follow [DZ93, Excercise 4.1.10]. If the space E is locally compact, then an easier proof is possible, see [DZ93, 1.2.19].

Let d be a metric generating the topology on E such that (E, d) is complete, and let  $B_r(x)$  denote the open ball (w.r.t. this metric) of radius r around x. Since E is separable, we can choose a dense sequence  $(x_k)_{k\geq 1}$  in E. Then, for every  $\delta > 0$ , the open sets  $O_{\delta,m} := \bigcup_{k=1}^m B_{\delta}(x_k)$  increase to E. By Lemma 1.8 (c),  $\|1_{O_{\delta,m}^c}\|_{\infty,I} \downarrow 0$ . Thus, for each  $\varepsilon, \delta > 0$  we can choose an  $m \geq 1$  such that

$$\limsup_{n \to \infty} \| 1_{O_{\delta,m}^{c}} \|_{s_{n},\mu_{n}} \le \| 1_{O_{\delta,m}^{c}} \|_{\infty,I} \le \varepsilon.$$

### 1.9. EXPONENTIAL TIGHTNESS

In particular, for any  $\varepsilon > 0$ , we can choose  $(m_k)_{k \ge 1}$  such that

$$\limsup_{n \to \infty} \| \mathbb{1}_{O_{1/k, m_k}^c} \|_{s_n, \mu_n} \le 2^{-k} \varepsilon \qquad (k \ge 1).$$

It follows that

$$\limsup_{n \to \infty} \| 1_{\bigcup_{k=1}^{\infty} O_{1/k, m_k}^c} \|_{s_n, \mu_n} \leq \limsup_{n \to \infty} \sum_{k=1}^{\infty} \| 1_{O_{1/k, m_k}^c} \|_{s_n, \mu_n}$$
$$\leq \sum_{k=1}^{\infty} \limsup_{n \to \infty} \| 1_{O_{1/k, m_k}^c} \|_{s_n, \mu_n} \leq \sum_{k=1}^{\infty} 2^{-k} \varepsilon = \varepsilon.$$

Here

$$\bigcup_{k=1}^{\infty} O_{1/k,m_k}^{\mathbf{c}} = \left(\bigcap_{k=1}^{\infty} O_{1/k,m_k}\right)^{\mathbf{c}} = \left(\bigcap_{k=1}^{\infty} \bigcup_{l=1}^{m_k} B_{1/k}(x_l)\right)^{\mathbf{c}}.$$

Let K be the closure of  $\bigcap_{k=1}^{\infty} O_{1/k,m_k}$ . We claim that K is compact. Recall that a subset A of a metric space (E, d) is totally bounded if for every  $\delta > 0$  there exist a finite set  $\Delta \subset A$  such that  $A \subset \bigcup_{x \in \Delta} B_{\delta}(x)$ . It is well-known [Dud02, Thm 2.3.1] that a subset A of a metric space (E, d) is compact if and only if it is complete and totally bounded. In particular, if (E, d) is complete, then A is compact if and only if A is closed and totally bounded. In light of this, it suffices to show that K is totally bounded. But this is obvious from the fact that  $K \subset \bigcup_{l=1}^{m_k} B_{2/k}(x_l)$  for each  $k \geq 1$ . Since

$$\limsup_{n \to \infty} \| \mathbb{1}_{K^{c}} \|_{s_{n}, \mu_{n}} \leq \limsup_{n \to \infty} \| \mathbb{1}_{\left(\bigcap_{k=1}^{\infty} O_{1/k, m_{k}}\right)^{c}} \|_{s_{n}, \mu_{n}} \leq \varepsilon$$

and  $\varepsilon > 0$  is arbitrary, this proves the exponential tightness of the  $\mu_n$ .

The following theorem generalizes Proposition 1.22 to non-compact spaces. This result is due to O'Brian and Verwaat [OV91] and Puhalskii [Puk91]; see also the treatment in Dupuis and Ellis [DE97, Theorem 1.3.7].

**Theorem 1.24 (Exponential tightness implies LDP along a subsequence)** Let E be a Polish space, let  $\mu_n$  be probability measures on E and let  $s_n$  be positive constants converging to infinity. Assume that the  $\mu_n$  are exponentially tight with speed  $s_n$ . Then there exists  $n(m) \to \infty$  and a normalized rate function I such that the  $\mu_{n(m)}$  satisfy the large deviation principle with speed  $s_{n(m)}$  and good rate function I. We will derive Theorem 1.24 from Proposition 1.22 using compactification techniques. For this, we need to recall some general facts about compactifications of metrizable spaces.

If  $(E, \mathcal{O})$  is a topological space (with  $\mathcal{O}$  the collection of open subsets of E) and  $E' \subset E$  is any subset of E, then E' is also naturally equipped with a topology given by the collection of open subsets  $\mathcal{O}' := \{O \cap E' : O \in \mathcal{O}\}$ . This topology is called the *induced* topology from E. If  $x_n, x \in E'$ , then  $x_n \to x$  in the induced topology on E' if and only if  $x_n \to x$  in E.

If  $(E, \mathcal{O})$  is a topological space, then a *compactification* of E is a compact topological space  $\overline{E}$  such that E is a dense subset of  $\overline{E}$  and the topology on E is the induced topology from  $\overline{E}$ . If  $\overline{E}$  is metrizable, then we say that  $\overline{E}$  is a *metrizable compactification* of E. It turns out that each separable metrizable space E has a metrizable compactification [Cho69, Theorem 6.3].

A topological space E is called *locally compact* if for every  $x \in E$  there exists an open set O and compact set C such that  $x \in O \subset C$ . We cite the following proposition from [Eng89, Thms 3.3.8 and 3.3.9].

**Proposition 1.25 (Compactification of locally compact spaces)** Let E be a metrizable topological space. Then the following statements are equivalent.

- (i) E is locally compact and separable.
- (ii) There exists a metrizable compactification  $\overline{E}$  of E such that E is an open subset of  $\overline{E}$ .
- (iii) For each metrizable compactification  $\overline{E}$  of E, E is an open subset of  $\overline{E}$ .

A subset  $A \subset E$  of a topological space E is called a  $G_{\delta}$ -set if A is a countable intersection of open sets (i.e., there exist  $O_i \in \mathcal{O}$  such that  $A = \bigcap_{i=1}^{\infty} O_i$ . The following result can be found in [Bou58, §6 No. 1, Theorem. 1]. See also [Oxt80, Thms 12.1 and 12.3].

**Proposition 1.26 (Compactification of Polish spaces)** Let E be a metrizable topological space. Then the following statements are equivalent.

- (i) E is Polish.
- (ii) There exists a metrizable compactification  $\overline{E}$  of E such that E is a  $G_{\delta}$ -subset of  $\overline{E}$ .

(iii) For each metrizable compactification  $\overline{E}$  of E, E is a  $G_{\delta}$ -subset of  $\overline{E}$ .

Moreover, a subset  $F \subset E$  of a Polish space E is Polish in the induced topology if and only if F is a  $G_{\delta}$ -subset of E.

**Lemma 1.27 (Restriction principle)** Let E be a Polish space and let  $F \subset E$ be a  $G_{\delta}$ -subset of E, equipped with the induced topology. Let  $(\mu_n)_{n\geq 1}$  be finite measures on E such that  $\mu_n(E \setminus F) = 0$  for all  $n \geq 1$ , let  $s_n$  be positive constants converging to infinity and let I be a good rate function on E such that  $I(x) = \infty$  for all  $x \in E \setminus F$ . Let  $\mu_n|_F$  and  $I|_F$  denote the restrictions of  $\mu_n$  and I, respectively, to F. Then  $I|_F$  is a good rate function on F and the following statements are equivalent.

- (i) The  $\mu_n$  satisfy the large deviation principle with speed  $s_n$  and rate function I.
- (ii) The  $\mu_n|_F$  satisfy the large deviation principle with speed  $s_n$  and rate function  $I|_F$ .

**Proof** Since the level sets of I are compact in E and contained in F, they are also compact in F, hence  $I|_F$  is a good rate function. To complete the proof, by Proposition 1.7, it suffices to show that the large deviations upper and lower bounds for the  $\mu_n$  and  $\mu_n|_F$  are equivalent. A subset of F is open (resp. closed) in the induced topology if and only if it is of the form  $O \cap F$  (resp.  $C \cap F$ ) with Oan open subset of E (resp. C a closed subset of E). The equivalence of the upper bounds now follows from the observation that for each closed  $C \subset E$ ,

$$\limsup_{n \to \infty} \frac{1}{s_n} \log \mu_n \big|_F (C \cap F) = \limsup_{n \to \infty} \frac{1}{s_n} \log \mu_n(C)$$

and

$$\inf_{x \in C} I(x) = \inf_{x \in C \cap F} I\big|_F(x).$$

In the same way, we see that the large deviations lower bounds for the  $\mu_n$  and  $\mu_n|_F$  are equivalent.

**Exercise 1.28 (Weak convergence and the induced topology)** Let E be a Polish space and let  $\overline{E}$  be a metrizable compactification of E. Let d be a metric generating the topology on  $\overline{E}$ , and denote the restriction of this metric to E also

by d. Let  $\mathcal{C}_u(E)$  denote the class of functions  $f : E \to \mathbb{R}$  that are uniformly continuous w.r.t. the metric d, i.e.,

$$\forall \varepsilon > 0 \ \exists \delta > 0 \ \text{s.t.} \ d(x, y) \le \delta \quad \text{implies} \quad |f(x) - f(y)| \le \varepsilon.$$

Let  $(\mu_n)_{n\geq 1}$  and  $\mu$  be probability measures on E. Show that the following statements are equivalent:

- (i)  $\int f d\mu_n \to \int f d\mu$  for all  $f \in \mathcal{C}_b(E)$ ,
- (ii)  $\int f d\mu_n \to \int f d\mu$  for all  $f \in \mathcal{C}_u(E)$ ,

(iii)  $\mu_n \Rightarrow \mu$  where  $\Rightarrow$  denotes weak convergence of probability measures on E,

(iv)  $\mu_n \Rightarrow \mu$  where  $\Rightarrow$  denotes weak convergence of probability measures on  $\overline{E}$ .

Hint: Identify  $\mathcal{C}_u(E) \cong \mathcal{C}(\overline{E})$  and apply Proposition 1.1.

We note that compactifications are usually not unique, i.e., it is possible to construct many different compactifications of one and the same space E. If E is locally compact (but not compact), however, then we may take E such that  $E \setminus E$  consists one a single point (usually denoted by  $\infty$ ). This one-point compactification is (up to homeomorphisms) unique. For example, the one-point compactification of  $[0,\infty)$  is  $[0,\infty]$  and the one-point compactification of  $\mathbb{R}$  looks like a circle. Another useful compactification of  $\mathbb{R}$  is of course  $\overline{\mathbb{R}} := [-\infty, \infty]$ . To see an example of a compactification of a Polish space that is not locally compact, consider the space  $E := \mathcal{M}_1(\mathbb{R})$  of probability measures on  $\mathbb{R}$ , equipped with the topology of weak convergence. A natural compactification of this space is the space  $\overline{E} := \mathcal{M}_1(\overline{\mathbb{R}})$  of probability measures on  $\overline{\mathbb{R}}$ . Note that  $\mathcal{M}_1(\mathbb{R})$  is not an open subset<sup>3</sup> of  $\mathcal{M}_1(\overline{\mathbb{R}})$ , which by Proposition 1.25 proves that  $\mathcal{M}_1(\mathbb{R})$  is not locally compact. On the other hand, since by Excercise 1.28,  $\mathcal{M}_1(\mathbb{R})$  is Polish in the induced topology, we can conclude by Proposition 1.26 that  $\mathcal{M}_1(\mathbb{R})$  must be a  $G_{\delta}$ -subset  $\mathcal{M}_1(\mathbb{R})$ . (Note that in particular, this is a very quick way of proving that  $\mathcal{M}_1(\mathbb{R})$  is a measurable subset of  $\mathcal{M}_1(\mathbb{R})$ .)

Note that in all these examples, though the *topology* on E coincides with the (induced) topology from  $\overline{E}$ , the *metrics* on E and  $\overline{E}$  may be different. Indeed, if d is a metric generating the topology on  $\overline{E}$ , then E will never be complete in this metric (unless E is compact).

<sup>&</sup>lt;sup>3</sup>Indeed  $(1-n^{-1})\delta_0 + n^{-1}\delta_\infty \in \mathcal{M}_1(\overline{\mathbb{R}}) \setminus \mathcal{M}_1(\mathbb{R})$  converge to  $\delta_0 \in \mathcal{M}_1(\mathbb{R})$  which show that the complement of  $\mathcal{M}_1(\mathbb{R})$  is not closed.

**Proof of Theorem 1.24** Let  $\overline{E}$  be a metrizable compactification of E. By Proposition 1.22, there exists  $n(m) \to \infty$  and a normalized rate function  $I : \overline{E} \to [0, \infty]$  such that the  $\mu_{n(m)}$  (viewed as probability measures on  $\overline{E}$ ) satisfy the large deviation principle with speed  $s_{n(m)}$  and rate function I.

We claim that for each  $a < \infty$ , the level set  $L_a := \{x \in E : I(x) \leq a\}$  is a compact subset of E (in the induced topology). To see this, choose  $a < b < \infty$ . By exponential tightness, there exists a compact  $K \subset E$  such that

$$\limsup_{m \to \infty} \frac{1}{s_{n(m)}} \log \mu_{n(m)}(K^{c}) \le -b.$$
(1.11)

Note that since the identity map from E into  $\overline{E}$  is continuous, and the continuous image of a compact set is compact, K is also a compact subset of  $\overline{E}$ . We claim that  $L_a \subset K$ . Assume the converse. Then we can find some  $x \in L_a \setminus K$  and open subset O of  $\overline{E}$  such that  $x \in O$  and  $O \cap K = \emptyset$ . Since the  $\mu_{n(m)}$  satisfy the LDP on  $\overline{E}$ , by Proposition 1.7 (ii),

$$\liminf_{m \to \infty} \frac{1}{s_{n(m)}} \log \mu_{n(m)}(O) \ge -\inf_{x \in O} I(x) \ge -a,$$

contradicting (1.11). This shows that  $L_a \subset K$ . Since  $L_a$  is a closed subset of  $\overline{E}$ , it follows that  $L_a$  is a compact subset of  $\overline{E}$  (in the induced topology). In particular, our arguments show that  $I(x) = \infty$  for all  $x \in \overline{E} \setminus E$ . The statement now follows from the restriction principle (Lemma 1.27) and the fact that the  $\mu_{n(m)}$  viewed as probability measures on  $\overline{E}$  satisfy the large deviation principle with speed  $s_{n(m)}$ and rate function I.

## **1.10** Applications of exponential tightness

By definition, we say that a set of functions  $\mathcal{D} \subset \mathcal{C}_{b,+}(E)$  is rate function determining if for any two normalized good rate functions I, J,

$$||f||_{\infty,I} = ||f||_{\infty,J} \quad \forall f \in \mathcal{D} \quad \text{implies} \quad I = J.$$

By combining Lemma 1.23 and Theorem 1.24, we obtain the following analogue of Lemma 1.20. Note that by Lemma 1.23, the conditions (i) and (ii) below are clearly necessary for the measures  $\mu_n$  to satisfy a large deviation principle.

**Proposition 1.29 (Conditions for LDP)** Let E be a Polish space, let  $\mu_n$  be probability measures on E, and let  $s_n$  be positive constants converging to infinity. Assume that  $\mathcal{D} \subset \mathcal{C}_{b,+}(E)$  is rate function determining and that:

- (i) The sequence  $(\mu_n)_{n\geq 1}$  is exponentially tight with speed  $s_n$ .
- (ii) The limit  $\Lambda(f) = \lim_{n \to \infty} ||f||_{s_n,\mu_n}$  exists for all  $f \in \mathcal{D}$ .

Then there exists a good rate function I on E which is uniquely characterized by the requirement that  $\Lambda(f) = ||f||_{\infty,I}$  for all  $f \in \mathcal{D}$ , and the  $\mu_n$  satisfy the large deviation principle with speed  $s_n$  and rate function I.

**Proof** By exponential tightness and Theorem 1.24, there exist  $n(m) \to \infty$  and a normalized rate function I such that the  $\mu_{n(m)}$  satisfy the large deviation principle with speed  $s_{n(m)}$  and good rate function I. It follows that

$$\Lambda(f) = \lim_{m \to \infty} \|f\|_{s_{n(m)}, \mu_{n(m)}} = \|f\|_{\infty, I} \qquad (f \in \mathcal{D}),$$

which characterizes I uniquely by the fact that  $\mathcal{D}$  is rate function determining. Now imagine that the  $\mu_n$  do not satisfy the large deviation principle with speed  $s_n$  and good rate function I. Then we can find some  $n'(m) \to \infty$  and  $g \in \mathcal{C}_{b,+}(E)$  such that

$$\left| \|g\|_{s_{n'(m)},\mu_{n'(m)}} - \|g\|_{\infty,I} \right| \ge \varepsilon \qquad (m \ge 1).$$
(1.12)

By exponential tightness, there exists a further subsequence n''(m) and good rate function J such that the  $\mu_{n''(m)}$  satisfy the large deviation principle with speed  $s_{n''(m)}$  and good rate function J. It follows that  $||f||_{\infty,J} = \Lambda(f) = ||f||_{\infty,I}$  for all  $f \in \mathcal{D}$  and therefore, by the fact that  $\mathcal{D}$  is rate function determining, J = I. Since this contradicts (1.12), we conclude that the  $\mu_n$  satisfy the large deviation principle with speed  $s_n$  and good rate function I.

A somewhat weaker version of Proposition 1.29 where  $\mathcal{D}$  is replaced by  $\mathcal{C}_{b,+}$  is known as Bryc's theorem [Bry90], which can also be found in [DZ93, Theorem 4.4.2].

In view of Proposition 1.29, we are interested in finding sufficient conditions for a set  $\mathcal{D} \subset \mathcal{C}_{b,+}$  to be rate function determining. The following simple observation is useful.

Lemma 1.30 (Sufficient conditions to be rate function determining) Let E be a Polish space,  $\mathcal{D} \subset \mathcal{C}_{b,+}(E)$ , and assume that for each  $x \in E$  there exist  $f_k \in \mathcal{D}$  such that  $f_k \downarrow 1_{\{x\}}$ . Then  $\mathcal{D}$  is rate function determining. **Proof** If  $f_k \downarrow 1_{\{x\}}$ , then, by Lemma 1.8,  $||f_k||_{\infty,I} \downarrow ||1_{\{x\}}||_{\infty,I} = e^{-I(x)}$ .

Proposition 1.29 shows that in the presence of exponential tightness, it is possible to prove large deviation principles by showing that the limit  $\lim_{n\to\infty} ||f||_{s_n,\mu_m}$  exists for sufficiently many continuous functions f. Often, it is more convenient to prove that the large deviations upper and lower bounds from Proposition 1.7 hold for sufficiently many closed and open sets.

Let  $\mathcal{A}$  be a collection of measurable subsets of some Polish space E. We say that  $\mathcal{A}$  is *rate function determining* if for any pair I, J of normalized good rate functions on E, the condition

$$\inf_{x \in \overline{A}} I(x) \le \inf_{x \in \text{int}(A)} J(x) \quad \forall A \in \mathcal{A}$$
(1.13)

implies that  $I \leq J$ . A set  $\mathcal{O}' \subset \mathcal{O}$  is a *basis for the topology* if every  $O \in \mathcal{O}$  can be written as a (possibly uncountable) union of sets in  $\mathcal{O}'$ . Equivalently, this says that for each  $x \in E$  and open set  $O \ni x$ , there exists some  $O' \in \mathcal{O}'$  such that  $x \in O' \subset O$ . For example, in any metric space, the open balls form a basis for the topology.

**Lemma 1.31 (Rate function determining sets)** Let  $\mathcal{A}$  be a collection of measurable subsets of a Polish space E. Assume that  $\{int(A) : A \in \mathcal{A}\}$  is a basis for the topology. Then  $\mathcal{A}$  is rate function determining.

**Proof** Choose  $\varepsilon_k \downarrow 0$ . Since  $\{ \operatorname{int}(A) : A \in \mathcal{A} \}$  is a basis for the topology, for each  $z \in E$  and k there exists some  $A_k \in \mathcal{A}$  such that  $z \in \operatorname{int}(A_k) \subset B_{\varepsilon_k}(z)$ . Since I is a good rate function, it assumes its minimum over  $\overline{A}_k$ , so (1.13) implies that there exist  $z_k \in \overline{A}_k$  such that  $I(z_k) \leq \operatorname{inf}_{x \in \operatorname{int}(A_k)} J(x) \leq J(z)$ . Since  $z_k \to z$ , the lower semi-continuity of I implies that  $I(z) \leq \liminf_{k \to \infty} I(z_k) \leq J(z)$ .

**Theorem 1.32 (Conditions for LDP)** Let E be a Polish space, let  $\mu_n$  be probability measures on E, let  $s_n$  be positive constants converging to infinity, let I be a normalized good rate function on E, and let  $\mathcal{A}_{up}$ ,  $\mathcal{A}_{low}$  be collections of measurable subsets of E that are rate function determining. Then the  $\mu_n$  satisfy the large deviation principle with speed  $s_n$  and rate function I if and only if the following three conditions are satisfied.

(i) 
$$\limsup_{n \to \infty} \frac{1}{s_n} \log \mu_n(A) \le -\inf_{x \in \overline{A}} I(x) \quad \forall A \in \mathcal{A}_{up},$$
  
(ii) 
$$\liminf_{n \to \infty} \frac{1}{s_n} \log \mu_n(A) \ge -\inf_{x \in int(A)} I(x) \quad \forall A \in \mathcal{A}_{low},$$

### (iii) the $\mu_n$ are exponentially tight.

**Proof** The necessity of the conditions (i)–(iii) follows from Remark 1 below Proposition 1.7 and Lemma 1.23. To prove sufficiency, we observe that by Theorem 1.24, exponential tightness implies that going to a subsequence if necessary, we can assume that the  $\mu_n$  satisfy a large deviations principle with speed  $s_n$  and some good rate function J. By the argument used in the proof of Proposition 1.29, if we can show that for each such subsequence, J = I, then it follows that the  $\mu_n$  satisfy the large deviations principle with speed  $s_n$  and rate function I.

In view of this, it suffices to show that if the  $\mu_n$  satisfy a large deviations principle with speed  $s_n$  and some good rate function J and conditions (i) and (ii) are satisfied, then J = I. Indeed, condition (i) and the large deviation principle for J imply that for any  $A \in \mathcal{A}_{up}$ ,

$$-\inf_{x\in \operatorname{int}(A)} J(x) \le \liminf_{n\to\infty} \frac{1}{s_n} \log \mu_n(\operatorname{int}(A)) \le \limsup_{n\to\infty} \frac{1}{s_n} \log \mu_n(A) \le -\inf_{x\in\overline{A}} I(x),$$

which by the assumption that  $\mathcal{A}_{up}$  is rate function determining implies that  $I \leq J$ . Similarly, using (ii) instead of (i), we find that for any  $A \in \mathcal{A}_{low}$ ,

$$-\inf_{x\in \operatorname{int}(A)} I(x) \le \liminf_{n\to\infty} \frac{1}{s_n} \log \mu_n(A) \le \limsup_{n\to\infty} \frac{1}{s_n} \log \mu_n(\overline{A}) \le -\inf_{x\in\overline{A}} J(x),$$

which by the assumption that  $\mathcal{A}_{up}$  is rate function determining implies that  $J \leq I$ .

**Remark** In Theorem 1.32, instead of assuming that  $\mathcal{A}_{low}$  is rate function determining, it suffices to assume that

$$\forall \varepsilon > 0 \text{ and } z \in E \text{ s.t. } I(z) < \infty, \ \exists A \in \mathcal{A}_{\text{low}} \text{ s.t. } z \in A \subset B_{\varepsilon}(z).$$
 (1.14)

Indeed, the proof of Lemma 1.31 shows that if (1.13) holds with I and J interchanged, and we moreover have (1.14), then  $J(z) \leq I(z)$  for all  $z \in E$  such that  $I(z) < \infty$ . Trivially, this also holds if  $I(z) = \infty$ , and the proof proceeds as before.

The next lemma shows that in Theorem 1.32, instead of assuming that  $\mathcal{A}_{up}$  is rate function determining, we can also take for  $\mathcal{A}_{up}$  the set of all compact subsets of E. If E is locally compact, then {int(K) : K compact} is a basis for the topology, so in view of Lemma 1.31 this does not add anything new. However, if E is not locally compact, then {int(K) : K compact} is never a basis for the topology. In fact, there exist Polish spaces in which every compact set has empty interior. Clearly, in such spaces, the compact sets are not rate function determining and hence the lemma below does add something new. **Lemma 1.33 (Upper bound for compact sets)** Let E be a Polish space, let  $\mu_n$  be finite measures on E, let  $s_n$  be positive constants converging to infinity, and let I be a good rate function on E. Assume that

- (i) The sequence  $(\mu_n)_{n\geq 1}$  is exponentially tight with speed  $s_n$ .
- (ii)  $\limsup_{n \to \infty} \frac{1}{s_n} \log \mu_n(K) \le -\inf_{x \in K} I(x) \qquad \forall K \ compact.$

Then

$$\limsup_{n \to \infty} \frac{1}{s_n} \log \mu_n(C) \le -\inf_{x \in C} I(x) \qquad \forall C \ closed.$$

**Remark** If  $I : E \to (-\infty, \infty]$  is lower semi-continuous and not identically  $\infty$ , but not necessarily has compact level sets, and if  $\mu_n$  are measures and  $s_n \to \infty$  constants such that

(i) 
$$\limsup_{n \to \infty} \frac{1}{s_n} \log \mu_n(K) \le -\inf_{x \in K} I(x) \quad \forall K \text{ compact.}$$
  
(ii) 
$$\liminf_{n \to \infty} \frac{1}{s_n} \log \mu_n(O) \le -\inf_{x \in O} I(x) \quad \forall O \text{ open,}$$

then one says that the  $\mu_n$  satisfy the weak large deviation principle with speed  $s_n$  and rate function I. Thus, a weak large deviation principle is basically a large deviation principle without exponential tightness. The theory of weak large deviation principles is much less elegant than for large deviation principles. For example, the contraction principle (Proposition 1.14 below) may fail for measures satisfying a weak large deviation principle.

**Proof of Lemma 1.33** By exponential tightness, for each  $M < \infty$  we can find a compact  $K \subset E$  such that

$$\limsup_{n \to \infty} \frac{1}{s_n} \log \mu_n(E \setminus K) \le -M.$$

By (1.4), it follows that, for any closed  $C \subset E$ ,

$$\limsup_{n \to \infty} \frac{1}{s_n} \log \mu_n(C) = \limsup_{n \to \infty} \frac{1}{s_n} \log \left( \mu_n(C \cap K) + \mu_n(C \setminus K) \right) \\= \left( \limsup_{n \to \infty} \frac{1}{s_n} \log \mu_n(C \cap K) \right) \lor \left( \limsup_{n \to \infty} \frac{1}{s_n} \log \mu_n(C \setminus K) \right) \\\leq - \left( M \wedge \inf_{x \in C \cap K} I(x) \right) \leq - \left( M \wedge \inf_{x \in C} I(x) \right) \xrightarrow[M \to \infty]{} - \inf_{x \in C} I(x).$$

Let E and F be sets and let  $(f_{\gamma})_{\gamma \in \Gamma}$  be a collection of functions  $f : E \to F$ . By definition, we say that  $(f_{\gamma})_{\gamma \in \Gamma}$  separates points if for each  $x, y \in E$  with  $x \neq y$ , there exists a  $\gamma \in \Gamma$  such that  $f_{\gamma}(x) \neq f_{\gamma}(y)$ . The following theorem is a sort of 'inverse' of the contraction principle, in the sense that a large deviation principle for sufficiently many image measures implies a large deviation principle for the original measures. For weak convergence, the analogous statement is that if we have a sequence  $X^{(n)}$  of discrete-time processes  $(X_i^{(n)})_{i\in\mathbb{N}}$ , then weak convergence of the finite dimensional distributions implies weak convergence in law of the processes.

**Theorem 1.34 (Projective limit)** Let E and F be Polish spaces, let  $\mu_n$  be probability measures on E, and let  $s_n$  be positive constants converging to infinity. Let  $(\psi_i)_{i \in \mathbb{N}_+}$  be continuous functions  $\psi_i : E \to F$ . For each  $m \ge 1$ , let  $\psi_m : E \to F^m$  be defined as  $\psi_m(x) = (\psi_1(x), \ldots, \psi_m(x))$   $(x \in E)$ . Assume that  $(\psi_i)_{i \in \mathbb{N}_+}$  separates points and that:

- (i) The sequence  $(\mu_n)_{n\geq 1}$  is exponentially tight with speed  $s_n$ .
- (ii) For each finite  $m \ge 1$ , there exists a good rate function  $I_m$  on  $F^m$ , equipped with the product topology, such that the measures  $\mu_n \circ \vec{\psi}_m^{-1}$  satisfy the large deviation principle with speed  $s_n$  and rate function  $I_m$ .

Then there exists a good rate function I on E which is uniquely characterized by the requirement that

$$I_m(y) = \inf_{x: \ \vec{\psi}_m(x) = y} I(x) \qquad (m \ge 1, \ y \in F^m).$$

Moreover, the measures  $\mu_n$  satisfy the large deviation principle with speed  $s_n$  and rate function I.

**Proof** Our assumptions imply that for each  $f \in \mathcal{C}_{b,+}(F^m)$ ,

$$\|f \circ \vec{\psi}_m\|_{s_n,\mu_n} = \|f\|_{s_n,\mu_n \circ \vec{\psi}_m^{-1}} \xrightarrow[n \to \infty]{} \|f\|_{\infty,I_m}.$$

We claim that the set

$$\mathcal{D} := \left\{ f \circ \vec{\psi}_m : m \ge 1, \ f \in \mathcal{C}_{b,+}(F^m) \right\}$$

is rate function determining. To see this, fix  $z \in E$  and define  $f_{i,k} \in \mathcal{D}$  by

$$f_{i,k}(x) := (1 - kd(\psi_i(x), \psi_i(z))) \vee 0 \quad (i, k \ge 1, y \in E),$$

where d is any metric generating the topology on F. We claim that

$$\mathcal{D} \ni \bigwedge_{i=1}^{m} f_{i,m} \downarrow 1_{\{z\}} \quad \text{as } m \uparrow \infty.$$

Indeed, since the  $(\psi_i)_{i \in \mathbb{N}_+}$  separate points, for each  $x \neq z$  there is an  $i \geq 1$  such that  $\psi_i(x) \neq \psi_i(z)$  and hence  $f_{i,m}(y) = 0$  for m large enough. By Lemma 1.30, it follows that  $\mathcal{D}$  is rate function determining.

Proposition 1.29 now implies that there exists a good rate function I on E such that the  $\mu_n$  satisfy the large deviation principle with speed  $s_n$  and rate function I. Moreover, I is uniquely characterized by the requirement that

$$||f \circ \vec{\psi}_m||_{\infty,I} = ||f||_{\infty,I_m} \qquad (m \ge 1, \ f \in \mathcal{C}_{b,+}(F^m)).$$
 (1.15)

 $\operatorname{Set}$ 

$$I'_m(y) := \inf_{x: \ \vec{\psi}_m(x) = y} I(x) \qquad (y \in F^m),$$

which by the contraction principle (Proposition 1.14) is a good rate function on  $F^m$ . Since

$$\|f \circ \vec{\psi}_m\|_{\infty,I} = \sup_{\substack{x \in E \\ x \in F^m}} e^{-I(x)} f(\vec{\psi}_m(x))$$
$$= \sup_{y \in F^m} e^{-\inf_{x: \ \vec{\psi}_m(x) = y} I(x)} f(y) = \|f\|_{\infty, I'_m}$$

formula (1.15) implies that  $||f||_{\infty,I'_m} = ||f||_{\infty,I_m}$  for all  $f \in \mathcal{C}_{b,+}(F^m)$ , which is in turn implies that  $I_m = I'_m$ .

The following lemma gives a more explicit expression for the rate function I from Theorem 1.34 in terms of the rate functions  $\vec{\psi}_m$ .

Lemma 1.35 (Formula for high-level rate function) In the set-up of Theorem 1.34,

$$I_m(\psi_m(x)) \uparrow I(x) \quad as \ m \uparrow \infty.$$

**Proof** We observe that

$$I_m(\vec{\psi}_m(x)) = \inf_{x' \in E: \ \vec{\psi}_m(y) = \vec{\psi}_m(x)} I(x').$$

The sets  $C_m := \{x' \in E : \vec{\psi}_m(y) = \vec{\psi}_m(x)\}$  are closed and decrease to  $\{x\}$  as  $m \uparrow \infty$  by the fact that the  $\psi_i$  separate points. Therefore, by Lemma 1.8 (c),  $\inf_{x' \in C_m} I(x') \uparrow I(x)$  as  $\uparrow \infty$ .

# Chapter 2

# Sums of i.i.d. random variables

## 2.1 The Legendre transform

In order to prepare for the proof of Cramér's theorem (Theorem 0.1), and especially Lemma 0.2, we start by studying the way the rate function I in (0.3) is defined. It turns out that I is the Legendre transform of  $\log Z$ , which we now define more generally.

For any function  $f : \mathbb{R} \to [-\infty, \infty]$ , we define the Legendre transform (sometimes also called Legendre-Fenchel transform or Fenchel-Legendre transform, to honour Fenchel who first studied the transformation for non-smooth functions) of f as

$$f^*(y) := \sup_{x \in \mathbb{R}} \left[ yx - f(x) \right] \qquad (y \in \mathbb{R}).$$



Figure 2.1: The Legendre transform.

**Exercise 2.1** For  $a \in \mathbb{R}$ , let  $l_a$  denote the linear function  $l_a(x) := ax$ , and for any function  $f : \mathbb{R} \to [-\infty, \infty]$ , define  $T_a f(x) := f(x - a)$   $(x \in \mathbb{R})$ . Show that:

- (a)  $f \leq g \Rightarrow f^* \geq g^*$ .
- (b)  $(f+c)^* = f^* c$ .
- (c)  $(f + l_a)^* = T_a f^*$
- (d)  $(T_a f)^* = f^* + l_a$ .

Recall that a function  $f : \mathbb{R} \to (-\infty, \infty]$  is convex if  $f(px_1 + (1-p)x_2) \leq pf(x_1) + (1-p)f(x_2)$  for all  $0 \leq p \leq 1$  and  $x_1, x_2 \in \mathbb{R}$ . A convex function is always continuous on the interior of the interval  $\{x \in \mathbb{R} : f(x) < \infty\}$ ; in particular, convex functions taking values in  $\mathbb{R}$  are continuous everywhere. In general, for convex functions that may take the value  $+\infty$ , it will be convenient to assume that such functions are also lower semi-continuous.

**Exercise 2.2** Show that for any function  $f : \mathbb{R} \to (-\infty, \infty]$  that is not identically  $\infty$ , the Legendre transform  $f^*$  is a convex, lower semi-continuous function  $f^* : \mathbb{R} \to (-\infty, \infty]$ , regardless of whether f is convex or lower semi-continuous or not. Hint: show that the *epigraph* of  $f^*$  is given by

$$\left\{(y,z)\in\mathbb{R}^2:z\geq f^*(y)\right\}=\bigcap_{x:\,f(x)<\infty}\overline{H}_x,$$

where  $\overline{H}_x$  denotes the closed half-space  $\overline{H}_x := \{(y, z) : z \ge yx - f(x)\}.$ 

For any convex function  $f : \mathbb{R} \to (-\infty, \infty]$ , let us write

$$\mathcal{D}_f := \{x \in \mathbb{R} : f(x) < \infty\}$$
 and  $\mathcal{U}_f := \operatorname{int}(\mathcal{D}_f),$ 

where int(A) denotes the interior of a set A. We adopt the notation

$$\partial f(x) := \frac{\partial}{\partial x} f(x)$$
 and  $\partial^2 f(x) := \frac{\partial^2}{\partial x^2} f(x)$ .

We let  $\operatorname{Conv}_{\infty}$  denote the class of convex, lower semi-continuous functions  $f : \mathbb{R} \to (-\infty, \infty]$  such that

(i)  $\mathcal{U}_f \neq \emptyset$ ,

- (ii) f is  $\mathcal{C}^{\infty}$  on  $\mathcal{U}_f$ ,
- (iii)  $\partial^2 f(x) > 0$  for all  $x \in \mathcal{U}_f$ ,
- (iv) If  $x_+ := \sup \mathcal{U}_f < \infty$ , then  $\lim_{x \uparrow x_+} \partial f(x) = \infty$ , and if  $x_- := \inf \mathcal{U}_f > -\infty$ , then  $\lim_{x \downarrow x_-} \partial f(x) = -\infty$ .

**Proposition 2.3 (Legendre transform)** Let  $f \in \text{Conv}_{\infty}$ . Then:

(a)  $f^* \in \operatorname{Conv}_{\infty}$ .

(b) 
$$f^{**} = f$$
.

- (c)  $\partial f: \mathcal{U}_f \to \mathcal{U}_{f^*}$  is a bijection, and  $(\partial f)^{-1} = \partial f^*$ .
- (d) For each  $y \in \mathcal{U}_{f^*}$ , the function  $x \mapsto yx f(x)$  assumes its unique maximum in  $x_\circ = \partial f^*(y)$ .

**Proof** Set  $\mathcal{U}_f =: (x_-, x_+)$  and

$$y_{-} := \lim_{x \to x_{-}} \partial f(x),$$
  

$$y_{+} := \lim_{x \to x_{+}} \partial f(x).$$
(2.1)

Since  $\partial^2 f > 0$ , the function  $\partial f : (x_-, x_+) \to (y_-, y_+)$  is strictly increasing, hence a bijection. It follows from assumption (iv) in the definition of  $\text{Conv}_{\infty}$  that

$$f^*(y) = \infty$$
  $(y \in \mathbb{R} \setminus [y_-, y_+])$ 

which proves that  $\mathcal{U}_{f^*} = (y_-, y_+)$ . For each  $y \in (y_-, y_+)$ , the function  $x \mapsto yx - f(x)$  assumes its maximum in a unique point  $x_\circ = x_\circ(y) \in (x_-, x_+)$ , which is characterized by the requirement

$$\frac{\partial}{\partial x} (yx - f(x)) \Big|_{x = x_{\circ}} = 0$$
  
$$\Leftrightarrow \partial f(x_{\circ}) = y.$$

In other words, this says that

$$x_{\circ}(y) = (\partial f)^{-1}(y) \qquad (y \in (y_{-}, y_{+})).$$
 (2.2)

It follows that  $x_{\circ}$  depends smoothly on  $y \in (y_{-}, y_{+})$  and hence the same is true for  $f^{*}(y) = yx_{\circ}(y) - f(x_{\circ}(y))$ . Moreover,

$$\partial f^*(y) = \frac{\partial}{\partial y} (yx_\circ(y) - f(x_\circ(y))) = x_\circ(y) + y\frac{\partial}{\partial y}x_\circ(y) - \partial f(x_\circ(y))\frac{\partial}{\partial y}x_\circ(y)$$
$$= x_\circ(y) + y\frac{\partial}{\partial y}x_\circ(y) - y\frac{\partial}{\partial y}x_\circ(y) = x_\circ(y),$$

where we have used (2.2). See Figure 2.2 for a more geometric proof of this fact. It follows that

$$\partial f^*(y) = x_{\circ}(y) = (\partial f)^{-1}(y) \qquad \left(y \in (y_-, y_+)\right),$$

which completes the proof of parts (c) and (d).

We next wish to show that the double Legendre transform  $f^{**} = (f^*)^*$  equals fon  $\mathcal{U}_f$ . Let  $x_o \in (x_-, x_+)$  and  $y_o \in (y_-, y_+)$  be related by  $x_o = \partial f^*(y_o)$ , and hence  $y_o = \partial f(x_o)$ . Then the function  $y \mapsto x_o y - f^*(y)$  assumes its maximum in  $y = y_o$ , hence

$$f^{**}(x_{\circ}) = x_{\circ}y_{\circ} - f^{*}(y_{\circ}) = f(x_{\circ}),$$

(see Figure 2.1). This proves that  $f^{**} = f$  on  $\mathcal{U}_f$ . We have already seen that  $f^{**} = \infty = f$  on  $\mathbb{R} \setminus [x_-, x_+]$ , so by symmetry, it remains to show that  $f^{**}(x_+) = f(x_+)$  if  $x_+ < \infty$ . But this follows from the fact that  $\lim_{x \uparrow x_+} f(x) = x_+$  by the lower semicontinuity of f, and the same holds for  $f^{**}$  by Excercise 2.2. This completes the proof of part (b).

It remains to prove part (a). We have already shown that  $f^*$  is infinitely differentiable on  $\mathcal{U}_{f^*} \neq \emptyset$ . Since  $\partial f : (x_-, x_+) \to (y_-, y_+)$  is strictly increasing, the same is true for its inverse, which proves that  $\partial^2 f^* > 0$  on  $\mathcal{U}_{f^*}$ . Finally, by assumption (iv) in the definition of  $\operatorname{Conv}_{\infty}$ , we can have  $y_+ < \infty$  only if  $x_+ = \infty$ , hence  $\lim_{y \uparrow y_+} \partial f^*(y) = x_+ = \infty$  in this case. By symmetry, an analogue statement holds for  $y_-$ .

**Remark** Proposition 2.3 can be generalized to general convex, lower semi-continuous functions  $f : \mathbb{R} \to (-\infty, \infty]$ . In particular, the statement that  $f^{**} = f$  holds in this generality, but the other statements of Proposition 2.3 need modifications. Let us say that a function  $f : \mathbb{R} \to (-\infty, \infty]$  admits a *supporting line* with slope  $a \in \mathbb{R}$  at x if there exists some  $c \in \mathbb{R}$  such that

$$l_a(x) + c = f(x)$$
 and  $l_a(x') + c \le f(x') \quad \forall x' \in \mathbb{R}.$ 

Obviously, if f is convex and continuously differentiable in x, then there is a unique supporting line at x, with slope  $\partial f(x)$ . For general convex, lower semi-continuous



Figure 2.2: Proof of the fact that  $x_{\circ}(y) = \partial f^{*}(y)$ . Note that since  $x_{\circ}(y + \varepsilon) = x_{\circ}(y) + O(\varepsilon)$ , one has  $f(x_{\circ}(y + \varepsilon)) = f(x_{\circ}(y)) + O(\varepsilon^{2})$ .



Figure 2.3: Definition of the rate function in Cramér's theorem. The functions below are derivatives of the functions above, and inverses of each other.



Figure 2.4: Legendre transform of a non-smooth function.

functions  $f : \mathbb{R} \to (-\infty, \infty]$ , we may define  $\partial f(x)$  as a 'multi-valued' function, whose values in x are the slopes of all supporting lines at x. See Figure 2.4 for an example. For any function f (not necessarily convex), let us define the *convex hull* h of f as the largest convex, lower semi-continuous function such that  $h \leq f$ . It can be shown that in general,  $f^* = h^*$  and therefore  $f^{**} = h$ . We refer to [Roc70] for details.

Lemma 2.4 (Smoothness of logarithmic moment generating function) Let  $Z(\lambda)$  be the function defined in (0.1), let  $\mu$  be the law of  $X_1$ , and for  $\lambda \in \mathbb{R}$ , let  $\mu_{\lambda}$  denote the tilted law

$$\mu_{\lambda}(\mathrm{d}x) := \frac{1}{Z(\lambda)} e^{\lambda x} \mu(\mathrm{d}x) \qquad (\lambda \in \mathbb{R}).$$

Then  $\lambda \mapsto \log Z(\lambda)$  is infinitely differentiable and

(i) 
$$\frac{\partial}{\partial \lambda} \log Z(\lambda) = \langle \mu_{\lambda} \rangle,$$
  
(ii)  $\frac{\partial^2}{\partial \lambda^2} \log Z(\lambda) = \operatorname{Var}(\mu_{\lambda})$ 
  
 $\left\{ \lambda \in \mathbb{R} \right\}$ 

where  $\langle \mu_{\lambda} \rangle$  and  $\operatorname{Var}(\mu_{\lambda})$  denote the mean and variance of  $\mu_{\lambda}$ .

#### 2.1. THE LEGENDRE TRANSFORM

**Proof** Let  $\mu$  be the common law of the i.i.d. random variables  $(X_k)_{k\geq 1}$ . We claim that  $\lambda \mapsto Z(\lambda)$  is infinitely differentiable and

$$\left(\frac{\partial}{\partial\lambda}\right)^n Z(\lambda) = \int x^n e^{\lambda x} \mu(\mathrm{d}x).$$

To justify this, we must show that the interchanging of differentiation and integral is allowed. By symmetry, it suffices to prove this for  $\lambda \geq 0$ . We observe that

$$\frac{\partial}{\partial\lambda}\int x^n e^{\lambda x}\mu(\mathrm{d}x) = \lim_{\varepsilon \to 0}\int x^n \varepsilon^{-1} (e^{(\lambda+\varepsilon)x} - e^{\lambda x})\mu(\mathrm{d}x),$$

where

$$|x|^{n}\varepsilon^{-1}\left|e^{(\lambda+\varepsilon)x}-e^{\lambda x}\right|=|x|^{n}\left|\varepsilon^{-1}\int_{\lambda}^{\lambda+\varepsilon}xe^{\kappa x}\mathrm{d}\kappa\right|\leq|x|^{n+1}e^{(\lambda+1)x}\qquad(x\in\mathbb{R},\ \varepsilon\leq1).$$

It follows from the existence of all exponential moments (assumption (0.1)) that this function is integrable, hence we may use dominated convergence to interchange the limit and integral.

It follows that

(i) 
$$\frac{\partial}{\partial\lambda}\log Z(\lambda) = \frac{\partial}{\partial\lambda}\log\int e^{\lambda x}\mu(\mathrm{d}x) = \frac{\int xe^{\lambda x}\mu(\mathrm{d}x)}{\int e^{\lambda x}\mu(\mathrm{d}x)} = \langle\mu_{\lambda}\rangle,$$
  
(ii)  $\frac{\partial^{2}}{\partial\lambda^{2}}\log Z(\lambda) = \frac{Z(\lambda)\int x^{2}e^{\lambda x}\mu(\mathrm{d}x) - (\int xe^{\lambda x}\mu(\mathrm{d}x))^{2}}{Z(\lambda)^{2}}$   
 $= \int x^{2}\mu_{\lambda}(\mathrm{d}x) - \left(\int x\mu_{\lambda}(\mathrm{d}x)\right)^{2} = \operatorname{Var}(\mu_{\lambda}).$ 
(2.3)

We next turn our attention to the proof of Lemma 0.2. See Figure 2.3 for an illustration.

**Proof of Lemma 0.2** By Lemma 2.4,  $\lambda \mapsto \log Z(\lambda)$  is an element of the function class  $\operatorname{Conv}_{\infty}$ , so by Proposition 2.3 (a) we see that  $I \in \operatorname{Conv}_{\infty}$ . This immediately proves parts (i) and (ii). It is immediate from the definition of  $Z(\lambda)$  that Z(0) = 1 and hence  $\log Z(0) = 0$ . By Proposition 2.3 (b),  $\log Z$  is the Legendre transform of I. In particular, this shows that

$$0 = \log Z(0) = \sup_{y \in \mathbb{R}} [0y - I(y)] = -\inf_{y \in \mathbb{R}} I(y),$$

proving part (iii). By Lemma 2.4,  $\partial \log Z(0) = \langle \mu \rangle =: \rho$ , which means that  $\lambda \mapsto \rho \lambda$  is a tangent to the function  $\log Z$  in the point  $\lambda = 0$ . By the concavity of  $\log Z$ , it follows that  $I(\rho) = \sup_{\lambda \in \mathbb{R}} [\rho \lambda - \log Z(\lambda)] = 0$ . By the fact that  $I \in \text{Conv}$  this implies that I assumes its unique minimum in  $\rho$ , proving part (iv). By Lemma 2.4, Excercise 2.5 below, and (2.1), it follows that  $\mathcal{U}_I = (y_-, y_+)$ , proving part (v). Since  $I \in \text{Conv}_{\infty}$ , this also proves part (vi). Part (vii) follows from the fact that, by Proposition 2.3 (c),  $\partial I : (y_-, y_+) \to \mathbb{R}$  is a bijection. The fact that I'' > 0 on  $\mathcal{U}_I$  follows from the fact that  $I \in \text{Conv}_{\infty}$ . We recall that if f is smooth and strictly increasing and f(x) = y, then  $\frac{\partial}{\partial x}f(x) = 1/(\frac{\partial}{\partial y}f^{-1}(y))$ . Therefore, Proposition 2.3 (c), the fact that  $\partial \log Z(0) = \rho$ , and Lemma 2.4 imply that  $\partial^2 I(\rho) = 1/(\partial^2 \log Z(0)) = 1/\sigma^2$ , proving part (viii). To prove part (ix), finally, by symmetry it suffices to prove the statement for  $y_+$ . If  $y_+ < \infty$ , then

$$\begin{split} e^{-I(y_{+})} &= \inf_{\lambda \in \mathbb{R}} \left[ e^{\log Z(\lambda) - y_{+}\lambda} \right] = \inf_{\lambda \in \mathbb{R}} e^{-y_{+}\lambda} Z(\lambda) \\ &= \inf_{\lambda \in \mathbb{R}} e^{-y_{+}\lambda} \int e^{\lambda y} \mu(\mathrm{d}y) = \inf_{\lambda \in \mathbb{R}} \int e^{\lambda(y - y_{+})} \mu(\mathrm{d}y) \\ &= \lim_{\lambda \to \infty} \int e^{\lambda(y - y_{+})} \mu(\mathrm{d}y) = \mu(\{y_{+}\}), \end{split}$$

which completes our proof.

**Exercise 2.5 (Maximal and minimal mean of tilted law)** Let  $\mu_{\lambda}$  be defined as in Lemma 2.4. Show that

$$\lim_{\lambda \to -\infty} \langle \mu_{\lambda} \rangle = y_{-} \quad \text{and} \quad \lim_{\lambda \to +\infty} \langle \mu_{\lambda} \rangle = y_{+},$$

where  $y_{-}, y_{+}$  are defined as in Lemma 0.2.

### 2.2 Cramér's theorem

**Proof of Theorem 0.1** By symmetry, it suffices to prove (0.2) (i). In view of the fact that  $1_{[0,\infty)}(z) \leq e^z$ , we have, for each  $y \in \mathbb{R}$  and  $\lambda \geq 0$ ,

$$\mathbb{P}\Big[\frac{1}{n}\sum_{k=1}^{n}X_{k} \ge y\Big] = \mathbb{P}\Big[\frac{1}{n}\sum_{k=1}^{n}(X_{k}-y)\ge 0\Big] = \mathbb{P}\Big[\lambda\sum_{k=1}^{n}(X_{k}-y)\ge 0\Big]$$
$$\le \mathbb{E}\Big[e^{\lambda\sum_{k=1}^{n}(X_{k}-y)}\Big] = \prod_{k=1}^{n}\mathbb{E}\Big[e^{\lambda(X_{k}-y)}\Big] = e^{-n\lambda y}\mathbb{E}\Big[e^{\lambda X_{1}}\Big]^{n}$$
$$= e^{(\log Z(\lambda) - \lambda y)n}.$$

### 2.2. CRAMÉR'S THEOREM

If  $y > \rho$ , then, by Lemma 2.4,  $\frac{\partial}{\partial \lambda} [\log Z(\lambda) - \lambda y]|_{\lambda=0} = \rho - y < 0$ , so, by the convexity of the function  $\lambda \mapsto [\log Z(\lambda) - \lambda y]$ ,

$$\inf_{\lambda \ge 0} [\log Z(\lambda) - \lambda y] = \inf_{\lambda \in \mathbb{R}} [\log Z(\lambda) - \lambda y] =: -I(y).$$

Together with our previous formula, this shows that

$$\mathbb{P}\left[\frac{1}{n}\sum_{k=1}^{n}X_{k}\geq y\right]\leq e^{-nI(y)}\qquad(y>\rho),$$

and hence, in particular,

$$\limsup_{n \to \infty} \frac{1}{n} \log \mathbb{P} \big[ T_n \ge y \big] \le -I(y) \qquad (y > \rho)$$

To estimate the limit inferior from below, we distinguish three cases. If  $y > y_+$ , then  $\mathbb{P}[T_n \ge y] = 0$  for all  $n \ge 1$  while  $I(y) = \infty$  by Lemma 0.2 (v), so (0.2) (i) is trivially fulfilled. If  $y = y_+$ , then  $\mathbb{P}[T_n \ge y] = \mathbb{P}[X_1 = y_+]^n$  while  $I(y_+) = -\log \mathbb{P}[X_1 = y_+]$  by Lemma 0.2 (ix), hence again (0.2) (i) holds.

If  $y < y_+$ , finally, then by Proposition 2.3 (c) and (d),  $I(y) = \sup_{\lambda \in \mathbb{R}} [y\lambda - \log Z(\lambda)] = y\lambda_{\circ} - \log Z(\lambda_{\circ})$ , where  $\lambda_{\circ} = (\partial \log Z)^{-1}(y)$ . In other words, recalling Lemma 2.4, this says that  $\lambda_{\circ}$  is uniquely characterized by the requirement that

$$\langle \mu_{\lambda_{\circ}} \rangle = \partial \log Z(\lambda_{\circ}) = y.$$

We observe that if  $(\hat{X}_k)_{k\geq 1}$  are i.i.d. random variables with common law  $\mu_{\lambda_o}$ , and  $\hat{T}_n := \frac{1}{n} \sum_{k=1}^n \hat{X}_k$ , then  $\lim_{n\to\infty} \mathbb{P}[\hat{T}_n \geq y] = \frac{1}{2}$  by the central limit theorem and therefore  $\lim_{n\to\infty} \frac{1}{n} \log \mathbb{P}[\hat{T}_n \geq y] = 0$ . The idea of the proof is to replace the law  $\mu$  of the  $(X_k)_{k\geq 1}$  by  $\mu_{\lambda_o}$  at an exponential cost of size I(y). More precisely, we estimate

$$\mathbb{P}[T_{n} \geq y] = \mathbb{P}[\sum_{k=1}^{n} (X_{k} - y) \geq 0] = \int \mu(dx_{1}) \cdots \int \mu(dx_{n}) \mathbf{1}\{\sum_{k=1}^{n} (x_{k} - y) \geq 0\} \\
= Z(\lambda_{\circ})^{n} \int e^{-\lambda_{\circ}x_{1}} \mu_{\lambda_{\circ}}(dx_{1}) \cdots \int e^{-\lambda_{\circ}x_{n}} \mu_{\lambda_{\circ}}(dx_{n}) \mathbf{1}\{\sum_{k=1}^{n} (x_{k} - y) \geq 0\} \\
= Z(\lambda_{\circ})^{n} e^{-n\lambda_{\circ}y} \int \mu_{\lambda_{\circ}}(dx_{1}) \cdots \int \mu_{\lambda_{\circ}}(dx_{n}) \\
\times e^{-\lambda_{\circ}\sum_{k=1}^{n} (x_{k} - y)} \mathbf{1}\{\sum_{k=1}^{n} (x_{k} - y) \geq 0\} \\
= e^{-nI(y)} \mathbb{E}[e^{-n\lambda_{\circ}(\hat{T}_{n} - y)} \mathbf{1}\{\hat{T}_{n} - y \geq 0\}].$$
(2.4)

By the central limit theorem,

$$\mathbb{P}\left[y \leq \hat{T}_n \leq y + \sigma n^{-1/2}\right] \xrightarrow[n \to \infty]{} \frac{1}{\sqrt{2\pi}} \int_0^1 e^{-z^2/2} \mathrm{d}z =: \theta > 0.$$

Since

$$\mathbb{E}\left[e^{-n\lambda_{\circ}(\hat{T}_{n}-y)} 1_{\{\hat{T}_{n}-y\geq 0\}}\right] \geq \mathbb{P}\left[y\leq \hat{T}_{n}\leq y+\sigma n^{-1/2}\right]e^{-\sqrt{n}\sigma\lambda_{\circ}},$$

this implies that

$$\liminf_{n \to \infty} \frac{1}{n} \log \mathbb{E} \left[ e^{-n\lambda_{\circ}(\hat{T}_{n} - y)} \mathbf{1}_{\{\hat{T}_{n} - y \ge 0\}} \right]$$
  
$$\geq \liminf_{n \to \infty} \frac{1}{n} \log \left( \theta e^{-\sqrt{n}\sigma\lambda_{\circ}} \right) = -\liminf_{n \to \infty} \frac{1}{n} \left( \log \theta + \sqrt{n}\sigma\lambda_{\circ} \right) = 0.$$

Inserting this into (2.4) we find that

$$\liminf_{n \to \infty} \frac{1}{n} \log \mathbb{P} \big[ T_n \ge y \big] \ge -I(y) \qquad (y > \rho).$$

**Remark** Our proof of Cramér's theorem actually shows that for any  $\rho < y < y_+$ ,

$$e^{-nI(y) - O(\sqrt{n})} \le \mathbb{P}[T_n \ge y] \le e^{-nI(y)}$$
 as  $n \to \infty$ .

Here the term of order  $\sqrt{n}$  in the lower bound comes from the central limit theorem. A simpler method to obtain a more crude lower bound is to use the weak law of large numbers instead. For each  $\lambda_* > \lambda_{\circ}$ , the calculation in (2.4) shows that

$$\mathbb{P}[T_n \ge y] = e^{-n[\lambda_* y - \log Z(\lambda_*)]} \mathbb{E}[e^{-n\lambda_* (\hat{T}_n - y)} 1_{\{\hat{T}_n - y \ge 0\}}],$$

where  $\hat{T}_n$  now denotes the mean of n i.i.d. random variables with common law  $\mu_{\lambda_*}$ , instead of  $\mu_{\lambda_\circ}$ . Let  $\varepsilon := \langle \mu_{\lambda_*} \rangle - \langle \mu_{\lambda_\circ} \rangle = \langle \mu_{\lambda_*} \rangle - y$ . By the weak law of large numbers

$$\mathbb{P}\big[y \le \hat{T}_n \le y + 2\varepsilon\big] \underset{n \to \infty}{\longrightarrow} 1.$$

Inserting this into our previous formula yields

$$\mathbb{P}[T_n \ge y] \ge e^{-n[\lambda_* y - \log Z(\lambda_*)]} e^{-n2\varepsilon\lambda_*},$$

and hence

$$\liminf_{n \to \infty} \mathbb{P}[T_n \ge y] \ge \lambda_* y - \log Z(\lambda_*) - 2\varepsilon \lambda_*.$$

Since  $\varepsilon \downarrow 0$  as  $\lambda_* \downarrow \lambda_{\circ}$ , taking the limit, we obtain that

$$\liminf_{n \to \infty} \mathbb{P}[T_n \ge y] \ge \lambda_\circ y - \log Z(\lambda_\circ) = I(y)$$

**Remark** Using Theorem 0.1, it is not hard to show that indeed, the laws  $\mathbb{P}[T_n \in \cdot]$  satisfy a large deviation principle with speed n and good rate function I. We will postpone this until we treat the multidimensional case in Theorem 2.17. Theorem 0.1 is in fact a bit stronger than the large deviation principle. Indeed, if  $y_+ < \infty$  and  $\mu(\{y_+\}) > 0$ , then the large deviation principle tells us that

$$\limsup_{n \to \infty} \mu_n([y_+, \infty)) \le - \inf_{y \in [y_+, \infty)} I(y) = -I(y_+),$$

but, as we have seen in Excercise 1.11, the complementary statement for the limit inferior does not follow from the large deviation principle since  $[y_+, \infty)$  is not an open set.

**Remark** Theorem 0.1 remains true if the assumption that  $Z(\lambda) < \infty$  for all  $\lambda \in \mathbb{R}$  is replaced by the weaker condition that  $Z(\lambda) < \infty$  for  $\lambda$  in some open interval containing the origin. See [DZ98, Section 2.2.1].

**Remark** For  $\rho < y < y_+$ , it can be shown that for fixed  $m \ge 1$ ,

$$\mathbb{P}\left[X_1 \in \mathrm{d}x_1, \dots, X_m \in \mathrm{d}x_m \mid \frac{1}{n} \sum_{k=1}^n X_k \ge y\right] \underset{n \to \infty}{\Longrightarrow} \mu_{\lambda_o}(\mathrm{d}x_1) \cdots \mu_{\lambda_o}(\mathrm{d}x_m),$$

where  $\mu_{\lambda}$  denotes a tilted law as in Lemma 2.4 and  $\lambda_{\circ}$  is defined by the requirement that  $\langle \mu_{\lambda_{\circ}} \rangle = y$ . This means that conditioned on the rare event  $\frac{1}{n} \sum_{k=1}^{n} X_k \ge y$ , in the limit  $n \to \infty$ , the random variables  $X_1, \ldots, X_n$  are approximately distributed as if they are i.i.d. with common law  $\mu_{\lambda_{\circ}}$ .

### 2.3 The multi-dimensional Legendre transform

At first sight, one might have the impression that the theory of the Legendre transform, as described in Section 2.1, is very much restricted to one dimension.

We will see shortly that this impression is incorrect. Indeed, Proposition 2.3 has an almost straightforward generalization to convex functions on  $\mathbb{R}^d$ .

We denote a vector in  $\mathbb{R}^d$  as  $x = (x(1), \ldots, x(d))$  and let

$$\langle x,y\rangle:=\sum_{i=1}x(i)y(j)\qquad (x,y\in\mathbb{R}^d)$$

denote the usual inner product. For any function  $f : \mathbb{R}^d \to [-\infty, \infty]$ , we defined the *Legendre transform* as

$$f^*(y) := \sup_{x \in \mathbb{R}^d} \left[ \langle y, x \rangle - f(x) \right] \qquad (y \in \mathbb{R}^d).$$

Recall that a function  $f : \mathbb{R}^d \to (-\infty, \infty]$  is *convex* if  $f(px_1 + (1-p)x_2) \leq pf(x_1) + (1-p)f(x_2)$  for all  $0 \leq p \leq 1$  and  $x_1, x_2 \in \mathbb{R}^d$ . By induction, this implies that

$$f\left(\sum_{k=1}^{n} p_k x_k\right) \le \sum_{k=1}^{n} p_k f(x_k)$$

for all  $x_1, \ldots, x_n \in \mathbb{R}^d$  and  $p_1, \ldots, p_n \ge 0$  such that  $\sum_{k=1}^n p_k = 1$ .

**Exercise 2.6** For  $a \in \mathbb{R}^d$ , let  $l_a$  denote the linear function  $l_a(x) := \langle a, x \rangle$ , and for any function  $f : \mathbb{R}^d \to [-\infty, \infty]$ , define  $T_a f(x) := f(x-a)$   $(x \in \mathbb{R}^d)$ . Show that:

- (a)  $f \leq g \Rightarrow f^* \geq g^*$ .
- (b)  $(f+c)^* = f^* c$ .

(c) 
$$(f+l_a)^* = T_a f^*$$

(d)  $(T_a f)^* = f^* + l_a$ .

**Exercise 2.7** Show that for any function  $f : \mathbb{R}^d \to (-\infty, \infty]$  that is not identically  $\infty$ , the Legendre transform  $f^*$  is a convex, lower semi-continuous function  $f^* : \mathbb{R}^d \to (-\infty, \infty]$ , regardless of whether f is convex or lower semi-continuous or not.

For any function  $f : \mathbb{R}^d \to (-\infty, \infty]$ , we write

$$\mathcal{D}_f := \{x \in \mathbb{R} : f(x) < \infty\}$$
 and  $\mathcal{U}_f := \operatorname{int}(\mathcal{D}_f).$ 

For any open set  $O \subset \mathbb{R}^d$ , we let

$$\partial O := \overline{O} \backslash O$$

denote the *boundary* of O. For smooth functions, we adopt the notation

$$\partial_i f(x) := \frac{\partial}{\partial x(i)} f(x)$$
 and  $\nabla f(x) := \left(\frac{\partial}{\partial x(d)} f(x), \dots, \frac{\partial}{\partial x(d)} f(x)\right).$ 

We call the function  $\mathbb{R}^d \ni x \mapsto \nabla f(x) \in \mathbb{R}^d$  the gradient of f. We note that for any  $y \in \mathbb{R}^d$ ,

$$\langle y, \nabla f(x) \rangle = \lim_{\varepsilon \to 0} \varepsilon^{-1} (f(x + \varepsilon y) - f(x))$$

is the *directional derivative* of f at x in the direction y. Likewise,

$$\frac{\partial^2}{\partial \varepsilon^2} f(x + \varepsilon y) \Big|_{\varepsilon=0} = \sum_{i=1}^d y(i) \frac{\partial}{\partial x(i)} \left( \sum_{j=1}^d y(j) \frac{\partial}{\partial x(j)} f(x) \right)$$
$$= \sum_{i,j=1}^d y(i) \partial_i \partial_j f(x) y(j) = \langle y, D^2 f(x) y \rangle,$$

where  $D^2 f(x)$  denotes the  $d \times d$  matrix

$$D_{ij}^2 f(x) := \partial_i \partial_j f(x).$$

We let  $\operatorname{Conv}_{\infty}(\mathbb{R}^d)$  denote the class of convex, lower semi-continuous functions  $f: \mathbb{R}^d \to (-\infty, \infty]$  such that

- (i)  $\mathcal{U}_f \neq \emptyset$ ,
- (ii) f is  $\mathcal{C}^{\infty}$  on  $\mathcal{U}_f$ ,
- (iii)  $\langle y, D^2 f(x)y \rangle > 0$  for all  $x \in \mathcal{U}_f, 0 \neq y \in \mathbb{R}^d$ ,
- (iv)  $|\nabla f(x_n)| \to \infty$  for any  $\mathcal{U}_f \ni x_n \to x \in \partial \mathcal{U}_f$ .

**Proposition 2.8 (Legendre transform in more dimensions)** Assume that  $f \in \text{Conv}_{\infty}(\mathbb{R}^d)$ . Then:

- (a)  $f^* \in \operatorname{Conv}_{\infty}(\mathbb{R}^d)$ .
- (b)  $f^{**} = f$ .

- (c)  $\nabla f : \mathcal{U}_f \to \mathcal{U}_{f^*}$  is a bijection, and  $(\nabla f)^{-1} = \nabla f^*$ .
- (d) For each  $y \in \mathcal{U}_{f^*}$ , the function  $x \mapsto \langle y, x \rangle f(x)$  assumes its unique maximum in  $x_\circ = \nabla f^*(y)$ .

**Proof** Let  $\mathcal{O}$  be the image of  $\mathcal{U}_f$  under the gradient mapping  $x \mapsto \nabla f(x)$ . It follows from assumption (iv) in the definition of  $\operatorname{Conv}_{\infty}(\mathbb{R}^d)$  that

$$f^*(y) = \infty$$
  $(y \in \mathbb{R}^d \setminus \overline{O}).$ 

On the other hand, by the strict concavity of f (assumption (iii)), for each  $y \in O$ , the function  $x \mapsto \langle y, x \rangle - f(x)$  assumes its maximum in a unique point  $x_{\circ} = x_{\circ}(y) \in \mathcal{U}_{f}$ , which is characterized by the requirement that

$$\nabla f(x_\circ) = y$$

This proves that  $f^*(y) < \infty$  for  $y \in \mathcal{O}$  and therefore  $\mathcal{O} = \mathcal{U}_{f^*}$ . Moreover, we see from this that the gradient map  $\mathcal{U}_f \ni x \mapsto \nabla f(x) \in \mathcal{U}_{f^*}$  is a bijection and

$$x_{\circ}(y) = (\nabla f)^{-1}(y) \qquad (y \in \mathcal{U}_{f^*}).$$

The proof of parts (c) and (d) now proceeds completely analogous to the onedimensional case. The proof that  $f^{**} = f$  is also the same, where we observe that the values of a function  $f \in \operatorname{Conv}_{\infty}(\mathbb{R}^d)$  on  $\partial \mathcal{U}_f$  are uniquely determined by the restriction of f to  $\mathcal{U}_f$  and lower semi-continity.

It remains to show that  $f^* \in \operatorname{Conv}_{\infty}(\mathbb{R}^d)$ . The fact that  $f^*$  is infinitely differentiable on  $\mathcal{U}_{f^*}$  follows as in the one-dimensional case. To prove that  $f^*$  satisfies condition (iii) in the definition of  $\operatorname{Conv}_{\infty}(\mathbb{R}^d)$ , we observe that by the fact that  $\nabla f$ and  $\nabla f^*$  are each other's inverses, we have  $y(i) = \partial_i f(\nabla f^*(y))$  and therefore, by the chain rule

$$1_{\{i=j\}} = \frac{\partial}{\partial y(j)} \partial_i f(\nabla f^*(y)) = \sum_{k=1}^d \partial_k \partial_i f(\nabla f^*(y)) \partial_j \partial_k f^*(y).$$

In matrix notation, this says that

$$1 = D^2 f(\nabla f^*(y)) D^2 f^*(y),$$

i.e., the (symmetric) matrix  $D^2 f^*(y)$  is the inverse of the strictly positive, symmetric matrix  $D^2 f(\nabla f^*(y))$ . Recall that any symmetric real matrix can be diagonalized with respect to an orthonormal basis and that such a matrix is strictly

positive if and only if all its eigenvalues are. Then  $D^2 f^*(y)$  can be diagonalized with respect to the same orthonormal basis as  $D^2 f(\nabla f^*(y))$  and its eigenvalues are the inverses of those of the latter. In particular,  $D^2 f^*(y)$  is strictly positive.

To complete the proof, we must show that  $f^*$  satisfies condition (iv) from the definition of  $\operatorname{Conv}_{\infty}(\mathbb{R}^d)$ . Choose  $\mathcal{U}_{f^*} \ni y_n \to y \in \partial \mathcal{U}_{f^*}$  and let  $x_n := \nabla f^*(y_n)$  and hence  $y_n = \nabla f(x_n)$ . By going to a subsequence if necessary, we may assume that the  $x_n$  converge, either to a finite limit or  $|x_n| \to \infty$ . (I.e., the  $x_n$  converge in the one-point compactification of  $\mathbb{R}^d$ .) If the  $x_n$  converge to a finite limit  $x \in \partial \mathcal{U}_f$ , then by assumption (iv) in the definition of  $\operatorname{Conv}_{\infty}(\mathbb{R}^d)$ , we must have  $|y_n| = |\nabla f(x_n)| \to \infty$ , contradicting our assumption that  $y_n \to y \in \partial \mathcal{U}_{f^*}$ . It follows that  $|\nabla f^*(y_n)| = |x_n| \to \infty$ , which shows that  $f^*$  satisfies condition (iv).

Lemma 2.4 also generalizes in a straightforward way to the multi-dimensional setting. For any probability measure  $\mu$  on  $\mathbb{R}^d$  which has at least finite first, respectively second moments, we let

$$\langle \mu \rangle(i) := \int \mu(\mathrm{d}x) x(i),$$
  
$$\operatorname{Cov}_{ij}(\mu) := \int \mu(\mathrm{d}x) x(i) x(j) - \left(\int \mu(\mathrm{d}x) x(i)\right) \left(\int \mu(\mathrm{d}x) x(j)\right)$$

 $(i, j = 1, \ldots, d)$  denote the mean and covariance matrix of  $\mu$ .

Lemma 2.9 (Smoothness of logarithmic moment generating function) Let  $\mu$  be a probability measure on  $\mathbb{R}^d$  and let Z be given by

$$Z(\lambda) := \int e^{\langle \lambda, x \rangle} \mu(\mathrm{d}x) \qquad (\lambda \in \mathbb{R}^d).$$
(2.5)

Assume that  $Z(\lambda) < \infty$  for all  $\lambda \in \mathbb{R}^d$  and for  $\lambda \in \mathbb{R}$ , let  $\mu_{\lambda}$  denote the tilted law

$$\mu_{\lambda}(\mathrm{d}x) := \frac{1}{Z(\lambda)} e^{\langle \lambda, x \rangle} \mu(\mathrm{d}x) \qquad (\lambda \in \mathbb{R}^d).$$
(2.6)

Then  $\lambda \mapsto \log Z(\lambda)$  is infinitely differentiable and

(i) 
$$\frac{\partial}{\partial\lambda(i)}\log Z(\lambda) = \langle \mu_{\lambda} \rangle(i),$$
  
(ii)  $\frac{\partial^{2}}{\partial\lambda(i)\partial\lambda(j)}\log Z(\lambda) = \operatorname{Cov}_{ij}(\mu_{\lambda})$ 
  
 $\left\{ \lambda \in \mathbb{R}^{d}, i, j = 1, \dots, d \right\}.$ 

**Proof** Analogue to the proof of Lemma 2.4.

We also need a multi-dimensional analogue of Lemma 0.2. We will be satified with a less detailed statement than in the one-dimensional case. Also, we will not list those properties that are immediate consequences of Proposition 2.8. We need a bit of convex analysis. For any set  $A \subset \mathbb{R}^d$ , the *convex hull* of A is defined as

$$C(A) := \Big\{ \sum_{k=1}^{n} p_k x_k : n \ge 1, \ x_1, \dots, x_n \in A, \ p_1, \dots, p_k \ge 0, \ \sum_{k=1}^{n} p_k = 1 \Big\}.$$

The closed convex hull  $\overline{\mathbb{C}}(A)$  of A is the closure of  $\mathbb{C}(A)$ . There is another characterization of the closed convex hull of a set that will be of use to us. By definition, we will call any set of the form

$$H = \{ x \in \mathbb{R}^d : \langle y, x \rangle > c \} \quad \text{resp.} \quad \overline{H} = \{ x \in \mathbb{R}^d : \langle y, x \rangle \ge c \}$$
(2.7)

with  $0 \neq y \in \mathbb{R}^d$  and  $c \in \mathbb{R}$  an *open*, resp. *closed half-space*, and we let  $\mathcal{H}_{open}$  resp.  $\mathcal{H}_{closed}$  denote the collection of all open, resp. closed half-spaces of  $\mathbb{R}^d$ . We claim that, for any set  $A \subset \mathbb{R}^d$ ,

$$\overline{\mathcal{C}}(A) = \bigcap \{ \overline{H} \in \mathcal{H}_{\text{closed}} : A \subset \overline{H} \}.$$
(2.8)

We will skip the proof of this rather inuitive but not entirely trivial fact. A formal proof may easily be deduced from [Roc70, Theorem 11.5] or [Dud02, Thm 6.2.9].

**Lemma 2.10 (Properties of the rate function)** Let  $\mu$  be a probability measure on  $\mathbb{R}^d$ . Assume that the moment generating function Z defined in (2.5) is finite for all  $\lambda \in \mathbb{R}^d$  and that

$$\langle y, \operatorname{Cov}(\mu)y \rangle > 0 \qquad (0 \neq y \in \mathbb{R}^d).$$

For  $\lambda \in \mathbb{R}^d$ , define  $\mu_{\lambda}$  as in (2.6) and let  $\langle \mu \rangle$  resp.  $\langle \mu_{\lambda} \rangle$  be the mean of  $\mu$  and  $\mu_{\lambda}$ . Let  $I : \mathbb{R}^d \to (-\infty, \infty]$  be the Legendre transform of  $\log Z$ . Then:

- (i)  $I \in \operatorname{Conv}_{\infty}(\mathbb{R}^d)$ .
- (ii)  $I(\langle \mu \rangle) = 0$  and I(y) > 0 for all  $y \neq \langle \mu \rangle$ .
- (iii) I is a good rate function.
- (iv)  $\mathcal{U}_I = \{ \langle \mu_\lambda \rangle : \lambda \in \mathbb{R}^d \}.$
- (v)  $\overline{\mathcal{U}}_I$  is the closed convex hull of support( $\mu$ ).
- (vi) For each  $y_{\circ} \in \mathcal{U}_{I}$ , the function  $\mathbb{R}^{d} \ni \lambda \mapsto \langle y_{\circ}, \lambda \rangle \log Z(\lambda)$  assumes its maximum in a unique point  $\lambda_{\circ} \in \mathbb{R}^{d}$ , which is uniquely characterized by the requirement that  $\langle \mu_{\lambda_{\circ}} \rangle = y_{\circ}$ .

**Proof** The fact that  $I \in \operatorname{Conv}_{\infty}(\mathbb{R}^d)$  follows from Proposition 2.8 and the fact that  $\log Z \in \operatorname{Conv}_{\infty}(\mathbb{R}^d)$ , which follows from Lemma 2.4 and the assumption that the matrix  $\operatorname{Cov}(\mu)$  is strictly positive. The fact that  $I(\langle \mu \rangle) = 0$  and I(y) > 0 for all  $y \neq \langle \mu \rangle$  can be proved exactly as in the finite-dimensional case. Since  $I \in \operatorname{Conv}_{\infty}(\mathbb{R}^d)$ , the function I is lower semi-continuous, while part (ii) and the convexity of I imply that the level sets of I are bounded, hence I is a good rate function.

Property (iv) is immediate from Proposition 2.8 (c) and Lemma 2.9. Proposition 2.8 (d) moreover tells us that for each  $y_{\circ} \in \mathcal{U}_{I}$ , the function  $\mathbb{R}^{d} \ni \lambda \mapsto \langle y_{\circ}, \lambda \rangle - \log Z(\lambda)$  assumes its maximum in a unique point  $\lambda_{\circ} \in \mathbb{R}^{d}$ , which is given by  $\lambda_{\circ} = \nabla I(y_{\circ})$ . By Proposition 2.8 (c), the function  $\lambda \mapsto \nabla \log Z(\lambda)$  is the inverse of  $y \mapsto \nabla I(y)$ , so the condition  $\lambda_{\circ} = \nabla I(y_{\circ})$  is equivalent to  $\nabla \log Z(\lambda_{\circ}) = y_{\circ}$ . By Lemma 2.9, this says that  $\langle \mu_{\lambda_{\circ}} \rangle = y_{\circ}$ , proving (vi).

It remains to prove (v). Since  $\operatorname{support}(\mu_{\lambda}) = \operatorname{support}(\mu)$  for all  $\lambda \in \mathbb{R}^d$ , it is easy to see that if H is an open half-space such that  $H \cap \operatorname{support}(\mu) = \emptyset$ , then  $\langle \mu_{\lambda} \rangle \notin H$ . Since by (2.8), the complement of  $\overline{\mathbb{C}}(\operatorname{support}(\mu))$  is the union of all open half-spaces that do not intersect  $\operatorname{support}(\mu)$ , this proves the inclusion  $\mathcal{U}_I \subset \overline{\mathbb{C}}(\operatorname{support}(\mu))$ .

On the other hand, if  $H = \{y \in \mathbb{R}^d : \langle \lambda, y \rangle > c\}$  is an open half-space such that  $H \cap \operatorname{support}(\mu) \neq \emptyset$ , then, in the same way as in Exercise 2.5, one can check that there exists some r > 0 large enough such that  $\langle \mu_{r\lambda} \rangle \in H$ . This proves that  $\overline{C}(\mathcal{U}_I) \supset \overline{C}(\operatorname{support}(\mu))$ . Since I is convex, so is  $\mathcal{U}_I$ , and therefore the closed convex hull of  $\mathcal{U}_I$  is just the closure of  $\mathcal{U}_I$ . Thus, we have  $\overline{\mathcal{U}}_I \supset \overline{C}(\operatorname{support}(\mu))$ , completing our proof.

In Theorem 2.17 below, we see that Cramér's theorem generalizes to the multidimensional case in a more or less straightforward manner. In particular, in the multi-dimensional case, the rate function is the function I of Lemma 2.10. Before we prove this, it will be convenient to broaden our horizon a bit and already start preparing for the proof of Boltzmann-Sanov Theorem (Theorem 0.7 from the introduction) and its generalization to infinite spaces, Sanov's Theorem (Theorem 2.18 below). We will see that Cramér's rate function I can be derived from Sanov's rate function H by the contraction principle (see Figure 2.5). This leads to new way of looking at I that will also be useful for proving Cramér's theorem.



Figure 2.5: Two levels of large deviation principles: relation between Sanov's and Cramér's theorem.

## 2.4 Relative entropy

Let E be a Polish space and let  $\mathcal{M}_1(E)$  be the space of probability measures on E, equipped with the topology of weak convergence, under which  $\mathcal{M}_1(E)$  is Polish. Recall that by the Radon-Nikodym theorem, if  $\nu, \mu \in \mathcal{M}_1(E)$ , then  $\nu$  has a density w.r.t.  $\mu$  if and only if  $\nu$  is *absolutely continuous* w.r.t.  $\mu$ , i.e.,  $\nu(A) = 0$ for all  $A \in \mathcal{B}(E)$  such that  $\mu(A) = 0$ . We denote this as  $\nu \ll \mu$  and let  $\frac{d\nu}{d\mu}$  denote the density of  $\nu$  w.r.t.  $\mu$ , which is uniquely defined up to a.s. equality w.r.t.  $\mu$ . For any  $\nu, \mu \in \mathcal{M}_1(E)$ , we define the *relative entropy*  $H(\nu|\mu)$  of  $\nu$  w.r.t.  $\mu$  as

$$H(\nu|\mu) := \begin{cases} \int \log\left(\frac{\mathrm{d}\nu}{\mathrm{d}\mu}\right) \mathrm{d}\nu = \int \frac{\mathrm{d}\nu}{\mathrm{d}\mu} \log\left(\frac{\mathrm{d}\nu}{\mathrm{d}\mu}\right) \mathrm{d}\mu & \text{if } \nu \ll \mu, \\ \infty & \text{otherwise.} \end{cases}$$

Note that if  $\nu \ll \mu$ , then a.s. equality w.r.t.  $\mu$  implies a.s. equality w.r.t.  $\nu$ , which shows that the first formula for  $H(\nu|\mu)$  is unambiguous.

The following lemma gives some more properties of the relative entropy.

**Lemma 2.11 (Properties of the relative entropy)** For each  $\mu \in \mathcal{M}_1(E)$ , the function  $H(\cdot | \mu)$  has the following properties.

(i) 
$$H(\mu|\mu) = 0$$
 and  $H(\nu|\mu) > 0$  for all  $\nu \neq \mu$ .
### 2.4. RELATIVE ENTROPY

- (ii) The map  $\mathcal{M}_1(E) \ni \nu \mapsto H(\nu|\mu)$  is convex.
- (iii)  $H(\cdot | \mu)$  is a good rate function.

**Proof** Define  $\phi : [0, \infty) \to \mathbb{R}$  by

$$\phi(r) := \begin{cases} r \log r & (r > 0), \\ 0 & (r = 0). \end{cases}$$

Then  $\phi$  is continuous at 0 and

$$\phi'(r) = \log r + 1$$
 and  $\phi''(r) = r^{-1}$   $(r > 0).$ 

In particular,  $\phi$  is strictly convex, so by Jensen's inequality

$$H(\nu|\mu) = \int \phi\left(\frac{\mathrm{d}\nu}{\mathrm{d}\mu}\right) \mathrm{d}\mu \ge \phi\left(\int \frac{\mathrm{d}\nu}{\mathrm{d}\mu} \mathrm{d}\mu\right) = 1\log 1 = 0,$$

with equality if and only if  $d\nu/d\mu$  is equal to a constant a.s. w.r.t.  $\mu$ . This proves part (i).

To prove part (ii), fix  $\nu_1, \nu_2 \in \mathcal{M}_1(E)$  and  $0 \leq p \leq 1$ . We wish to show that

$$H(p\nu_1 + (1-p)\nu_2|\mu) \ge pH(\nu_1|\mu) + (1-p)H(\nu_2|\mu).$$

If either  $\nu_1 \not\ll \mu$  or  $\nu_2 \not\ll \mu$  (or both), then the statement is obvious. Otherwise, setting  $f_i = d\nu_i/d\mu$ , we have

$$H(p\nu_1 + (1-p)\nu_2|\mu) = \int \phi(pf_1 + (1-p)f_2)d\mu$$
  

$$\geq \int (p\phi(f_1) + (1-p)\phi(f_2))d\mu = pH(\nu_1|\mu) + (1-p)H(\nu_2|\mu)$$

by the convexity of  $\phi(r) = r \log r$ .

To prove part (iii), finally, we must show that for each  $r < \infty$ , the level set

$$L_r := \left\{ \nu \in \mathcal{M}_1(E) : H(\nu|\mu) \le r \right\}$$

is a compact subset of  $\mathcal{M}_1(E)$ . Let  $L^1(\mu)$  be the Banach space consisting of all equivalence classes of w.r.t.  $\mu$  a.e. equal, absolutely integrable functions, equipped with the norm  $||f||_1 := \int |f| d\mu$ . Then, identifying a measure with its density, we have

$$\{\nu \in \mathcal{M}(E) : \nu \ll \mu\} \cong \{f \in L^1(\mu) : f \ge 0\},\$$

and we we may identify  $L_r$  with the set

$$L_r \cong \left\{ f \in L^1_+(\mu) : \int f d\mu = 1, \ \int f \log f d\mu \le r \right\}.$$

We note that  $L_r$  is convex by the convexity of  $H(\cdot | \mu)$ . We start by showing that  $L_r$  is closed with respect to the norm  $\|\cdot\|_1$ . Let  $f_n \in L_r$ ,  $f \in L^1(\mu)$  be such that  $\|f_n - f\|_1 \to 0$ . By going to a subsequence if necessary, we may assume that  $f_n \to f$  a.s. Since the function  $\phi$  is bounded from below, it follows from Fatou's lemma that

$$\int \phi(f) \mathrm{d}\mu \le \liminf_{n \to \infty} \int \phi(f_n) \mathrm{d}\mu \le r,$$

which shows that  $f \in L_r$ .

We recall that for any real Banach space  $(V, \| \cdot \|)$ , the dual  $V^*$  is the space of all continuous linear forms on V, i.e., the space of all continuous linear maps  $l: V \to \mathbb{R}$ . The weak topology on V is the weakest topology on V that makes all the maps  $\{l: l \in V^*\}$  continuous, i.e., it is the topology on V generated by the open sets  $\{l^{-1}(O): O \subset \mathbb{R} \text{ open}, l \in V^*\}$ . The weak topology is usually weaker than the norm topology on V. Some care is needed when dealing with weak topologies since they are often not metrizable.

In particular, it is known that the dual of  $L^1(\mu)$  is isomorphic to the space  $L^{\infty}(\mu)$  of equivalence classes of w.r.t.  $\mu$  a.e. equal, bounded measurable functions, equipped with the essential supremumnorm  $||f||_{\infty} := \inf\{R < \infty : |f| \leq R \text{ a.s.}\}$ . In particular, this means that the weak topology on  $L^1(\mu)$  is the weakest topology that makes the linear forms

$$f \mapsto l_g(f) := \int fg \,\mathrm{d}\mu \quad (g \in B_b(E))$$

continuous. We now need two facts from functional analysis.

- 1. Let V be a Banach space and let  $C \subset V$  be convex and norm-closed. Then C is also closed with respect to the weak topology on V.
- 2. (Dunford-Pettis) A subset  $C \subset L^1(\mu)$  is relatively compact in the weak topology if and only if C is uniformly integrable.

Here, a set C is called *relatively compact* if its closure  $\overline{C}$  is compact, and we recall that a set  $C \subset L^1(\mu)$  is *uniformly integrable* if for each  $\varepsilon > 0$  there exists a  $K < \infty$  such that

$$\sup_{f \in C} \int \mathbb{1}_{\{|f| \ge K\}} |f| \mathrm{d}\mu \le \varepsilon.$$

### 2.4. RELATIVE ENTROPY

A sufficient condition for uniform integrability is the existence of a nonnegative, increasing, convex function  $\psi : [0, \infty) \to [0, \infty)$  such that  $\lim_{r \to \infty} \psi(r)/r = \infty$  and

$$\sup_{f\in C}\int\psi(|f|)\mathrm{d}\mu<\infty.$$

(In fact, by the De la Valle-Poussin theorem, this condition is also necessary, but we will not need this deeper converse.) Applying this to  $\psi = \phi + 1$ , we see that the set  $L_r$  is relatively compact in the weak topology by the Dunford-Pettis theorem. Since, by 1.,  $L_r$  is also closed with respect to the weak topology, we conclude that  $L_r$  is compact with respect to the weak topology.

If E is any topological space with collection of open sets  $\mathcal{O}$ , and  $F \subset E$  is a subset of E, then the *induced topology* on F is the topology defined by the collection of open sets  $\mathcal{O}' := \{O \cap F : O \in \mathcal{O}\}$ . We have just seen that  $L_r$  is a compact space in the induced topology obtained by viewing  $L_r$  as a subset of  $L^1(\mu)$ , equipped with the weak topology.

Viewing  $L_r$  as a subset of  $\mathcal{M}_1(E)$ , we observe that the topology of weak convergence of probability measures on E induces on  $L_r$  the weakest topology that makes the linear forms

Recall that identifying a probability measure with its density with respect to  $\mu$ , we can view  $L_r$  both as a subset of  $\mathcal{M}_1(E)$  and as a subset of  $L^1(\mu)$ . On  $\mathcal{M}_1(E)$ , we have the topology of weak convergence of probability measures, which on  $L_r$ induces a topology that is the weakest topology that makes the linear forms

$$f \mapsto l_g(f) := \int fg \,\mathrm{d}\mu \quad (g \in \mathcal{C}_b(E)).$$

continuous. We already know that  $L_r$  is compact in the weak topology on  $L^1(\mu)$ , which is the weakest topology to make the linear forms  $l_g$  with  $g \in B_b(E)$  measurable. Since it is defined by a smaller collection of linear forms, the topology of weak convergence of probability measures is even weaker than the weak topology. We claim that as a result, every set that is compact in the weak topology is also compact in the topology of weak convergence of probability measures. Applying this to  $L_r$  then shows that it is a compact subset of  $\mathcal{M}_1(E)$ .

Indeed, if  $\mathcal{T}_1, \mathcal{T}_2$  are topologies defined by their collections of open sets  $\mathcal{O}_1, \mathcal{O}_2$ , then  $\mathcal{T}_1$  is weaker than  $\mathcal{T}_2$  if  $\mathcal{O}_1 \subset \mathcal{O}_2$ . We claim that this implies that the set of all  $\mathcal{T}_2$ -compact sets is a subset of the set of all  $\mathcal{T}_1$ -compact sets. Indeed, if a set K is  $\mathcal{T}_2$ -compact and a collection of  $\mathcal{T}_1$ -open sets covers K, then these sets are also  $\mathcal{T}_2$ -open and hence there exists a finite subcover, proving that K is also  $\mathcal{T}_1$ -compact.

**Remark** The proof of Lemma 2.11, part (iii) is quite complicated. A seemingly much shorter proof can be found in [DE97, Lemma 1.4.3 (b) and (c)], but this proof depends on a variational formula which has a rather long and complicated proof that is deferred to an appendix. The proof in [DS89, Lemma 3.2.13] also depends on a variational formula. The proof of [DZ93, Lemma 6.2.16], on the other hand, is very similar to our proof above.

The following lemma prepares for the proof that Cramér's rate function I is the contraction of Sanov's rate function  $H(\cdot | \mu)$  (see Figure 2.5). Note that if  $E = \mathbb{R}^d$  and f is the identity function, then the measures  $\mu_{\lambda}$  defined in (2.9) are the same as those defined in (2.6). The next lemma says that among all measures  $\nu$  such that  $\int f \, d\nu = \int f \, d\mu_{\lambda}$ , the measure  $\mu_{\lambda}$  stands out since it has the lowest relative entropy with respect to  $\mu$ .

**Lemma 2.12 (Minimizers of the entropy)** Let E be a Polish space, let  $\mu$  be a probability measure on E, and let  $f : E \to \mathbb{R}^d$  be a measurable function. Assume that

$$Z(\lambda) := \int e^{\langle \lambda, f(x) \rangle} \mu(\mathrm{d}x) < \infty \qquad (\lambda \in \mathbb{R}^d),$$

and that the covariance matrix of  $\mu \circ f^{-1}$  is strictly positive. Let

$$I(y) := \sup_{\lambda \in \mathbb{R}^d} \left[ \langle y, \lambda \rangle - \log Z(\lambda) \right] \qquad (y \in \mathbb{R}^d)$$

be the Legendre transform of log Z. For each  $\lambda \in \mathbb{R}^d$ , let  $\mu_{\lambda}$  be the probability measure on E defined by

$$\mu_{\lambda}(\mathrm{d}x) = \frac{1}{Z(\lambda)} e^{\langle \lambda, f(x) \rangle} \mu(\mathrm{d}x).$$
(2.9)

Let

$$\mathcal{M}_1^f(E) := \{ \nu \in \mathcal{M}_1(E) : \int |f| \, \mathrm{d}\nu < \infty \}.$$

Then, for all  $y_{\circ} \in \mathcal{U}_{I}$ , the function  $H(\cdot | \mu)$  assumes its minimum over the set  $\{\nu \in \mathcal{M}_{1}^{f}(E) : \int \nu(\mathrm{d}x)f(x) = y_{\circ}\}$  in the unique point  $\mu_{\lambda_{\circ}}$  given by the requirement that

$$\int \mu_{\lambda_{\circ}}(\mathrm{d}x)f(x) = y_{\circ}$$

Moreover, one has  $H(\mu_{\lambda_{\circ}}|\mu) = I(y_{\circ})$ .

76

#### 2.4. RELATIVE ENTROPY

**Proof** We wish to find the minimum of  $H(\cdot | \mu)$  over the set of all  $\nu \in \mathcal{M}_1^f(E)$  subject to the constraint  $\int f \, d\nu = y_\circ$ , which are really d constraints since  $y_\circ \in \mathbb{R}^d$ . We use the method of Lagrange multipliers: we first try to find the minimum of the function  $H(\nu | \mu) - \langle \lambda, \int f \, d\nu \rangle$  for general  $\lambda \in \mathbb{R}^d$ , and then try to choose  $\lambda$  in such a way that the minimizer satisfies the constraints.

We start by proving that for any  $\lambda \in \mathbb{R}^d$  and  $\nu \in \mathcal{M}_1^f(E)$ ,

$$H(\nu|\mu) \ge \int \nu(\mathrm{d}x) \langle \lambda, f(x) \rangle - \log Z(\lambda) \qquad \left(\nu \in \mathcal{M}_1^f(E)\right), \tag{2.10}$$

where equality holds for a given value of  $\lambda$  if and only if  $\nu = \mu_{\lambda}$ . The inequality is trivial if  $H(\nu|\mu) = \infty$  so we may assume that  $\nu \ll \mu$  and  $H(\nu|\mu) = \int \log(d\nu/d\mu)d\nu < \infty$ . We can split the measure  $\mu$  in an absolutely continuous and singular part w.r.t.  $\nu$ , i.e., we can find a measurable set A and nonnegative measurable function h such that  $\nu(A) = 0$  and

$$\mu(\mathrm{d}x) = 1_A(x)\mu(\mathrm{d}x) + h(x)\nu(\mathrm{d}x).$$

Weighting the measures on both sides of this equation with the density  $d\nu/d\mu$ , which is zero on A a.s. w.r.t.  $\mu$ , we see that

$$\nu(\mathrm{d}x) = \frac{\mathrm{d}\nu}{\mathrm{d}\mu}(x)h(x)\nu(\mathrm{d}x),$$

which shows that  $h(x) = (d\nu/d\mu)^{-1}$  a.s. with respect to  $\nu$ . Since  $r \mapsto \log(r)$  is a strictly concave function, Jensen's inequality gives

$$\int \nu(\mathrm{d}x)\langle\lambda, f(x)\rangle - H(\nu|\mu) = \int \nu(\mathrm{d}x) \Big(\log\left(e^{\langle\lambda, f(x)\rangle}\right) - \log\left(\frac{\mathrm{d}\nu}{\mathrm{d}\mu}(x)\right)\Big)$$
$$= \int \nu(\mathrm{d}x)\log\left(e^{\langle\lambda, f(x)\rangle}(\frac{\mathrm{d}\nu}{\mathrm{d}\mu})^{-1}(x)\right) \le \log\left(\int \nu(\mathrm{d}x)e^{\langle\lambda, f(x)\rangle}h(x)\right)$$
$$\le \log\left(\int \mu(\mathrm{d}x)e^{\langle\lambda, f(x)\rangle}\right) = \log Z(\lambda).$$

This proves (2.10). Since the logarithm is a strictly concave function, the first inequality here (which is an application of Jensen's inequality) is an equality if and only if the function  $e^{\langle \lambda, f \rangle} (\frac{d\nu}{d\mu})^{-1}$  is a.s. constant w.r.t.  $\nu$ . Since the logarithm is a strictly increasing function and  $e^{\langle \lambda, f \rangle}$  is strictly positive, the second inequality is an equality if and only if  $\mu = h\nu$ , i.e., if  $\mu \ll \nu$ . Thus, we have equality in (2.10) if and only if  $\mu \ll \nu$  and

$$\nu(\mathrm{d}x) = \frac{1}{Z} e^{\langle \lambda, f(x) \rangle} \mu(\mathrm{d}x),$$

where Z is some constant. Since  $\nu$  is a probability measure, we must have  $Z = Z(\lambda)$ .

Our arguments so far imply that for each  $y_{\circ} \in \mathbb{R}^d$ , one has

$$H(\nu|\mu) \ge \langle \lambda, y_{\circ} \rangle - \log Z(\lambda) \qquad \forall \lambda \in \mathbb{R}^{d}, \ \nu \in \mathcal{M}_{1}^{f}(E) \text{ s.t. } \int f \, \mathrm{d}\nu = y_{\circ}, \quad (2.11)$$

with equality if and only if  $\lambda$  has the property that  $\int f d\mu_{\lambda} = y_{\circ}$  and  $\nu = \mu_{\lambda}$ . To complete the proof, we must show that if  $y_{\circ} \in \mathcal{U}_{I}$ , then there exists a unique  $\lambda_{\circ}$  such that  $\int f d\mu_{\lambda_{\circ}} = y_{\circ}$ , and  $H(\mu_{\lambda_{\circ}}|\mu) = I(y_{\circ})$ .

Note that  $Z(\lambda)$  is the moment generating function of  $\mu \circ f^{-1}$ , i.e.,

$$Z(\lambda) = \int (\mu \circ f^{-1}) (\mathrm{d}x) e^{\langle \lambda, x \rangle}.$$

Moreover, the image under f of the measure  $\mu_{\lambda}$  defined in (2.9) is the measure

$$\mu_{\lambda} \circ f^{-1}(\mathrm{d}y) = \frac{1}{Z(\lambda)} e^{\langle \lambda, y \rangle} (\mu \circ f^{-1})(\mathrm{d}y),$$

i.e., this is  $(\mu \circ f^{-1})_{\lambda}$  in the notation of formula (2.6). Note that we are assuming that the the covariance matrix of  $\mu \circ f^{-1}$  is strictly positive, so Lemma 2.10 is applicable. Now, if  $y_{\circ} \in \mathcal{U}_{I}$ , then by Lemma 2.10 (vi), the supremum

$$I(y_{\circ}) = \sup_{\lambda \in \mathbb{R}^d} \left[ \langle y_{\circ}, \lambda \rangle - \log Z(\lambda) \right]$$

is attained in a unique point  $\lambda_{\circ} \in \mathbb{R}^d$  which is uniquely characterized by the requirement that  $\int f d\mu_{\lambda_{\circ}} = \langle (\mu \circ f^{-1})_{\lambda_{\circ}} \rangle = y_{\circ}$ . Comparing with (2.11), we see that  $I(y_{\circ}) = H(\mu_{\lambda_{\circ}} | \mu)$ .

We now prove that Cramér's rate function I is the contraction of Sanov's rate function  $H(\cdot | \mu)$  (see Figure 2.5).

**Proposition 2.13 (Contracted rate function)** Let E be a Polish space, let  $\mu$  be a probability measure on E, and let  $f : E \to \mathbb{R}^d$  be a measurable function. Assume that

$$Z(\lambda) := \int e^{\langle \lambda, f(x) \rangle} \mu(\mathrm{d}x) < \infty \qquad (\lambda \in \mathbb{R}^d),$$

and let

$$I(y) := \sup_{\lambda \in \mathbb{R}^d} \left[ \langle y, \lambda \rangle - \log Z(\lambda) \right] \qquad (y \in \mathbb{R}^d)$$

#### 2.4. RELATIVE ENTROPY

be the Legendre transform of  $\log Z$ . Then

$$I(y) = \inf_{\substack{\nu \in \mathcal{M}_1^f(E) \\ \int f \, \mathrm{d}\nu = y}} H(\nu \mid \mu).$$
(2.12)

**Remark** Using Lemma 2.12, it is easy to prove (2.12) when the covariance matrix of  $\mu \circ f^{-1}$  is positive and  $y \in \mathcal{U}_I$  or  $y \notin \overline{\mathcal{U}}_I$ . By going to a suitable subspace, it is easy to get rid of the condition on the covariance matrix. Thus, it only remains to prove (2.12) when y lies on the boundary of  $\mathcal{U}_I$ . This seems to be surprisingly hard. One can try to use a continuity argument,<sup>1</sup> using that both sides of (2.12) are convex and lower-semicontinuous in y. Convexity is easy, but proving lowersemicontinuity for the right-hand side seems to be hard. If f is bounded, then this follows from the (nontrivial) fact that the level sets of  $H(\cdot | \mu)$  are compact in the (non-Polish) topology of strong convergence of measures, but the general case seems hard. A different approach is to approximate  $\mu$  with other, nicer measures, for which  $\mathcal{U}_I = \mathbb{R}^d$ . Again, one runs into the problem that convergence of the right-hand side of (2.12) seems to be difficult to prove. The proof below is by brute force, explicitly identifying the unique minimizer of the right-hand side of (2.12) for each value of y where the infimum is not  $\infty$ . This proof is probably best skipped at a first reading.

**Proof of Proposition 2.13** Let us write  $Z_{\mu}(\lambda)$  and  $I_{\mu}(y)$  to make the dependence of these quantities on  $\mu$  explicit. We observe that the formulas for  $Z_{\mu}(\lambda)$ ,  $I_{\mu}(y)$ , and  $H(\nu \mid \mu)$  still make sense if  $\mu$  is a finite measure but not necessarily a probability measure. Moreover, for any nonnegative constant r, one has  $Z_{r\mu}(\lambda) = rZ_{\mu}(\lambda)$  and hence

$$I_{r\mu}(y) = I_{\mu}(y) - \log r \quad \text{and} \quad H(\nu | r\mu) = H(\nu | \mu) - \log r.$$
 (2.13)

In view of this, if (2.12) holds for probability measures  $\mu$ , then it holds more generally when  $\mu$  is a finite measure, and vice versa. We will prove the statement immediately for finite measures.

Using the scaling relations (2.13), we see that (2.11) holds more generally if  $\mu$  is a finite measure. Taking the supremum over  $\lambda \in \mathbb{R}^d$ , this implies that

$$I(y) \le \inf_{\substack{\nu \in \mathcal{M}_1^f(E) \\ \int f \, d\nu = y}} H(\nu \mid \mu).$$

 $<sup>^1\</sup>mathrm{I}$  made such a claim in the previous version of the lecture notes, but the argument I used is not correct.

To prove the opposite inequality, by the definition of I(y), we must show that there exists  $\lambda_n \in \mathbb{R}^d$  and  $\nu \in \mathcal{M}_1^f(E)$  with  $\int f \, d\nu = y$  such that

$$\langle y, \lambda_n \rangle - \log Z(\lambda_n) \xrightarrow[n \to \infty]{} H(\nu \mid \mu).$$
 (2.14)

For any finite nonzero measure  $\mu \in \mathcal{M}(E)$  and  $\lambda \in \mathbb{R}^d$ , we define  $\mu_{\lambda}$  by (2.9), which is a probability measure even if  $\mu$  is not. We have seen in (2.10) that

$$H(\mu_{\lambda} | \mu) = \int \mu_{\lambda}(\mathrm{d}x) \langle \lambda, f(x) \rangle - \log Z(\lambda).$$
(2.15)

In the proof of Lemma 2.12, we have seen that if y lies in the interior of the support of  $\mu \circ f^{-1}$ , then there exists a unique  $\lambda_{\circ} \in \mathbb{R}^d$  such that  $\int f \, d\mu_{\lambda_{\circ}} = y$ . By (2.15), we then see that (2.14) is satisfied for  $\lambda_n := \lambda_{\circ}$  and  $\nu := \mu_{\lambda_{\circ}}$ .

For general y, we have to proceed more carefully. Consider the set

$$C := \{ \int f \, \mathrm{d}\nu : \nu \in \mathcal{M}_1^f(E), \ \nu \ll \mu \}.$$

It is not hard to see that C is a convex set. For  $y \in C$ , let

$$F_y := \left\{ z \in \mathbb{R}^d : \exists \varepsilon > 0 \text{ s.t. } y - \varepsilon z \in C \text{ and } y + \varepsilon z \in C \right\}.$$

It follows from the convexity of C that  $F_y$  is a linear subspace of  $\mathbb{R}^d$  (possibly of dimension zero). For example, if C is a closed cube, then for a point y that lies in the interior of C, in the interior of a face of C, in the interior of an edge of C, or on a corner of C, the dimension of  $F_y$  is 3, 2, 1, or 0, respectively. Since C may in general be neither open nor closed, its structure can be quite complicated. For example, it is possible that the closure of C is a cube, but for a given face of this cube, only a part of the face lies inside C.

It is clear that the right-hand side of (2.12) is  $\infty$  if  $y \notin C$ . We will show that also  $I(y) = \infty$  for  $y \notin C$ . On the other hand, we will show that for each  $y \in C$ , the infimum on the right-hand side of (2.12) is attained in a unique probability measure  $\nu$ , and we will show that there exists  $\lambda_n$  such that (2.14) holds for this  $\nu$ .

Let  $L_y$  denote the affine space  $L_y := \{y + z : z \in F_y\}$ , let  $E' := f^{-1}(L_y)$  and let  $\mu'$ denote the restriction of  $\mu$  to E'. Then  $\mu' \circ f^{-1}$  is the restriction of  $\mu \circ f^{-1}$  to  $L_y$ . If  $y \in C$ , then  $\mu' \circ f^{-1}$  must be nonzero and y lies in the interior of support $(\mu' \circ f^{-1})$ , viewed as a subset of  $L_y$ . Since  $\nu \ll \mu$  and  $\int f \, d\nu = y$  imply that  $\nu \ll \mu'$ , the right-hand side of (2.12) can be rewritten as

$$\inf_{\substack{\nu \in \mathcal{M}_1^f(E') \\ \int f \, \mathrm{d}\nu = y}} H(\nu \mid \mu').$$
(2.16)

Note that  $\mu'$  may fail to be a probability measure even if  $\mu$  is one. Defining  $f': E' \to F_y$  by f'(x) := f(x) - y, we can rewrite (2.16) as

$$\inf_{\substack{\nu \in \mathcal{M}_1^f(E') \\ \int f' \, d\nu = 0}} H(\nu \mid \mu').$$
(2.17)

For each  $\lambda' \in F_y$ , we define  $Z'(\lambda') := \int e^{\langle \lambda', f'(x) \rangle} \mu'(\mathrm{d}x)$  and we define tilted measures  $\mu'_{\lambda'}$  as in (2.9). Since 0 lies in the interior of  $\mathrm{support}(\mu' \circ f'^{-1})$ , viewed as a subset of  $F_y$ , the proof of Lemma 2.12 tells us that there exists a unique  $\lambda'_o \in F_y$  such that  $\int f' \,\mathrm{d}\mu'_{\lambda'_o} = 0$ , and the infimum in (2.17) is attained in the unique point  $\nu = \mu'_{\lambda'_o}$ . By (2.15),

$$H(\mu_{\lambda_{\circ}'}' \mid \mu) = 0 - \log Z'(\lambda_{\circ}').$$

We will show that (2.14) is satisfied for  $\nu = \mu'_{\lambda'_o}$ . By a change of basis, we can without loss of generality assume that  $F_y = \{\lambda \in \mathbb{R}^d : \lambda(i) = 0 \ \forall i = d' + 1, \dots, d\}$  and that

$$C \subset \{y + z : z \in \mathbb{R}^d : z(i) \le 0 \ \forall i = d' + 1, \dots, d\}.$$
 (2.18)

In (2.14), we choose  $\lambda_n$  in such a way that

$$\lambda_n(i) = \lambda'_{\circ}(i) \quad (i = 1, \dots, d'), \qquad \lambda_n(i) \to \infty \quad (i = d' + 1, \dots, d).$$

Then

$$H(\mu_{\lambda_{\circ}'}' \mid \mu) = -\log Z'(\lambda_{\circ}') = -\log \int_{\{x: f(x) \in L_y\}} e^{\langle \lambda_{\circ}', f(x) - y \rangle} \mu(\mathrm{d}x).$$

On the other hand, the left-hand side of (2.14) can be written as

$$-\log \int_{E} e^{\langle \lambda_n, f(x) - y \rangle} \mu(\mathrm{d}x).$$

To prove (2.14), we need to show that

$$\int_{\mathbb{R}^d} e^{\langle \lambda_n, z \rangle} \mu \circ f'^{-1}(\mathrm{d} z) \xrightarrow[n \to \infty]{} \int_{F_y} e^{\langle \lambda'_o, z \rangle} \mu \circ f'^{-1}(\mathrm{d} z).$$

By (2.18), the measure  $\mu \circ f'^{-1}$  is concentrated on  $\{z \in \mathbb{R}^d : z(i) \leq 0 \ \forall i = d' + 1, \ldots, d\}$ . Since  $e^{\langle \lambda_n, z \rangle} \downarrow 0$  if z(i) > 0 for some  $i \in \{d' + 1, \ldots, d\}$ , in the limit, only the integral over  $F_y$  remains and we see that (2.14) is satisfied.

To complete the proof, we must show that  $I(y) = \infty$  for  $y \notin C$ . In this case, by a change of basis, we can without loss of generality assume that  $\mu \circ f^{-1}$  is concentrated on  $\{y + z : z(i) < 0 \ \forall i = 1, ..., d\}$ . Choosing  $\lambda_n(i) \to \infty$  for all i = 1, ..., d, setting f'(x) := f(x) - y as before, one finds that

$$\langle y, \lambda_n \rangle - \log Z(\lambda_n) = -\log \int_{\mathbb{R}^d} e^{\langle \lambda_n, z \rangle} \mu \circ f'^{-1}(\mathrm{d}z) \xrightarrow[n \to \infty]{} \infty$$

proving that  $I(y) = \infty$ 

## 2.5 Cramér's theorem in more dimensions

In the present section, we will prove the multi-dimensional version of Cramér's theorem. We first prove an abstract result for convex good rate functions on  $\mathbb{R}^d$ .

As in (2.7), we let  $\mathcal{H}_{\text{open}}(\mathbb{R}^d)$  resp.  $\mathcal{H}_{\text{closed}}(\mathbb{R}^d)$  denote the collection of all open, resp. closed half-spaces of  $\mathbb{R}^d$ . The following theorem says that if a good rate function is convex, then it suffices to check the large deviations upper bound for half-spaces.

**Theorem 2.14 (LDP with convex rate function** Let  $d \ge 1$ , let  $\mu_n$  be a sequence of probability measures on  $\mathbb{R}^d$ , let  $s_n$  be positive constants, converging to  $\infty$ , and let I be a convex, normalized, good rate function on  $\mathbb{R}^d$ . Then the  $\mu_n$  satisfy the large deviations principle with speed  $s_n$  and rate function I if and only if the following two conditions are satisfied.

- (i)  $\limsup_{n \to \infty} \frac{1}{s_n} \log \mu_n(\overline{H}) \le \inf_{x \in \overline{H}} I(x) \text{ for every closed half space } \overline{H} \subset \mathbb{R}^d,$
- (ii)  $\liminf_{n \to \infty} \frac{1}{s_n} \log \mu_n(B) \ge \inf_{x \in B} I(x)$  for every open ball  $B \subset \mathbb{R}^d$ .

**Proof** We will apply Theorem 1.32. We start by proving exponential tightness. Consider the open half-spaces

$$H_{i,R}^{-} := \{ x \in \mathbb{R}^{d} : x(i) < -R \} \text{ and } H_{i,R}^{+} := \{ x \in \mathbb{R}^{d} : x(i) > R \}$$

 $(i = 1, \ldots, d)$ . Then  $K_R := \{x \in \mathbb{R}^d : |x(i)| \le R \ \forall i = 1, \ldots, d\}$  is compact and, by Lemma 1.10,

$$\limsup_{n \to \infty} \frac{1}{s_n} \log \mu_n(\mathbb{R}^d \setminus K_R) \le \sup_{i=1,\dots,d} \sup_{\sigma \in \{-,+\}} -\inf_{x \in \overline{H}_{i,R}^\sigma} I(x)$$



Figure 2.6: Open half-plane H separating the convex sets B and  $D_r$ .

Since I has compact level sets, for each  $M \in \mathbb{R}$  we can choose R sufficiently large such that

$$\inf_{x \in \overline{H}_{i,R}^{\sigma}} I(x) \ge M \quad \forall i = 1, \dots, d, \ \sigma \in \{-, +\},$$

which proves exponential tightness.

We next show that

$$\limsup_{n \to \infty} \frac{1}{s_n} \log \mu_n(C) \le -\inf_{x \in C} I(x)$$
(2.19)

for every closed convex set  $C \subset \mathbb{R}^d$ . Set  $r := \inf_{x \in C} I(x)$ . If r = 0, then (2.19) follows from the assumption that the  $\mu_n$  are probability measures. If  $0 < r \leq \infty$ , then the convex sets  $D_r = \{x \in \mathbb{R}^d : I(x) < r\}$  and C are disjoint. By a wellknown separation theorem [Roc70, Theorem 11.3], there exists a closed half-space H such that  $L_r \cap H = \emptyset$  and  $C \subset H$  (see Figure 2.6). It follows that  $I(x) \geq r$  for all  $x \in H$  and therefore

$$\limsup_{n \to \infty} \frac{1}{s_n} \log \mu_n(C) \le \limsup_{n \to \infty} \frac{1}{s_n} \log \mu_n(H) \le -\inf_{x \in H} I(x) = -\inf_{x \in C} I(x),$$

proving (2.19). In particular, (2.19) holds for closed balls. Since the open balls form a basis for the topology, by Lemma 1.31, this shows that conditions (i) and (ii) of Theorem 1.32 are satisfied.

**Remark** If the rate function I is strictly convex, then it seems that we can replace condition (ii) of Theorem 2.14 by the condition that the large deviations lower

bound holds for all open half spaces  $H \subset \mathbb{R}^d$ . The basic idea is that by looking at differences of open and closed half spaces, using the strict convexity of I and Lemma 2.15 below, it should be possible to derive the large deviations lower bound for a collection of open sets satisfying (1.14). Filling in the details is quite technical, however, which is why we do not pursue this idea further here.

**Lemma 2.15 (Difference of sets)** Let E be a measurable space, let  $A, B \subset E$ be measurable, let  $I : E \to (-\infty, \infty]$  be a function, let  $s_n$  be positive constants such that  $s_n \to \infty$ , and let  $\mu_n$  be finite measures on E. Assume that

- (i)  $\limsup_{n \to \infty} \frac{1}{s_n} \log \mu_n(B) \le -\inf_{x \in B} I(x),$
- (ii)  $\liminf_{n \to \infty} \frac{1}{s_n} \log \mu_n(A) \ge -\inf_{x \in A} I(x).$
- (iii)  $\inf_{x \in A} I(x) < \inf_{x \in B} I(x),$

Then

$$\liminf_{n \to \infty} \frac{1}{s_n} \log \mu_n(A \setminus B) \ge -\inf_{x \in A \setminus B} I(x).$$

**Proof** Since

$$\limsup_{n \to \infty} \frac{1}{s_n} \log \left( \frac{\mu_n(B)}{\mu_n(A)} \right) = \limsup_{n \to \infty} \frac{1}{s_n} \left( \log \mu_n(B) - \log \mu_n(A) \right)$$
$$\leq \inf_{x \in A} I(x) - \inf_{x \in B} I(x) < 0,$$

we see that  $\mu_n(B)/\mu_n(A) \to 0$  exponentially fast. It follows that

$$\begin{split} \liminf_{n \to \infty} \frac{1}{s_n} \log \mu_n(A \setminus B) &= \liminf_{n \to \infty} \frac{1}{s_n} \log \left( \mu_n(A) - \mu_n(A \cap B) \right) \\ \geq \liminf_{n \to \infty} \frac{1}{s_n} \log \left( \mu_n(A) - \mu_n(B) \right) &= \liminf_{n \to \infty} \frac{1}{s_n} \log \left[ \mu_n(A) \left( 1 - \frac{\mu_n(B)}{\mu_n(A)} \right) \right] \\ &= \liminf_{n \to \infty} \frac{1}{s_n} \left[ \log \mu_n(A) + \log \left( 1 - \frac{\mu_n(B)}{\mu_n(A)} \right) \right] \\ &= \liminf_{n \to \infty} \frac{1}{s_n} \log \mu_n(A) = -\inf_{x \in A} I(x) = -\inf_{x \in A \setminus B} I(x), \end{split}$$

where in the last step we have again used that  $\inf_{x \in A} I(x) < \inf_{x \in B} I(x)$ .

To prepare for the multidimensional version of Cramér's theorem, we need one more lemma, which says that the rate functions of Cramér's theorem in different dimensions are consistent in a way that one would expect from the contraction principle. (Note that since we have not proved Cramér's theorem in more dimensions yet, we cannot use the contraction principle to prove this.)

**Lemma 2.16 (Contracted rate function)** Let X be an  $\mathbb{R}^d$ -valued random variable, let  $l : \mathbb{R}^d \to \mathbb{R}^{d'}$  be a linear function, and set X' := l(X). Let

$$Z(\lambda) := \mathbb{E}[e^{\langle \lambda, X \rangle}] \qquad (\lambda \in \mathbb{R}^d),$$
$$Z'(\lambda') := \mathbb{E}[e^{\langle \lambda', X' \rangle}] \qquad (\lambda \in \mathbb{R}^{d'})$$

be the moment generating functions of X and X', and assume that  $Z(\lambda) < \infty$  for all  $\lambda \in \mathbb{R}^d$ . Let

$$I(y) := \sup_{\lambda \in \mathbb{R}^d} \left[ \langle y, \lambda \rangle - \log Z(\lambda) \right] \qquad (y \in \mathbb{R}^d), I'(y') := \sup_{\lambda' \in \mathbb{R}^{d'}} \left[ \langle y, \lambda' \rangle - \log Z'(\lambda') \right] \qquad (y' \in \mathbb{R}^{d'}),$$

be the Legendre transforms of  $\log Z$  and  $\log Z'$ , respectively. Then

$$I'(y') := \inf_{y: \, l(y)=y'} I(y) \qquad (y' \in \mathbb{R}^{d'}).$$

**Proof** It is possible to prove this directly from the definitions of I and I', but the proof is surprisingly tricky. We give here a much faster proof based on Proposition 2.13. Applying that proposition to  $E = \mathbb{R}^d$  and  $f : \mathbb{R}^d \to \mathbb{R}^d$  the identity function, we can express the rate function I as

$$I(y) = \inf \left\{ H(\nu \mid \mu) : \nu \in \mathcal{M}_1^f(\mathbb{R}^d), \ \int x \,\nu(\mathrm{d}x) = y \right\} \qquad (y \in \mathbb{R}^d).$$

Applying Lemma 2.12 to  $E = \mathbb{R}^d$  and f = l, we can express the rate function I' as

$$I'(y') = \inf \left\{ H(\nu \mid \mu) : \nu \in \mathcal{M}_1^f(\mathbb{R}^d), \ \int l(x) \,\nu(\mathrm{d}x) = y' \right\} \qquad (y' \in \mathbb{R}^{d'}).$$

It follows that

$$I'(y') = \inf_{\substack{y: \, l(y)=y' \\ \int x \, d\nu=y}} \inf_{\substack{\nu \in \mathcal{M}_1^f(E) \\ \int x \, d\nu=y}} H(\nu \mid \mu) = \inf_{\substack{y: \, l(y)=y' \\ y \in I(y)=y' \\ I$$

as required.

**Theorem 2.17 (Multidimensional Cramér's theorem)** Let  $(X_k)_{k\geq 1}$  be i.i.d.  $\mathbb{R}^d$ -valued random variables with common law  $\mu$ . Assume that the moment generating function  $Z(\lambda)$  defined in (2.5) is finite for all  $\lambda \in \mathbb{R}^d$ . Then the probability measures

$$\mu_n := \mathbb{P}\Big[\frac{1}{n} \sum_{k=1}^n X_k \in \cdot\,\Big] \qquad (n \ge 1)$$

satisfy the large deviation principle with speed n and rate function I given by

$$I(y) := \sup_{\lambda \in \mathbb{R}^d} \left[ \langle \lambda, y \rangle - \log Z(\lambda) \right].$$

**Proof** We apply Theorem 2.14, using Cramér's original one-dimensional theorem to get the upper bound for close half-spaces.

A general closed half-space  $\overline{H}$  is of the form  $\overline{H} = \{y \in \mathbb{R}^d : l_z(y) \geq c\}$  where  $c \in \mathbb{R}, 0 \neq z \in \mathbb{R}^d$ , and  $l_z$  is the linear form  $l_z(y) := \langle z, y \rangle \ (y \in \mathbb{R}^d)$ . Then

$$\mu_n(\overline{H}) = \mathbb{P}\Big[\frac{1}{n}\sum_{k=1}^n l(X_k) \ge c\Big].$$

If  $E[X_1] \in \overline{H}$ , then  $\inf_{y \in \overline{H}} I(y) = 0$  and the large deviations upper bound is trivial since the  $\mu_n$  are probability measures. If  $E[X_1] \notin \overline{H}$  or equivalently  $\mathbb{E}[l(X_1)] < c$ , then Cramér's theorem (Theorem 0.1) tells us that

$$\lim_{n \to \infty} \mu_n(\overline{H}) = -I'(c),$$

where I' is Cramér's rate function for the i.i.d. random variables  $(l(X_k))_{k\geq 1}$ . By Lemma 2.16,

$$I'(c) = \inf_{y: l(y)=c} I(y) = \inf_{y \in \overline{H}} I(y),$$

where in the last step we have used that I is convex and assumes its minimum in the point  $E[X_1] \notin \overline{H}$ . This proves the large deviations upper bound for closed half spaces.

In a similar way, we could get the large deviations lower bound for open half spaces, which by the remark below Theorem 2.14 and the strict convexity of I is probably sufficient to prove the large deviations principle for the  $\mu_n$ . Since this approach is a bit technical, we proceed differently, adapting the corresponding part of the proof of Theorem 0.1 to the multi-dimensional setting. In order not to repeat too much, we will use a slight modification of the argument. Let B be an open ball. We first consider the case that  $B \cap \mathcal{U}_I \neq \emptyset$ . Then, for each  $\varepsilon > 0$ , we can choose  $y_{\circ} \in B \cap \mathcal{U}_I$  such that  $I(y_{\circ}) \leq \inf_{y \in B} I(y) + \varepsilon$ . By Lemma 2.10 (vi),

$$I(y_{\circ}) = \sup_{\lambda \in \mathbb{R}^d} \left[ \langle y_{\circ}, \lambda \rangle - \log Z(\lambda) \right] = \langle y_{\circ}, \lambda_{\circ} \rangle - \log Z(\lambda_{\circ}),$$

where  $\lambda_{\circ}$  is uniquely characterized by the requirement that  $y_{\circ} = \langle \mu_{\lambda_{\circ}} \rangle$ . Let  $(X_k)_{k \geq 1}$ are i.i.d. random variables with common law  $\mu_{\lambda_{\circ}}$ . Choose  $\delta > 0$  such that  $B_{\delta}(y_{\circ}) \subset B$  and  $\langle \lambda_{\circ}, y \rangle \leq \langle \lambda_{\circ}, y_{\circ} \rangle + \varepsilon$  for all  $y \in B_{\delta}(y_{\circ})$ . Then, in analogy with (2.4), we estimate

$$\mathbb{P}\Big[\frac{1}{n}\sum_{k=1}^{n}X_{k}\in B\Big] = \int \mu(\mathrm{d}x_{1})\cdots\int \mu(\mathrm{d}x_{n})\mathbf{1}_{\left\{\frac{1}{n}\sum_{k=1}^{n}x_{k}\in B\right\}}$$
$$= Z(\lambda_{\circ})^{n}\int \mu_{\lambda_{\circ}}(\mathrm{d}x_{1})\cdots\int \mu_{\lambda_{\circ}}(\mathrm{d}x_{n})e^{-n\langle\lambda_{\circ},\frac{1}{n}\sum_{k=1}^{n}x_{k}\rangle}\mathbf{1}_{\left\{\frac{1}{n}\sum_{k=1}^{n}x_{k}\in B\right\}}$$
$$\geq Z(\lambda_{\circ})^{n}e^{-n(\langle\lambda_{\circ},y_{\circ}\rangle+\varepsilon)}\int \mu_{\lambda_{\circ}}(\mathrm{d}x_{1})\cdots\int \mu_{\lambda_{\circ}}(\mathrm{d}x_{n})\mathbf{1}_{\left\{\frac{1}{n}\sum_{k=1}^{n}x_{k}\in B_{\delta}(y_{\circ})\right\}}$$
$$= e^{-n(I(y_{\circ})+\varepsilon)}\mathbb{P}\Big[\frac{1}{n}\sum_{k=1}^{n}\hat{X}_{k}\in B_{\delta}(y_{\circ})\Big].$$

By the weak law of large numbers, the probability on the right-hand side here tends to one, so taking logarithms and dividing by n we see that

$$\liminf_{n \to \infty} \frac{1}{n} \log \mathbb{P}\Big[\frac{1}{n} \sum_{k=1}^{n} X_k \in B\Big] \ge -\varepsilon - I(y_\circ) - \varepsilon \ge -2\varepsilon - \inf_{y \in B} I(y).$$

Since  $\varepsilon > 0$  is arbitrary, this completes the proof of the large deviations lower bound in the case that  $B \cap \mathcal{U}_I \neq \emptyset$ . By Lemma 1.9 (c), we also obtain the large deviations lower bound in case  $B \cap \mathcal{U}_I = \emptyset$  but  $\overline{B} \cap \overline{\mathcal{U}}_I \neq \emptyset$ . If  $\overline{B} \cap \overline{\mathcal{U}}_I = \emptyset$ , finally, then  $\inf_{y \in B} I(y) = \infty$  while  $\mathbb{P}[\frac{1}{n} \sum_{k=1}^n X_k \in B] = 0$  for each *n* by Lemma 2.10 (iv), so the bound is trivially fulfilled.

**Remark** In the proof of Theorem 0.1 in Section 2.2, we used the central limit theorem for the titled random variables  $(\hat{X}_k)_{k\geq 1}$  to obtain the large deviations lower bound. In the proof above, we have instead used the weak law of large numbers for the  $(\hat{X}_k)_{k\geq 1}$ . This proof is in a sense more robust, but on the other hand, if one is interested in exact estimates on the error term in formulas such as

$$\mathbb{P}\left[\frac{1}{n}\sum_{k=1}^{n}X_{k}\geq y\right]=e^{-nI(y)+o(n)},$$

then the proof based on the central limit theorem gives the sharpest estimates for o(n).

### 2.6 Sanov's theorem

The aim of this section is to prove the following result, which (at least in the case  $E = \mathbb{R}$ ) goes back to Sanov [San61]. As a simple application, we will also prove Theorem 0.7.

**Theorem 2.18 (Sanov's theorem)** Let  $(X_k)_{k\geq 0}$  be i.i.d. random variables taking values in a Polish space E, with common law  $\mu$ , and let

$$M_n := \frac{1}{n} \sum_{k=1}^n \delta_{X_k} \qquad (n \ge 1)$$

be the empirical laws of the  $(X_k)_{k\geq 0}$ . Then the laws  $\mu_n := \mathbb{P}[M_n \in \cdot]$ , viewed as probability laws on the Polish space  $\mathcal{M}_1(E)$  of probability measures on E, equipped with the topology of weak convergence, satisfy the large deviation principle with speed n and rate function  $H(\cdot|\mu)$ .

**Proof** We apply Theorem 1.34 about projective limits. We first consider the case that E is compact. In this case,  $\mathcal{M}_1(E)$  is also compact so exponential tightness comes for free.

Since  $\mathcal{C}(E)$  is separable, we may choose a countable dense set  $\{f_i : i \in \mathbb{N}_+\} \subset \mathcal{C}(E)$ . For each  $i \in \mathbb{N}_+$ , we define  $\psi_i : \mathcal{M}_1(E) \to \mathbb{R}$  by  $\psi_i(\nu) := \int f_i d\nu$ . The  $(\psi_i)_{i \in \mathbb{N}_+}$  are continuous by the definition of weak convergence of measures. We claim that they also separate points. To see this, imagine that  $\nu, \nu' \in \mathcal{M}_1(E)$  and  $\psi_i(\nu) = \psi_i(\nu')$  for all  $i \geq 1$ . Then  $\int f d\nu = \int f d\nu'$  for all  $f \in \mathcal{C}(E)$  by the fact that  $\{f_i : i \in \mathbb{N}_+\}$  is dense, and therefore  $\nu = \nu'$ .

Let

$$\vec{f}_d(x) := (f_1(x), \dots, f_d(x)) \qquad (x \in E, \ d \ge 1),$$

and

$$\vec{\psi}_d(\nu) := \left(\psi_1(\nu), \dots, \psi_d(\nu)\right) = \int \vec{f}_d \,\mathrm{d}\nu \qquad \left(\nu \in \mathcal{M}_1(E)\right)$$

By Theorem 2.17, for each  $d \ge 1$ , the laws  $\mu_n \circ \vec{\psi}_d^{-1}$  satisfy the large deviation principle with a good rate function  $I_d$ . By Proposition 2.13, this rate function is given by

$$I_d(y) = \inf_{\substack{\nu \in \mathcal{M}_1(E) \\ \int \vec{f}_d \, d\nu = y}} H(\nu \,|\, \mu) \qquad (y \in \mathbb{R}^d).$$

Theorem 1.34 now implies that the measures  $\mu_n$  satisfy the large deviation principle with rate function  $H(\cdot | \mu)$ . This completes the proof for compact E.

To prove the general statement, let  $\overline{E}$  be a metrizable compactification of E. By Proposition 1.26, such a compactification exists and E is a  $G_{\delta}$ -subset of E. By what we have already proved, the laws  $\mu_n$ , viewed as probability laws on the Polish space  $\mathcal{M}_1(\overline{E})$  of probability measures on  $\overline{E}$ , equipped with the topology of weak convergence, satisfy the large deviation principle with speed n and rate function  $H(\cdot | \mu)$ .

We view  $\mathcal{M}_1(E)$  as a subset of  $\mathcal{M}_1(\overline{E})$ . By Exercise 1.28, the topology on  $\mathcal{M}_1(E)$ is the induced topology from  $\mathcal{M}_1(\overline{E})$ . Since  $\mathcal{M}_1(E)$  is Polish in this topology, it must be a  $G_{\delta}$ -subset of  $\mathcal{M}_1(\overline{E})$ . By the restriction principle (Lemma 1.27), using the fact that  $H(\cdot|\mu)$  is a good rate function (which has been proved in Lemma 2.11) and the fact that  $H(\cdot|\mu) = \infty$  on  $\mathcal{M}_1(\overline{E}) \setminus \mathcal{M}_1(E)$ , we conclude that the laws  $\mu_n$ , viewed as probability laws on  $\mathcal{M}_1(E)$ , satisfy the large deviation principle with speed n and rate function  $H(\cdot|\mu)$ .

**Remark** For some purposes, the topology of weak convergence on  $\mathcal{M}_1(E)$  is too weak. With some extra work, it is possible to improve Theorem 2.18 by showing that the emperical measures satisfy the large deviation principle with respect to the (much stronger) topology of strong convergence of measures; see [DS89, Section 3.2].

**Proof of Lemma 0.6 and Theorem 0.7** If in Theorem 2.18, E = S is a finite set and  $\mu(\{x\}) > 0$  for all  $x \in S$ , then the theorem and its proof simplify considerably. In this case, without loss of generality, we may assume that  $S = \{0, \ldots, d\}$  for some  $d \ge 1$ . We may identify  $\mathcal{M}_1(S)$  with the convex subset of  $\mathbb{R}^d$  given by

$$\mathcal{M}_1(S) = \{ x \in \mathbb{R}^d : x(i) \ge 0 \ \forall i = 1, \dots, d, \ \sum_{i=1}^d x(i) \le 1 \},\$$

where x(0) is determined by the condition  $\sum_{i=0}^{d} x(i) = 1$ . Thus, we may apply Cramér's theorem (Theorem 2.17) to the  $\mathbb{R}^{d}$ -valued random variables  $M_{n}$ . The fact that the rate function from Cramér's theorem is in fact  $H(\nu|\mu)$  follows from Lemma 2.12. Since  $\mu(\{x\}) > 0$  for all  $x \in S$ , it is easy to see that the covariance condition of Lemma 2.10 is fulfilled, so Lemma 0.6 follows from Lemma 2.10 and the observation that  $H(\nu|\mu) < \infty$  for all  $\nu \in \mathcal{M}_{1}(S)$ .

**Remark** There exists a nice combinatorical proof of Sanov's theorem for finite spaces (Theorem 0.7), in the spirit of our Section 3.2 below. See [Hol00, Section II.1].

**Exercise 2.19 (Joint continuity of relative entropy)** Let S be a finite set and let  $\mathcal{M}_1(S) := \{ \mu \in \mathcal{M}_1(S) : \mu(x) > 0 \ \forall x \in S \}$ . Prove the continuity of the map

$$\mathcal{M}_1(S) \times \mathcal{M}_1(S) \ni (\nu, \mu) \mapsto H(\nu|\mu).$$

**Exercise 2.20 (Convexity of relative entropy)** Let S be a finite set and let  $\mu \in \mathcal{M}_1(S)$ . Give a direct proof of the fact that

$$\mathcal{M}_1(S) \ni \nu \mapsto H(\nu|\mu)$$

is a lower semi-continuous, convex function.

## Chapter 3

# Markov chains

### 3.1 Basic notions

Let S be a finite set and let P be a probability kernel on S, i.e.,  $P: S \times S \to \mathbb{R}$  is a function such that

(i) 
$$P(x,y) \ge 0$$
  $(x,y \in S),$   
(ii)  $\sum_{y \in S} P(x,y) = 1$   $(x \in S).$ 

For any function  $f: S \to \mathbb{R}$ , we put

$$Pf(x) := \sum_{y \in S} P(x, y) f(y),$$

which defines a linear operator  $P : \mathbb{R}^S \to \mathbb{R}^S$ . For any measure  $\mu$  on S we write  $\mu(x) := \mu(\{x\})$  and for  $f : S \to \mathbb{R}$ , we let

$$\mu f(y) := \sum_{x \in S} \mu(x) f(x)$$

denote the expectation of f w.r.t.  $\mu$ . Viewing a measure  $\mu$  as a linear operator  $\mu : \mathbb{R}^S \to \mathbb{R}$ , we see that the composition of a probability kernel  $P : \mathbb{R}^S \to \mathbb{R}^S$  and a probability measure  $\mu : \mathbb{R}^S \to \mathbb{R}$  is an operator  $\mu P : \mathbb{R}^S \to \mathbb{R}$  that corresponds to the probability measure  $\mu P(y) = \sum_{x \in S} \mu(x) P(x, y)$ .

A Markov chain with state space S, transition kernel P and initial law  $\mu$  is a collection of S-valued random variables  $(X_k)_{k>0}$  whose finite-dimensional distributions

are characterized by

$$\mathbb{P}[X_0 = x_0, \dots, X_n = x_n] = \mu(x_0) P(x_0, x_1) \cdots P(x_{n-1}, x_n)$$

 $(n \ge 1, x_0, \ldots, x_n \in S)$ . Note that in particular, the law of  $X_n$  is given by  $\mu P^n$ , where  $P^n$  is the *n*-th power of the linear operator P. We also introduce the notation

$$\mu \otimes P(x_0, x_1) := \mu(x_0) \otimes P(x_0, x_1)$$

to denote the probability measure on  $S^2$  that is the law of  $(X_0, X_1)$ .

Write  $x \stackrel{P}{\longrightarrow} y$  if there exist  $n \ge 0$  such that  $P^n(x, y) > 0$  or equivalently, there exist  $x = x_0, \ldots, x_n = y$  such that  $P(x_{k-1}, x_k) > 0$  for each  $k = 1, \ldots, n$ . Then P is called *irreducible* if  $x \stackrel{P}{\longrightarrow} y$  for all  $x, y \in S$ . An *invariant law* of P is a probability measure  $\mu$  on S such that  $\mu P = \mu$ . Equivalently,  $\mu$  is invariant if the Markov chain  $(X_k)_{k\ge 0}$  with transition kernel P and initial law  $\mu$  is *stationary*, i.e.  $(X_k)_{k\ge 0}$  is equal in law to  $(Y_k)_{k\ge 0}$  defined as  $Y_k := X_{k+1}$  ( $k \ge 0$ ). The *period* of a state  $x \in S$  is the greatest common divisor of the set  $\{n \ge 1 : P^n(x, x) > 0\}$ . If P is irreducible, then all states have the same period. If all states have period one, then we say that P is *aperiodic*. Basic results of Markov chain theory tell us that an irreducible Markov chain with a finite state space S has a unique invariant law  $\mu$ , which has the property that  $\mu(x) > 0$  for all  $x \in S$ . If P is moreover aperiodic, then  $\nu P^n$  converges to  $\mu$  as  $n \to \infty$ , for each initial law  $\nu$ .

For any Markov chain  $X = (X_k)_{k>0}$ , we let

$$M_n^{(2)} := \frac{1}{n} N_n^{(2)}, \quad \text{where} \quad N_n^{(2)}(x) := \sum_{k=1}^n \mathbb{1}_{\{(X_{k-1}, X_k) = (x_1, x_2)\}}$$
(3.1)

 $(x \in S^2, n \ge 1)$  be the *pair empirical distribution* of the first n + 1 random variables. The  $M_n^{(2)}$  are random variables taking values in the space  $\mathcal{M}_1(S^2)$  of probability measures on  $S^2 := \{x = (x_1, x_2) : x_i \in S \; \forall i = 1, 2\}$ . If X is irreducible, then the  $M_n^{(2)}$  satisfy a strong law of large numbers.

**Proposition 3.1 (SLLN for Markov chains)** Let  $X = (X_k)_{k\geq 0}$  be an irreducible Markov chain with finite state space S, transition kernel P, and arbitrary initial law. Let  $(M_n^{(2)})_{n\geq 1}$  be the pair empirical distributions of X and let  $\mu$  be its invariant law. Then

$$M_n^{(2)} \xrightarrow[n \to \infty]{} \mu \otimes P \quad \text{a.s.}$$
 (3.2)

### 3.2. A LDP FOR MARKOV CHAINS

**Proof (sketch)** It suffices to prove the statement for deterministic starting points  $X_0 = z$ . Let  $\tau_0 := 0$  and  $\tau_N := \inf\{k > \tau_{N-1} : X_k = z\}$   $(N \ge 1)$  be the return times of X to z and define random variables  $(Y_N)_{N>1}$  by

$$Y_N(x) := \sum_{k=\tau_{N-1}+1}^{\tau_N} \mathbb{1}\{(X_{k-1}, X_k) = (x_1, x_2)\} \qquad (x \in S^2)$$

It is not hard to check that the  $(Y_N)_{N\geq 1}$  are i.i.d. with finite mean  $\mathbb{E}[Y_i(x_1, x_2)] = \mathbb{E}[\tau_1] \ \nu \otimes P(x_1, x_2) \ ((x_1, x_2) \in S^2)$ , and the  $(\tau_N - \tau_{N-1})_{N\geq 1}$  are i.i.d. with mean  $\mathbb{E}[\tau_1]$ . Therefore, by the ordinary strong law of large numbers

$$M_{\tau_N}^{(2)} = \frac{N}{\tau_N} \frac{1}{N} \sum_{M=1}^N Y_M \xrightarrow[N \to \infty]{} \nu \otimes P \quad \text{a.s.}$$

The final part of the proof is a bit technical. For each  $n \ge 0$ , let  $N(n) := \inf\{N \ge 1 : \tau_N \ge n\}$ . Using Borel-Cantelli, one can check that for each  $\varepsilon > 0$ , the event

$$\{|M_n^{(2)} - M_{\tau_{N(n)}}^{(2)}| \ge \varepsilon\}$$

occurs only for finitely many n. Using this and the a.s. convergence of the  $M_{\tau_{N(n)}}^{(2)}$  one obtains the a.s. convergence of the  $M_n^{(2)}$ .

We will be interested in large deviations away from (3.2).

### **3.2** A LDP for Markov chains

In this section, we prove a basic large deviation result for the empirical pair distribution of irreducible Markov chains. For concreteness, for any finite set S, we equip the space  $\mathcal{M}_1(S)$  of probability measures on S with the *total variation distance* 

$$d(\mu,\nu) := \sup_{A \subset S} |\mu(A) - \nu(A)| = \frac{1}{2} \sum_{x \in S} |\mu(x) - \nu(x)|,$$

where for simplicity we write  $\mu(x) := \mu(\{x\})$ . Note that since S is finite, convergence in total variation norm is equivalent to weak convergence or pointwise convergence (and in fact any reasonable form of convergence one can think of).

For any  $\nu \in \mathcal{M}_1(S^2)$ , we let

$$\nu^1(x_1) := \sum_{x_2 \in S} \nu(x_1, x_2) \quad \text{and} \quad \nu^2(x_2) := \sum_{x_1 \in S} \nu(x_1, x_2)$$

denote the first and second marginals of  $\nu$ , respectively, and we let

$$\mathcal{V} := \left\{ \nu \in \mathcal{M}_1(S^2) : \nu^1 = \nu^2 \right\}$$

denote the space of all probability measures on  $S^2$  whose first and second marginals agree. The main result of this section is the following theorem.

**Theorem 3.2 (LDP for Markov chains)** Let  $X = (X_k)_{k\geq 0}$  be a Markov chain with finite state space S, irreducible transition kernel P, and arbitrary initial law. Let  $(M_n^{(2)})_{n\geq 1}$  be the pair empirical distributions of X. Then the laws  $\mathbb{P}[M_n^{(2)} \in \cdot]$ satisfy the large deviation principle with speed n and rate function  $I^{(2)}$  given by

$$I^{(2)}(\nu) := \begin{cases} H(\nu|\nu^1 \otimes P) & \text{if } \nu \in \mathcal{V}, \\ \infty & \text{otherwise}, \end{cases}$$

where  $H(\cdot | \cdot)$  denotes the relative entropy of one measure w.r.t. another.

**Remark** By the contraction principle, Theorem 3.2 also gives us a large deviation principle for the 'usual' empirical distributions

$$M_n(x) := \frac{1}{n} \sum_{k=0}^{n-1} \mathbb{1}_{\{X_k = x\}} \qquad (x \in S, \ n \ge 1).$$

In this case, however, it is in general<sup>1</sup> not possible to write down a nice, explicit formula for the rate function. This is because pairs are the 'natural' object to look at for Markov processes.

The proof of Theorem 3.2 needs some preparations.

### Lemma 3.3 (Characterization as invariant measures) One has

$$\mathcal{V} = \left\{ \nu^1 \otimes P : \nu^1 \in \mathcal{M}_1(S), \ P \ a \ probability \ kernel \ on \ S, \ \nu^1 P = \nu^1 \right\}.$$

**Proof** If P is a probability kernel on S, and  $\nu^1 \in \mathcal{M}_1(S)$  satisfies  $\nu^1 P = \nu^1$  (i.e.,  $\nu^1$  is an invariant law for the Markov chain with kernel P), then  $(\nu^1 \otimes P)^2 = \nu^1 P = \nu^1$ , which shows that  $\nu^1 \otimes P \in \mathcal{V}$ . On the other hand, for any  $\nu \in \mathcal{V}$ , we may define a kernel P by setting

$$P(x_1, x_2) := \frac{\nu(x_1, x_2)}{\nu^1(x_1)},$$

<sup>&</sup>lt;sup>1</sup>An exception are continuous-time reversible Markov chains. See [Hol00, Thm. IV.14(b)].

whenever the denominator is nonzero, and choosing  $P(x_1, \cdot)$  in some arbitrary way if  $\nu^1(x_1) = 0$ . Then  $\nu^1 \otimes P(x_1, x_2) = \nu(x_1, x_2)$  and  $\nu^1 P = (\nu^1 \otimes P)^2 = \nu^2 = \nu^1$ by the fact that  $\nu \in \mathcal{V}$ .

For any  $z \in S$ , let us define

$$\mathcal{R}_{n,z} := \left\{ r \in \mathbb{N}^{S^2} : \exists (x_0, \dots, x_n) \in S^{n+1}, \ x_0 = z, \\ \text{s.t.} \ r(y_1, y_2) = \sum_{k=1}^n \mathbb{1}_{\left\{ (x_{k-1}, x_k) = (y_1, y_2) \right\}} \ \forall y \in S^2 \right\}$$

and  $\mathcal{R}_n := \bigcup_{z \in S} \mathcal{R}_{n,z}$ . Then the random variables  $N_n^{(2)}$  from (3.1) take values in  $\mathcal{R}_n$ . For the pair empirical distributions  $M_n^{(2)}$ , the relevant spaces are

$$\mathcal{V}_n := \{ n^{-1}r : r \in \mathcal{R}_n \}$$
 and  $\mathcal{V}_{n,z} := \{ n^{-1}r : r \in \mathcal{R}_{n,z} \}.$ 

For any  $U \subset S^2$ , we identify the space  $\mathcal{M}_1(U)$  of probability laws on U with the space

$$\left\{\nu \in \mathcal{M}_1(S^2) : \nu(x_1, x_2) = 0 \ \forall x \notin U\right\},\$$

and we define

$$\mathcal{V}(U) := \mathcal{V} \cap \mathcal{M}_1(U), \quad \mathcal{V}_n(U) := \mathcal{V}_n \cap \mathcal{M}_1(U), \text{ and } \mathcal{V}_{n,z}(U) := \mathcal{V}_{n,z} \cap \mathcal{M}_1(U).$$

We will need a lemma that says that for suitable  $U \subset S^2$ , the spaces  $\mathcal{V}_n(U)$  approximate  $\mathcal{V}(U)$  as  $n \to \infty$ . The typical example we have in mind is  $U = \{(x_1, x_2) \in S^2 : P(x_1, x_2) > 0\}$  where P is an irreducible probability kernel on S or some subset of S. For any  $U \subset S^2$ , let us write

$$\overline{U} := \{ x_1 \in S : (x_1, x_2) \in U \text{ for some } x_2 \in S \} \\ \cup \{ x_2 \in S : (x_1, x_2) \in U \text{ for some } x_1 \in S \}.$$
(3.3)

We will say that U is *irreducible* if for every  $x, y \in \overline{U}$  there exist  $n \ge 0$  and  $x = x_0, \ldots, x_n = y$  such that  $(x_{k-1}, x_k) \in U$  for all  $k = 1, \ldots, n$ .

### Lemma 3.4 (Limiting space of pair empirical distribution) One has

$$\lim_{n \to \infty} \sup_{\nu \in \mathcal{V}_n} d(\nu, \mathcal{V}) = 0.$$
(3.4)

Moreover, for each  $z \in S$  and  $\nu \in \mathcal{V}$  there exist  $\nu_n \in \mathcal{V}_{n,z}$  such that  $d(\nu_n, \nu) \to 0$ as  $n \to \infty$ . If  $U \subset S^2$  is irreducible, then moreover, for each  $z \in \overline{U}$  and  $\nu \in \mathcal{V}(U)$ there exist  $\nu_n \in \mathcal{V}_{n,z}(U)$  such that  $d(\nu_n, \nu) \to 0$  as  $n \to \infty$ . **Proof** We leave formula (3.4) as an excercise to the reader (Excercise 3.5 below). To prove that for any  $z \in S$  we can approximate arbitrary  $\nu \in \mathcal{V}$  with  $\nu_n \in \mathcal{V}_{n,z}$ , by a simple diagonal argument (Exercise 3.6 below), we can without loss of generality assume that  $\nu(x) > 0$  for all  $x \in S^2$ . By Lemma 3.3, there must exist some probability kernel P on S such that  $\nu = \nu^1 \otimes P$  and  $\nu^1 P = \nu^1$ . Since  $\nu(x) > 0$  for all  $x \in S^2$ , we must have  $P(x_1, x_2) > 0$  for all  $x \in S^2$ . In particular, this implies that P is irreducible and  $\nu^1$  is the unique invariant law of P. Let  $X = (X_k)_{k\geq 0}$  be a Markov chain with transition kernel P and initial state  $X_0 = z$ , and let  $(M_n^{(2)})_{n\geq 1}$  be its pair empirical measures. Then  $M_n^{(2)} \in \mathcal{V}_{n,z}$  for all  $n \geq 1$  while  $M_n^{(2)} \to \nu^1 \otimes P = \nu$  a.s. by Proposition 3.1. Since the empty set cannot have probability one, it follows that there must exist  $\nu_n \in \mathcal{V}_{n,z}$  such that  $d(\nu_n, \nu) \to 0$  as  $n \to \infty$ .

The same argument shows that if U is irreducible, then for any  $z \in \overline{U}$ , an arbitrary  $\nu \in \mathcal{V}(U)$  can be approximated with  $\nu_n \in \mathcal{V}_{n,z}(U)$ . In this case, by a diagonal argument, we may assume without loss of generality that  $\nu(x) > 0$  for all  $x \in U$ . By Lemma 3.3, there exists some probability kernel P on  $\overline{U}$  such that  $\nu = \nu^1 \otimes P$  and  $\nu^1 P = \nu^1$ . Since  $\nu(x) > 0$  for all  $x \in U$ , we must have  $P(x_1, x_2) > 0$  for all  $x \in U$ , hence P is irreducible. Using the strong law of large numbers for the Markov chain with transition kernel P, the argument then proceeds as before.

Exercise 3.5 (Marginals almost agree) Prove formula (3.4).

**Exercise 3.6 (Diagonal argument)** Let (E, d) be a metric space, let  $x_n, x \in E$  satisfy  $x_n \to x$  and for each n, let  $x_{n,m} \in E$  satisfy  $x_{n,m} \to x_n$  as  $m \to \infty$ . Then there exist  $m(n) \to \infty$  such that  $x_{n,m'(n)} \to x$  for all  $m'(n) \ge m(n)$ .

**Exercise 3.7 (Continuity of rate function)** Let P be a probability kernel on S and let  $U := \{(y_1, y_2) \in S^2 : P(y_1, y_2) > 0\}$ . Prove the continuity of the map

$$\mathcal{M}_1(U) \ni \nu \mapsto H(\nu | \nu^1 \otimes P).$$

Show that if  $U \neq S^2$ , then the map  $\mathcal{M}_1(S^2) \ni \nu \mapsto H(\nu | \nu^1 \otimes P)$  is not continuous.

**Proof of Theorem 3.2** If  $\mu_n, \mu'_n$  both satisfy a large deviation principle with the same speed and rate function, then any convex combination of  $\mu_n, \mu'_n$  also satisfies this large deviation principle. In view of this, it suffices to prove the claim for Markov chains started in a deterministic initial state  $X_0 = z$ .

### 3.2. A LDP FOR MARKOV CHAINS

We observe that for any  $r: S^2 \to \mathbb{N}$ , the pair counting process defined in (3.1) satisfies

$$\mathbb{P}[N_n^{(2)} = r] = \mathcal{C}_{n,z}(r) \prod_{(x_1, x_2) \in S^2} P(x_1, x_2)^{r(x_1, x_2)}, \qquad (3.5)$$

where

$$\mathcal{C}_{n,z}(r) := \left| \{ x \in S^{n+1} : x_0 = z, \ \sum_{k=1}^n \mathbb{1}_{\{ (x_{k-1}, x_k) = (y_1, y_2) \}} = r(y_1, y_2) \ \forall y \in S^2 \} \right|$$

is the number of different sequences  $X_0, \ldots, X_n$  that give rise to the same pair frequencies  $N_n^{(2)} = r$ . In order to estimate  $\mathcal{C}_{n,z}(r)$ , for a given  $r \in \mathcal{R}_{n,z}$ , we draw a directed graph whose vertex set is S and that has  $r(x_1, x_2)$  arrows pointing from  $x_1$  to  $x_2$ . Let  $\mathcal{W}_{n,z}(r)$  be the number of distinct walks in this graph that start at z and that use each arrow exactly once, where we distinguish between different arrows, i.e., if there are more arrows pointing from  $x_1$  to  $x_2$ , then we do care about which arrow is used first, which arrow next, and so on. Then

$$C_{n,z}(r) = \frac{\mathcal{W}_{n,z}(r)}{\prod_{(x_1, x_2) \in S^2} r(x_1, x_2)!}.$$
(3.6)

A simple combinatorical argument (see Lemma 3.8 below) shows that

$$\prod_{x_1:r^1(x_1)>0} (r^1(x_1)-1)! \le \mathcal{W}_{n,z}(r) \le \prod_{x_1\in S} r^1(x_1)! \qquad (r\in\mathcal{R}_n).$$
(3.7)

Combining (3.6), (3.7) and (3.5), we obtain the bounds

$$\frac{\prod_{x_1:r(x_1)>0} (r^1(x_1) - 1)!}{\prod_{(x_1,x_2)\in S^2} r(x_1,x_2)!} \prod_{\substack{(x_1,x_2)\in S^2 \\ x_1,x_2)\in S^2}} P(x_1,x_2) r(x_1,x_2) 
\leq \mathbb{P}[N_n^{(2)} = r] \leq \frac{\prod_{x_1\in S} r^1(x_1)!}{\prod_{(x_1,x_2)\in S^2} r(x_1,x_2)!} \prod_{(x_1,x_2)\in S^2} P(x_1,x_2) r(x_1,x_2)$$
(3.8)

 $(r \in \mathcal{R}_{n,z})$ . We recall that *Stirling's formula*<sup>2</sup> implies that

$$\log(n!) = n \log n - n + H(n)$$
 as  $n \to \infty$ ,

where we use the convention that  $0 \log 0 = 0$ , and the error term H(n) is of order  $\log n$  and can in fact uniformly be estimated as

$$|H(n)| \le C \log n \qquad (n \ge 0),$$

<sup>&</sup>lt;sup>2</sup>Recall that Stirling's formula says that  $n! \sim \sqrt{2\pi n} (n/e)^n$ .

with  $C < \infty$  some constant. It follows that the logarithm of the right-hand side of (3.8) is given by

$$\begin{split} &\sum_{x_1 \in S} \left( r^1(x_1) \log r^1(x_1) - r^1(x_1) + H(r^1(x_1)) \right) \\ &- \sum_{(x_1, x_2) \in S^2} \left( r(x_1, x_2) \log r(x_1, x_2) - r(x_1, x_2) + H(r(x_1, x_2)) \right) \\ &+ \sum_{(x_1, x_2) \in S^2} r(x_1, x_2) \log P(x_1, x_2) \\ &= \sum_{(x_1, x_2) \in S^2} r(x_1, x_2) \left( \log r^1(x_1) + \log P(x_1, x_2) - \log r(x_1, x_2) \right) + H'(r, n), \end{split}$$

where we have used that  $\sum_{x_1} r^1(x_1) = n = \sum_{(x_1,x_2) \in S^2} r(x_1,x_2)$  and H'(r,n) is an error term that can be estimated uniformly in r as

$$|H'(r,n)| \leq \sum_{x_1 \in S} C \log(r^1(x_1)) + \sum_{\substack{(x_1,x_2) \in S^2 \\ \leq C(|S| + |S|^2) \log n}} C \log r(x_1,x_2)$$

with the same constant C as before. Dividing by n, we find that

$$\begin{aligned} \frac{1}{n} \log \mathbb{P}[N_n^{(2)} &= r] \leq -\sum_{(x_1, x_2) \in S^2} \frac{r(x_1, x_2)}{n} \log \frac{r(x_1, x_2)}{r^1(x_1) P(x_1, x_2)} + \frac{1}{n} H'(r, n) \\ &= -H(\nu | \nu_r^1 \otimes P) + \frac{1}{n} H'(r, n), \end{aligned}$$

where  $\nu(x_1, x_2) := n^{-1}r(x_1, x_2)$ . Treating the left-hand side of (3.8) in much the same way, we find that

$$\frac{1}{n}\log\mathbb{P}[M_n^{(2)} = \nu] = -H(\nu|\nu^1 \otimes P) + O(n^{-1}\log n)$$
(3.9)

for all  $\nu \in \mathcal{V}_{n,z}$ , where the error term is of order  $n^{-1} \log n$  uniformly for all  $\nu \in \mathcal{V}_{n,z}$ . We are now almost done. Let  $U := \{(x_1, x_2) \in S^2 : P(x_1, x_2) > 0\}$ . Then obviously  $M_n^{(2)} \in \mathcal{M}_1(U)$  for all  $n \ge 1$ , hence by the restriction principle (Lemma 1.27) and the fact that  $H(\nu|\nu^1 \otimes P) = \infty$  for all  $\nu \notin \mathcal{M}_1(U)$ , instead of proving the large deviation principle on  $\mathcal{M}_1(S^2)$ , we may equivalently prove the large deviation principle on  $\mathcal{M}_1(U)$ . By Excercise 3.7, the map

$$\mathcal{M}_1(U) \ni \nu \mapsto H(\nu|\nu^1 \otimes P)$$

### 3.2. A LDP FOR MARKOV CHAINS

is continuous. (Note that we need the space  $\mathcal{M}_1(U)$  since the same is not true for  $\mathcal{M}_1(S^2) \ni \nu \mapsto H(\nu|\nu^1 \otimes P)$ .) Using the continuity of this map and Lemmas 1.15 and 1.18, we see that it suffices to show that the counting measures on  $\mathcal{V}_{n,z}(U)$ 

$$\rho_n := \sum_{\nu \in \mathcal{V}_{n,z}(U)} \delta_{\nu}$$

satisfy the large deviation principle on  $\mathcal{M}_1(U)$  with speed n and trivial rate function

$$J(\nu) := \begin{cases} 0 & \text{if } \nu \in \mathcal{V}(U), \\ \infty & \text{otherwise.} \end{cases}$$

We will prove the large deviations upper and lower bounds from Proposition 1.7. For the upper bound, we observe that if  $C \subset \mathcal{M}_1(U)$  is closed and  $C \cap \mathcal{V}(U) = \emptyset$ , then, since  $\mathcal{V}(U)$  is a compact subset of  $\mathcal{M}_1(U)$ , the distance  $d(C, \mathcal{V}(U))$  must be strictly positive. By Lemma 3.4, it follows that  $C \cap \mathcal{V}_n(U) = \emptyset$  for *n* sufficiently large and hence  $\limsup_{n\to\infty} \frac{1}{n} \log \mathbb{P}[M_n^{(2)} \in C] = \infty$ . If  $C \cap \mathcal{V}(U) \neq \emptyset$ , then we may use the fact that  $|\mathcal{V}_n| \leq n^{|S|^2}$ , to obtain the crude estimate

$$\limsup_{n \to \infty} \frac{1}{n} \log \rho_n(C) \le \limsup_{n \to \infty} \frac{1}{n} \log \rho_n(\mathcal{V}) \le \lim_{n \to \infty} \frac{1}{n} \log \left( n^{|S|^2} \right) = 0,$$

which completes our proof of the large deviations upper bound. To prove also the large deviations lower bound, let  $O \subset \mathcal{M}_1(U)$  be open and let  $O \cap \mathcal{V}(U) \neq \emptyset$ (otherwise the statement is trivial). Pick any  $\nu \in O \cap \mathcal{V}(U)$ . By Lemma 3.4, we can choose  $\nu_n \in \mathcal{V}_{n,z}(U)$  such that  $\nu_n \to \nu$ . It follows that  $\nu_n \in O$  for n sufficiently large, and hence

$$\liminf_{n \to \infty} \frac{1}{n} \log \rho_n(O) \ge \lim_{n \to \infty} \frac{1}{n} \log \rho_n(\{\nu_n\}) = 0,$$

as required.

We still need to prove the estimates (3.7). Let G = (V, E) be a finite directed graph with vertex set V and set of directed edges E. For each edge  $e \in E$  there is defined a starting vertex  $e^- \in V$  and endvertex  $e^+ \in V$ . We allow for the case that  $e^- = e^+$  (in this case, e is called a loop). We write

$$E_{x,\bullet} := \{ e \in E : e^- = x \}, \quad E_{\bullet,y} := \{ e \in E : e^+ = y \}, \text{ and } E_{x,y} := E_{x,\bullet} \cap E_{\bullet,y}$$

for the sets of all edges with a specified starting vertex, or endvertex, or both. We allow for the case that  $r(x, y) := |E_{x,y}|$  is larger than one.

By definition, a walk is an ordered collection of edges  $(e_1, \ldots, e_n)$  such that  $e_k^+ = e_{k+1}^-$  for  $k = 1, \ldots, n-1$ . We call  $e_1^-$  and  $e_n^+$  the starting vertex and endvertex of the walk. For any subset of edges  $F \subset E$ , we write  $x \rightsquigarrow_F y$  if x = y or there exists a walk using only edges from F with starting vertex x and endvertex y. By definition, a (directed) spanning tree rooted at  $z \in V$  is a collection of edges  $T \subset E$  such that  $|T \cap E_{x,\bullet}| = 1$  and  $x \rightsquigarrow_T z$  for all  $x \in V$ , i.e., from each vertex there is a unique directed path to the root.

**Lemma 3.8 (Walks that use all edges)** Let G = (V, E) be a finite directed graph and let  $y, z \in V$ . Write  $r(x_1, x_2) := |E_{x_1, x_2}|, r^1(x_1) := |E_{x_1, \bullet}|, and r^2(x_2) := |E_{\bullet, x_2}|$   $(x_1, x_2 \in S)$ . Assume that  $r^1(x) > 0$  for each  $x \in V$  and that

$$r^{1}(x) - r^{2}(x) = 1_{\{x=y\}} - 1_{\{x=z\}}$$
  $(x \in V).$  (3.10)

Let  $\mathcal{W}$  denote the number of walks in G that end in z and use each edge exactly once. Let  $\mathcal{T}$  denote the number of spanning trees rooted at z. Then

$$\mathcal{W} = \mathcal{T}r^1(z)\prod_{x\in V} \left(r^1(x) - 1\right)! \tag{3.11}$$

In particular, one has the estimates (3.7).

**Proof** Let W denote the set of all walks w in G that end in z and use each edge exactly once. It follows from (3.10) that each  $w \in W$  must start in y. We can encode such a walk by numbering, for each  $x \in V$ , the set of outgoing edges  $E_{x,\bullet}$ at x according to which edges is used first, second etc. Let  $\Pi$  be the collection of all functions  $\pi : E \to \mathbb{N}_+$  such that  $\pi : E_{x,\bullet} \to \{1,\ldots,r^1(x)\}$  is a bijection for each  $x \in V$ . We say that such a function  $\pi$  encodes a walk  $w \in W$  if for each  $x \in V$ and  $e \in E_{x,\bullet}$ , one has  $\pi(e) = k$  iff w leaves x for the k-th time using the edge e. Clearly,  $\mathcal{W} = |W| \leq |\Pi|$  which yields the upper bound in (3.7). For any  $\pi \in \Pi$ , let  $T_{\pi} := \bigcup_{x \in V \setminus \{z\}} \{e \in E_x : \pi(e) = r^1(x)\}$ . In particular, if  $\pi$  encodes a walk  $w \in W$ , then these are the arrows used when the walk leaves a vertex  $\neq z$  for the last time. We claim that:

• A function  $\pi \in \Pi$  encodes a walk  $w \in W$  if and only if  $T_{\pi}$  is a spanning tree rooted at z.

Indeed, given a walk  $w \in W$ , if for a vertex  $\neq z$ , we follow the arrow used when w last leaves this vertex, and so on for the next vertex, then we end up in z, proving that  $T_{\pi}$  is a spanning tree rooted at z. Conversely, if a function  $\pi \in \Pi$  has the

property that  $T_{\pi}$  is a spanning tree rooted at z, then, starting at y, we can walk around through the graph in such a way that if at a given moment we are in a vertex x, then we leave x using the outgoing edge  $e \in E_{x,\bullet}$  with the lowest number  $\pi(e)$  that has not yet been used. This process stops when we arrive at a vertex such that all outgoing edges at this vertex have been used. By (3.10), it follows that all incoming arrows have also been used, which is possible only if we are in z. We observe that if  $e \in T_{\pi}$  has been used, then all arrows in  $E_{e^-,\bullet}$  have been used and hence by (3.10) also all arrows in  $E_{\bullet,e^-}$  have been used. Since all arrows in  $E_{\bullet,z}$  have been used and  $T_{\pi}$  is a spanning tree rooted at z, it follows that all arrows in  $T_{\pi}$  have been used, which implies that all arrows in E have been used, i.e.,  $w \in W$ .

This completes the proof of (3.11). In particular, fixing one spanning tree rooted at z, in each vertex  $x \neq z$  we have  $(r^1(x) - 1)!$  ways to choose the order of the outgoing edges except for the one that is used last, which yields the lower bound in (3.7). (Note that in (3.7), we apply Lemma 3.8 to the subgraph consisting of all vertices of G that have been visited at least once.)

The proof of Theorem 3.2 yields a useful corollary. Below, we use the notation

$$H(\nu|\mu) := \sum_{x \in S} \nu(x) \log \frac{\nu(x)}{\mu(x)} = \sum_{x \in S} \mu(x) \frac{\nu(x)}{\mu(x)} \log \frac{\nu(x)}{\mu(x)},$$

even if  $\mu$  is not a probability measure. Note that below, the transition kernel P need not be irreducible!

**Corollary 3.9 (Restricted Markov process)** Let  $X = (X_k)_{k\geq 0}$  be a Markov chain with finite state space S, transition kernel P, and arbitrary initial law. Let

$$U \subset \{(x_1, x_2) \in S^2 : P(x_1, x_2) > 0\}$$

be irreducible and let  $X_0 \in \overline{U}$  a.s. Let  $(M_n^{(2)})_{n\geq 1}$  be the pair empirical distributions of X and let  $\tilde{P}$  denote the restriction of P to U. Then the restricted measures

$$\mathbb{P}[M_n^{(2)} \in \cdot ]\big|_{\mathcal{M}_1(U)}$$

satisfy the large deviation principle with speed n and rate function  $I^{(2)}$  given by

$$\tilde{I}^{(2)}(\nu) := \begin{cases} H(\nu|\nu^1 \otimes \tilde{P}) & \text{if } \nu \in \mathcal{V}(U), \\ \infty & \text{otherwise.} \end{cases}$$

**Proof** The restricted measures  $\mathbb{P}[M_n^{(2)} \in \cdot]|_{\mathcal{M}_1(U)}$  are no longer probability measures, but we have never used this in the proof of Theorem 3.2. In fact, a careful inspection reveals that the proof carries over without a change, where we only need the irreducibility of U (but not of P). In particular, formula (3.9) also holds for the restricted measures and the arguments below there work for any irreducible  $U \subset \{(x_1, x_2) \in S^2 : P(x_1, x_2) > 0\}.$ 

**Exercise 3.10 (Relative entropy and conditional laws)** Let S be a finite space, let  $\nu, \mu$  be probability measures on S and let Q, P be probability kernels on S. Show that

$$H(\nu \otimes Q|\mu \otimes P) = H(\nu|\mu) + \sum_{x_1 \in S} \nu(x_1) H(Q_{x_1}|P_{x_1}),$$

where  $Q_{x_1}(x_2) := Q(x_1, x_2)$  and  $P_{x_1}(x_2) := P(x_1, x_2)$   $((x_1, x_2) \in S^2)$ . In particular, if Q is a probability kernel such that  $\nu = \nu^1 \otimes Q$ , then

$$H(\nu|\nu^1 \otimes P) = \sum_{x_1 \in S} \nu^1(x_1) H(Q_{x_1}|P_{x_1}).$$

**Exercise 3.11 (Minimizer of the rate function)** Let P be irreducible. Show that the unique minimizer of the function  $\mathcal{V} \ni \nu \mapsto H(\nu|\nu^1 \otimes P)$  is given by  $\nu = \mu \otimes P$ , where  $\mu$  is the invariant law of P.

By definition, a *cycle* in S is an ordered collection  $C = (x_1, \ldots, x_n)$  of points in S such that  $x_1, \ldots, x_n$  are all different. We call two cycles equal if they differ only by a cyclic permutation of their points and we call  $|C| = n \ge 1$  the *length* of a cycle  $C = (x_1, \ldots, x_n)$ . We write  $(y_1, y_2) \in C$  if  $(y_1, y_2) = (x_{k-1}, x_k)$  for some  $k = 1, \ldots, n$ , where  $x_0 := x_n$ .

Recall that an element x of a convex set K is an *extremal element* if x cannot be written as a nontrivial convex combination of other elements of K, i.e., there do not exist  $y, z \in K, y \neq z$  and 0 such that <math>x = py + (1 - p)z. If  $K \subset \mathbb{R}^d$  is convex and compact, then it is known that for each element  $x \in K$  there exists a unique probability measure  $\rho$  on the set  $K_e$  of extremal elements of K such that  $x = \int y\rho(dy)$ .

**Exercise 3.12 (Cycle decomposition)** Prove that the extremal elements of the space  $\mathcal{V}$  are the probability measures of the form

$$\nu_C(y_1, y_2) := \frac{1}{|C|} \mathbf{1}_{\{(y_1, y_2) \in C\}},$$

### 3.2. A LDP FOR MARKOV CHAINS

where  $C = (x_1, \ldots, x_n)$  is a cycle in S. Hint: show that for each  $\nu \in \mathcal{V}$  and  $(y_1, y_2) \in S^2$  such that  $\nu(y_1, y_2) > 0$ , one can find a cycle  $C \in \mathcal{C}(S^2)$  and a constant c > 0 such that  $(y_1, y_2) \in C$  and  $c\nu_C \leq \nu$ . Use this to show that for each  $\nu \in \mathcal{V}_e$  there exists a cycle C such that  $\nu(y_1, y_2) = 0$  for all  $(y_1, y_2) \notin C$ .

Note Since  $\mathcal{V}$  is a finite dimensional, compact, convex set, Excercise 3.12 shows that for each  $\nu \in \mathcal{V}$ , there exists a unique probability law  $\rho$  on the set of all cycles in S such that

$$\nu(y_1, y_2) = \sum_C \rho(C) \nu_C(y_1, y_2),$$

where the sum rums over al cycles in S. Note that in Excercise 3.12, you are not asked to give an explicit formula for  $\rho$ .

**Exercise 3.13 (Convexity of rate function (!))** Let P be a probability kernel on S. Prove that

$$\mathcal{M}_1(S^2) \ni \nu \mapsto H(\nu|\nu^1 \otimes P)$$

is a convex, lower semi-continuous function.

**Important note** I do not know an elegant solution to this exercise. I originally copied this from [Hol00], who first gives the special case that  $P(x, y) = \mu(y)$  does not depend on x as Exercise II.12, and then in his Lemma IV.5 shows that the general case can easily be derived from this special case. Den Hollander probably based himself on Problems IX 6.1 and 6.2 from [Ell85]. These problems, however, are meant to be solved using deep theory from Chapter IX of [Ell85] that is not available here or in [Hol00].

Recall that if  $(X_k)_{k\geq 0}$  is a Markov chain with initial law  $\mu$  and transition kernel P, then  $\mu \otimes P$  is the joint law of  $(X_0, X_1)$ . More generally, let  $\mu \otimes^n P$  denote the joint law of  $(X_0, \ldots, X_{n-1})$ . Let  $\mu, \nu$  be invariant laws of probability kernels P, Q, respectively. I conjecture that if P is irreducible, then

$$H(\nu \otimes Q|\nu^1 \otimes P) = \lim_{n \to \infty} \frac{1}{n} H(\nu \otimes^n Q|\mu \otimes^n P).$$

Results in this spirit are proved in Chapter IX of [Ell85]. In particular, it is shown there that if one wishes to minimize the relative entropy density of a stationary measure with respect to a product measure, under the condition that the two-dimensional marginals are given by some  $\nu \in \mathcal{M}_1(S^2)$ , then the minimum is attained by the Markov chain that has these two-dimensional marginals. The convexity of  $\nu \mapsto H(\nu | \nu^1 \otimes P)$  then follows from this and the contraction principle. **Exercise 3.14 (Not strictly convex)** Let P be any probability kernel on  $S = \{1, 2\}$ . Define  $\mu, \nu \in \mathcal{M}_1(S^2)$  by

$$\begin{pmatrix} \mu(1,1) & \mu(1,2) \\ \mu(2,1) & \mu(2,2) \end{pmatrix} := \begin{pmatrix} 1 & 0 \\ 0 & 0 \end{pmatrix} \text{ and } \begin{pmatrix} \nu(1,1) & \nu(1,2) \\ \nu(2,1) & \nu(2,2) \end{pmatrix} := \begin{pmatrix} 0 & 0 \\ 0 & 1 \end{pmatrix}.$$

Define  $\nu_p := p\mu + (1-p)\nu$ . Show that

$$[0,1] \ni p \mapsto H(\nu_p | \nu_p^1 \otimes P)$$

is an affine function. Prove the same statement for

$$\mu := \begin{pmatrix} 0 & 0 & \frac{1}{2} & 0 \\ 0 & 0 & 0 & 0 \\ \frac{1}{2} & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix} \quad \text{and} \quad \nu := \begin{pmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & \frac{1}{2} \\ 0 & 0 & 0 & 0 \\ 0 & \frac{1}{2} & 0 & 0 \end{pmatrix}.$$

These examples show that  $\mathcal{M}_1(S^2) \ni \nu \mapsto H(\nu | \nu^1 \otimes P)$  is not strictly convex. Do you see a general pattern how to create such examples? Hint: Excercise 3.10.

**Exercise 3.15 (Probability to stay inside a set)** Let P be a probability kernel on  $\{0, 1, \ldots, n\}$   $(n \ge 1)$  such that P(x, y) > 0 for all  $1 \le x \le n$  and  $0 \le y \le n$  but P(0, y) = 0 for all  $1 \le y \le n$ . (In particular, 0 is a *trap* of the Markov chain with transition kernel P.) Show that there exists a constant  $0 < \lambda < \infty$  such that the Markov chain  $(X_k)_{k\ge 0}$  with transition kernel P and initial state  $X_0 = z \ge 1$  satisfies

$$\lim_{n \to \infty} \frac{1}{n} \log \mathbb{P}[X_n \ge 1] = -\lambda.$$

Give a (formal) expression for  $\lambda$  and show that  $\lambda$  does not depend on z. Hint: Corollary 3.9.

## 3.3 The empirical process

In this section, we return to the i.i.d. setting, but rather than looking at the (standard) empirical distributions as we did in Section 2.4, we will look at pair empirical distributions and more general at empirical distributions of k-tuples. Since i.i.d. sequences are a special case of Markov processes, our results from the previous section immediately give us the following theorem.

### 3.3. THE EMPIRICAL PROCESS

#### Theorem 3.16 (Sanov for pair empirical distributions)

(a) Let S be a finite set and let  $\mu$  be a probability measure on S such that  $\mu(x) > 0$ for all  $x \in S$ . Let  $(X_k)_{k\geq 0}$  be i.i.d. with common law  $\mu$  and let  $M_n^{(2)}$  be their pair empirical distributions as defined in (3.1). Then the laws  $\mathbb{P}[M_n^{(2)} \in \cdot]$  satisfy the large deviation principle with speed n and rate function  $I^{(2)}$  given by

$$I^{(2)}(\nu) := \begin{cases} H(\nu|\nu^1 \otimes \mu) & \text{if } \nu^1 = \nu^2, \\ \infty & \text{otherwise,} \end{cases}$$

where  $\nu^1$  and  $\nu^2$  denote the first and second marginal of  $\nu$ , respectively, and  $H(\cdot | \cdot)$  denotes the relative entropy of one measure w.r.t. another.

(b) More generally, if  $U \subset S^2$  is irreducible, then the restricted measures

$$\mathbb{P}[M_n^{(2)} \in \cdot]\big|_{\mathcal{M}_1(U)}$$

satisfy the large deviation principle with speed n and rate function  $I^{(2)}$  given by

$$I^{(2)}(\nu) := \begin{cases} H(\nu | [\nu^1 \otimes \mu]_U) & \text{if } \nu^1 = \nu^2, \\ \infty & \text{otherwise,} \end{cases}$$

where  $[\nu^1 \otimes \mu]_U$  denotes the restriction of the product measure  $\nu^1 \otimes \mu$  to U.

**Proof** Immediate from Theorem 3.2 and Corollary 3.9.

**Exercise 3.17 (Sanov's theorem)** Show that through the contraction principle, Theorem 3.16 (a) implies Sanov's theorem (Theorem 2.18) for finite state spaces.

Although Theorem 3.16, which is a statement about i.i.d. sequences only, looks more special that Theorem 3.2 and Corollary 3.9 which apply to general Markov chains, the two results are in fact more or less equivalent.

**Derivation of Theorem 3.2 from Theorem 3.16** We first consider the special case that  $P(x_1, x_2) > 0$  for all  $(x_1, x_2) \in S^2$ . Let  $\rho$  be the initial law of X, let  $\mu$  be any probability measure on S satisfying  $\mu(x) > 0$  for all  $x \in S$ , and let  $\hat{X} = (\hat{X}_k)_{k\geq 0}$  be independent random variables such that  $\hat{X}_0$  has law  $\rho$  and  $\hat{X}_k$  has law  $\mu$  for all  $k \geq 1$ . For any  $x = (x_k)_{k\geq 0}$  with  $x_k \in S$   $(k \geq 0)$ , let us define  $M_n^{(2)}(x) \in \mathcal{M}_1(S^2)$  by

$$M_n^{(2)}(x)(y_1, y_2) := \frac{1}{n} \sum_{k=1}^n \mathbb{1}_{\{(x_{k-1}, x_k) = (y_1, y_2)\}}$$

We observe that

$$\mathbb{P}[X_0 = x_0, \dots, X_n = x_n] = \rho(x_0) e^{\sum_{k=1}^n \log P(x_{k-1}, x_k)}$$
$$= \rho(x_0) e^{n \sum_{(y_1, y_2) \in S^2} \log P(y_1, y_2) M_n^{(2)}(x)(y_1, y_2)},$$

while

$$\mathbb{P}[\hat{X}_0 = x_0, \dots, \hat{X}_n = x_n] = \rho(x_0) e^{\sum_{k=1}^n \log \mu(x_k)}$$
$$= \rho(x_0) e^{n \sum_{(y_1, y_2) \in S^2} \log \mu(y_2) M_n^{(2)}(x)(y_1, y_2)}.$$

It follows that the Radon-Nikodym derivative of  $\mathbb{P}[M_n^{(2)}(X) \in \cdot]$  with respect to  $\mathbb{P}[M_n^{(2)}(\hat{X}) \in \cdot]$  is given by

$$\frac{\mathbb{P}[M_n^{(2)}(X) = \nu]}{\mathbb{P}[M_n^{(2)}(\hat{X}) = \nu]} = e^{n \sum_{(y_1, y_2) \in S^2} \left(\log P(y_1, y_2) - \log \mu(y_2)\right) \nu(y_1, y_2)}.$$

By Theorem 3.16 (a), the laws  $\mathbb{P}[M_n^{(2)}(\hat{X}) \in \cdot]$  satisfy the large deviation principle with speed n and rate function  $\hat{I}^{(2)}$  given by

$$\hat{I}^{(2)}(\nu) = \begin{cases} H(\nu|\nu^1 \otimes \mu) & \text{if } \nu^1 = \nu^2, \\ \infty & \text{if } \nu^1 \neq \nu^2. \end{cases}$$

Applying Lemma 1.15 to the function

$$F(\nu) := \sum_{(y_1, y_2) \in S^2} \left( \log P(y_1, y_2) - \log \mu(y_2) \right) \nu(y_1, y_2),$$

which is continuous by our assumption that  $P(y_1, y_2) > 0$  for all  $y_1, y_2 \in S$ , we find that the laws  $\mathbb{P}[M_n^{(2)}(\hat{X}) \in \cdot]$  satisfy the large deviation principle with speed  $s_n$  and rate function  $I^{(2)} = \hat{I}^{(2)} - F$ . Since

$$H(\nu|\nu^{1} \otimes \mu) - F(\nu)$$

$$= \sum_{(y_{1}, y_{2}) \in S^{2}} \nu(y_{1}, y_{2}) \left( \log \frac{\nu(y_{1}, y_{2})}{\nu^{1}(y_{1})\mu(y_{2})} + \log \mu(y_{2}) - \log P(y_{1}, y_{2}) \right)$$

$$= \sum_{(y_{1}, y_{2}) \in S^{2}} \nu(y_{1}, y_{2}) \log \frac{\nu(y_{1}, y_{2})}{\nu^{1}(y_{1})P(y_{1}, y_{2})} = H(\nu|\nu^{1} \otimes P),$$

this proves the theorem.

In the general case, when P is irreducible but not everywhere positive, the argument is the same but we need to apply Theorem 3.16 (b) to  $U := \{(x_1, x_2) \in S^2 : P(x_1, x_2) > 0\}$ , and we use that the function F restricted to  $\mathcal{M}_1(U)$  is continuous, hence Lemma 1.15 is applicable.

**Exercise 3.18 (Periodic boundary conditions)** Let  $(X_k)_{k\geq 0}$  be i.i.d. with common law  $\mu \in \mathcal{M}_1(S)$ . Let  $\mathcal{M}_n^{(2)}$  be the pair empirical distributions defined in (3.1) and set

$$\tilde{M}_{n}^{(2)} := \frac{1}{n} \tilde{N}_{n}^{(2)}, \quad \text{where}$$

$$\tilde{N}_{n}^{(2)}(x) := 1\{(X_{n}, X_{1}) = (x_{1}, x_{2})\} + \sum_{k=2}^{n} 1\{(X_{k-1}, X_{k}) = (x_{1}, x_{2})\}$$
(3.12)

Show that the random variables  $M_n^{(2)}$  and  $\tilde{M}_n^{(2)}$  are exponentially close in the sense of (1.8), hence by Proposition 1.17, proving a large deviation principle for the  $M_n^{(2)}$  is equivalent to proving one for the  $\tilde{M}_n^{(2)}$ .

**Remark** Den Hollander [Hol00, Thm II.8] who again follows [Ell85, Sect. I.5], gives a very nice and short proof of Sanov's theorem for the pair empirical distributions using periodic boundary conditions. The advantage of this approach is that the pair empirical distributions  $\tilde{M}_n^{(2)}$  defined in (3.12) automatically have the property that their first and second marginals agree, which means that one does not need to prove formula (3.4).

Based on this, along the lines of the proof above, Den Hollander [Hol00, Thm IV.3] then derives Theorem 3.2 in the special case that the transition kernel P is everywhere positive. In [Hol00, Comment (4) from Section IV.3], it is then claimed that the theorem still applies when P is not everywhere positive but irreducible and  $S^2$  is replaced by  $U := \{(x_1, x_2) \in S^2 : P(x_1, x_2) > 0\}$ , and 'the proof is easily adapted'. This last comment seems to be quite far from the truth. At least, I do not see any *easy* way to adapt his proof. The reason is that periodic boundary conditions do not work well anymore if  $S^2$  is replaced by a more general subset  $U \subset S^2$ . As a result, the technicalities needed to prove the analogue of Lemma 3.4 in a set-up with periodic boundary conditions become very unpleasant. Although a proof along these lines is possible, this seems to be more complicated than the approach used in these lecture notes.

The fact that Theorem 3.2 can rather easily be derived from Theorem 3.16 shows that the point of view that Chapter 2 is about large deviations of independent random variables while the present chapter is about large deviations of Markov chains is naive. With equal right, we might say that both chapters are concerned with large deviations of functions of i.i.d. random variables. The essential difference is in what kind of functions we consider. In Chapter 2, we considered the empirical distributions and functions thereof (such as the mean), while in the present chapter we consider the pair empirical distributions. By looking at yet different functions of i.i.d. random variables one can obtain a lot of very different, often difficult, but interesting large deviation principles.

There is no need to restrict ourselves to pairs. In fact, once we have a theorem for pairs, the step to general *m*-tuples is easy. (In contrast, there seems to be no easy way to derive the result for pairs from the large deviation principle for singletons.)

**Theorem 3.19 (Sanov for empirical distributions of** *m***-tuples)** Let *S* be a finite set and let  $\mu$  be a probability measure on *S* such that  $\mu(x) > 0$  for all  $x \in S$ . Let  $(X_k)_{k>1}$  be i.i.d. with common law  $\mu$  and for fixed  $m \ge 1$ , define

$$M_n^{(m)}(x) := \frac{1}{n} \sum_{k=0}^{n-1} \mathbb{1}_{\{(X_{k+1}, \dots, X_{k+m}) = x\}} \qquad (x \in S^m, \ n \ge 1).$$

Then the laws  $\mathbb{P}[M_n^{(m)} \in \cdot]$  satisfy the large deviation principle with speed n and rate function  $I^{(m)}$  given by

$$I^{(m)}(\nu) := \begin{cases} H(\nu | \nu^{\{1,...,m-1\}} \otimes \mu) & \text{if } \nu^{\{1,...,m-1\}} = \nu^{\{2,...,m\}}, \\ \infty & \text{otherwise,} \end{cases}$$

where  $\nu^{\{1,\dots,m-1\}}$  and  $\nu^{\{2,\dots,m\}}$  denote the projections of  $\nu$  on its first m-1 and last m-1 coordinates, respectively.

**Proof** The statement for m = 1, 2 has already been proved in Theorems 2.18 and 3.2, respectively, so we may assume that  $m \ge 3$ . Define a probability kernel  $P: S^{m-1} \to S^{m-1}$  by

$$P(x,y) := 1_{\{(x_2,\ldots,x_{m-1}) = (y_1,\ldots,y_{m-2})\}} \mu(y_{m-1}) \qquad (x,y \in S^{m-1}),$$

and set

$$\vec{X}_k := \left( X_{k+1}, \dots, X_{k+m-1} \right) \qquad (k \ge 0)$$

Then  $\vec{X} = (\vec{X}_k)_{k\geq 0}$  is a Markov chain with irreducible transition kernel *P*. By Theorem 3.2, the pair empirical distributions  $\vec{M}_n^{(2)}$  of  $\vec{X}$  satisfy a large deviation
principle. Here the  $\vec{M}_n^{(2)}$  take values in the space  $\mathcal{M}_1(S^{m-1} \times S^{m-1})$  and the rate function is given by

$$\vec{I}^{(2)}(\rho) := \begin{cases} H(\rho|\rho^1 \otimes P) & \text{if } \rho^1 = \rho^2, \\ \infty & \text{otherwise,} \end{cases}$$

where  $\rho^1$  and  $\rho^2$  denote the first and second marginals of  $\rho$ , respectively. (Note that  $\rho$  is a probability measure on  $S^{m-1} \times S^{m-1}$ , hence  $\rho^1$  and  $\rho^2$  are probability measures on  $S^{m-1}$ .)

Define a map  $\psi: S^m \to S^{m-1} \times S^{m-1}$  by

$$\psi(x_1,\ldots,x_m) := ((x_1,\ldots,x_{m-1}),(x_2,\ldots,x_m)).$$

The image of  $S^m$  under  $\psi$  is the set

$$U := \{ (x, y) \in S^{m-1} \times S^{m-1} : (x_2, \dots, x_{m-1}) = (y_1, \dots, y_{m-2}) \}$$
  
=  $\{ (x, y) \in S^{m-1} \times S^{m-1} : P(x, y) > 0 \}.$ 

It follows that  $\vec{I}^{(2)}(\rho) = \infty$  unless  $\rho \in \mathcal{M}_1(U)$ . Since  $\psi : S^m \to U$  is a bijection, each  $\rho \in \mathcal{M}_1(U)$  is the image under  $\psi$  of a unique  $\nu \in \mathcal{M}_1(S^m)$ . Moreover,  $\rho^1 = \rho^2$  if and only if  $\nu^{\{1,\dots,m-1\}} = \nu^{\{2,\dots,m\}}$ . Thus, by the contraction principle (Proposition 1.14), our claim will follow provided we show that if  $\nu \in \mathcal{M}_1(S^m)$ satisfies  $\nu^{\{1,\dots,m-1\}} = \nu^{\{2,\dots,m\}}$  and  $\rho = \nu \circ \psi^{-1}$  is the image of  $\nu$  under  $\rho$ , then

$$H(\nu|\nu^{\{1,\dots,m-1\}} \otimes \mu) = H(\rho|\rho^1 \otimes P).$$

Here

$$H(\rho|\rho^{1} \otimes P) = \sum_{\substack{x_{1},\dots,x_{m-1}\\y_{1},\dots,y_{m-1}}} \rho(x_{1},\dots,x_{m-1},y_{1},\dots,y_{m-1}) \times \Big(\log\rho(x_{1},\dots,x_{m-1},y_{1},\dots,y_{m-1}) - \log\rho^{1}(x_{1},\dots,x_{m-1}) - \log P(x_{1},\dots,x_{m-1},y_{1},\dots,y_{m-1})\Big),$$

where

$$\rho(x_1, \dots, x_{m-1}, y_1, \dots, y_{m-1}) = 1_{\{(x_2, \dots, x_{m-1}) = (y_1, \dots, y_{m-2})\}} \nu(x_1, \dots, x_{m-1}, y_{m-1}),$$
  
$$\rho^1(x_1, \dots, x_{m-1}) = \nu^{\{1, \dots, m-1\}}(x_1, \dots, x_{m-1}),$$

and

$$P(x_1,\ldots,x_{m-1},y_1,\ldots,y_{m-1}) = 1_{\{(x_2,\ldots,x_{m-1})=(y_1,\ldots,y_{m-2})\}}\mu(y_{m-1}).$$

It follows that

$$H(\rho|\rho^{1} \otimes P) = \sum_{x_{1},\dots,x_{m-1},y_{m-1}} \nu(x_{1},\dots,x_{m-1},y_{m-1})$$
  
  $\times \left(\log \nu(x_{1},\dots,x_{m-1},y_{m-1}) - \log \nu^{\{1,\dots,m-1\}}(x_{1},\dots,x_{m-1}) - \log \mu(y_{m-1})\right)$   
  $= H(\nu|\nu^{\{1,\dots,m-1\}} \otimes \mu).$ 

It is even possible to go one step further than Theorem 3.19 and prove a large deviations result for '*m*-tuples' with  $m = \infty$ . Let  $S^{\mathbb{N}}$  be the space of all infinite sequence  $x = (x_k)_{k\geq 0}$  with  $x \in S$ . Note that  $S^{\mathbb{N}}$ , equipped with the product topology, is a compact metrizable space. Define a shift operator  $\theta : S^{\mathbb{N}} \to S^{\mathbb{N}}$  by

$$(\theta x)_k := x_{k+1} \qquad (k \ge 0).$$

Let  $X = (X_k)_{k\geq 0}$  be i.i.d. random variables with values in S and common law  $\mu$  satisfying  $\mu(x) > 0$  for all  $x \in S$ . For each  $n \geq 1$ , we define a random measure  $M_n^{(\infty)}$  on  $S^{\mathbb{N}}$  by

$$M_n^{(\infty)} := \frac{1}{n} \sum_{k=0}^{n-1} \delta_{\theta^k X},$$

where  $\delta_x$  denotes the delta measure at a point x. We call  $M_n^{(\infty)}$  the *empirical* process.

**Exercise 3.20 (Empirical process)** Sketch a proof of the fact that the laws  $\mathbb{P}[M_n^{(\infty)} \in \cdot]$  satisfy a large deviation principle. Hint: projective limit.

**Exercise 3.21 (First occurrence of a pattern)** Let  $(X_k)_{k\geq 0}$  be i.i.d. random variables with  $\mathbb{P}[X_k=0] = \mathbb{P}[X_k=1] = \frac{1}{2}$ . Give a formal expression for the limits

$$\lambda_{001} := \lim_{n \to \infty} \frac{1}{n} \log \mathbb{P} \big[ (X_k, X_{k+1}, X_{k+2}) \neq (0, 0, 1) \; \forall k = 1, \dots, n \big]$$

and

$$\lambda_{000} := \lim_{n \to \infty} \frac{1}{n} \log \mathbb{P} \big[ (X_k, X_{k+1}, X_{k+2}) \neq (0, 0, 0) \ \forall k = 1, \dots, n \big].$$

## **3.4** Perron-Frobenius eigenvalues

In excercises such as Excercise 3.21, we need an explicit way to determine the exponential rates associated with certain events or expectations of exponential functions in the spirit of Varadhan's lemma. In this section, we will see that such rates are given by the Perron-Frobenius eigenvalue of a suitably chosen irreducible, nonnegative matrix.

We start by recalling the classical Perron-Frobenius theorem. Let S be a finite set  $(S = \{1, \ldots, n\}$  in the traditional formulation of the Perron-Frobenius theorem) and let  $A : S \times S \to \mathbb{R}$  be a function. We view such functions a matrices, equipped with the usual matrix product, or equivalently we identify A with the linear operator  $A : \mathbb{R}^S \to \mathbb{R}^S$  given by  $Af(x) := \sum_{y \in S} A(x, y)f(y)$ . We say that A is nonnegative if  $A(x, y) \ge 0$  for all  $x, y \in S$ . A nonnegative matrix A is called *irreducible* if for each  $x, y \in S$  there exists an  $n \ge 1$  such that  $A^n(x, y) > 0$ . Note that for probability kernels, this coincides with our earlier definition of irreducibility. We let  $\sigma(A)$  denote the spectrum of A, i.e., the collection of (possibly complex) eigenvalues of A, and we let  $\rho(A)$  denote its spectral radius

$$\rho(A) := \sup\{|\lambda| : \lambda \in \sigma(A)\}.$$

If  $\|\cdot\|$  is any norm on  $\mathbb{R}^S$ , then we define the associated *operator norm*  $\|A\|$  of A as

$$||A|| := \sup\{||Af|| : f \in \mathbb{R}^S, ||f|| = 1\}.$$

It is well-known that for any such operator norm

$$\rho(A) = \lim_{n \to \infty} \|A^n\|^{1/n}.$$
(3.13)

We cite the following version of the Perron-Frobenius theorem from [Gan00, Section 8.3] (see also, e.g., [Sen73, Chapter 1]).

**Theorem 3.22 (Perron-Frobenius)** Let S be a finite set and let  $A : \mathbb{R}^S \to \mathbb{R}^S$  be a linear operator whose matrix is nonnegative and irreducible. Then

- (i) There exist an  $f: S \to \mathbb{R}$ , unique up to multiplication by positive constants, and a unique  $\alpha \in \mathbb{R}$  such that  $Af = \alpha f$  and  $f(x) \ge 0$ .
- (ii) f(x) > 0 for all  $x \in S$ .
- (iii)  $\alpha = \rho(A) > 0.$

(iv) The algebraic multiplicity of  $\alpha$  is one. In particular, if A is written in its Jordan normal form, then  $\alpha$  corresponds to a block of size  $1 \times 1$ .

**Remark** If A is moreover aperiodic, then there exists some  $n \ge 1$  such that  $A^n(x, y) > 0$  for all  $x, y \in S$ . Now Perron's theorem [Gan00, Section 8.2] implies that all other eigenvalues  $\lambda$  of A satisfy  $|\lambda| < \alpha$ . If A is not aperiodic, then it is easy to see that this statement fails in general. (This is stated incorrectly in [DZ98, Thm 3.1.1 (b)].)

We call the constant  $\alpha$  and function f from Theorem 3.22 the *Perron-Frobenius* eigenvalue and eigenfunction of A, respectively. We note that if  $A^{\dagger}(x, y) := A(y, x)$ denotes the transpose of A, then  $A^{\dagger}$  is also nonnegative and irreducible. It is wellknown that the spectra of a matrix and its transpose agree:  $\sigma(A) = \sigma(A^{\dagger})$ , and therefore also  $\rho(A) = \rho(A^{\dagger})$ , which implies that the Perron-Frobenius eigenvalues of A and  $A^{\dagger}$  are the same. The same is usually not true for the corresponding Perron-Frobenius eigenvectors. We call eigenvectors of A and  $A^{\dagger}$  also right and left eigenvectors, respectively.

The main aim of the present section is to prove the following result.

**Theorem 3.23 (Exponential rate as eigenvalue)** Let  $X = (X_k)_{k\geq 0}$  be a Markov chain with finite state space S, irreducible transition kernel P, and arbitrary initial law. Let  $\phi : S^2 \to [-\infty, \infty)$  be a function such that

$$U := \{(x,y) \in S^2 : \phi(x,y) > -\infty\} \subset \{(x,y) \in S^2 : P(x,y) > 0\}$$

is irreducible, and let  $\overline{U}$  be as in (3.3). Then, provided that  $X_0 \in \overline{U}$  a.s., one has

$$\lim_{n \to \infty} \frac{1}{n} \log \mathbb{E}\left[e^{\sum_{k=1}^{n} \phi(X_{k-1}, X_k)}\right] = r,$$

where  $e^r$  is the Perron-Frobenius eigenvalue of the nonnegative, irreducible matrix A defined by

$$A(x,y) := P(x,y)e^{\phi(x,y)} \qquad (x,y \in \overline{U}). \tag{3.14}$$

We start with some preparatory lemmas. The next lemma shows that there is a close connection between Perron-Frobenius theory and Markov chains.

**Lemma 3.24 (Perron-Frobenius Markov chain)** Let S be a finite set and let  $A : \mathbb{R}^S \to \mathbb{R}^S$  be a linear operator whose matrix is nonnegative and irreducible. Let

#### 3.4. PERRON-FROBENIUS EIGENVALUES

 $\alpha, \eta$  and h be its associated Perron-Frobenius eigenvalue and left and right eigenvectors, respectively, i.e.,  $\eta A = \alpha \eta$ ,  $Ah = \alpha h$ ,  $\eta, h > 0$ . Choose any normalization such that  $\sum_{x} h(x)\eta(x) = 1$ . Then the matrix

$$A_h(x,y) := \frac{A(x,y)h(y)}{\alpha h(x)} \qquad (x,y \in S)$$

$$(3.15)$$

is an irreducible probability kernel on S and  $h\eta$  is its unique invariant law.

**Proof** Recall from Theorem 3.22 that h is strictly positive, hence  $A_h$  is well-defined. Since

$$\sum_{y \in S} A_h(x, y) = \sum_{y \in S} \frac{A(x, y)h(y)}{\alpha h(x)} = \frac{\alpha h(x)}{\alpha h(x)} = 1 \quad (x \in S),$$

we see that  $A_h$  is a probability kernel. Since  $A_h(x, y) > 0$  if and only if A(x, y) > 0, the kernel  $A_h$  is irreducible. Since

$$\sum_{x \in S} h(x)\eta(x)A_h(x,y) = \sum_{x \in S} h(x)\eta(x)\frac{A(x,y)h(y)}{\alpha h(x)}$$
$$= \alpha^{-1}\sum_{x \in S} \eta(x)A(x,y)h(y) = \eta(y)h(y),$$

we see that  $h\eta$  is an invariant law for  $A_h$ , and the only such invariant law by the irreducibility of the latter.

The following lemma is not only the key to proving Theorem 3.23, it also provides a link between Perron-Frobenius eigenvectors and entropy. In particular, in some special cases (such as Excercise 3.27), the following lemma can actually be used to obtain Perron-Frobenius eigenvectors by minimizing a suitable functional.

**Lemma 3.25 (Minimizer of weighted entropy)** Let S be a finite set, let P be a probability kernel on S and let  $\phi: S^2 \to [-\infty, \infty)$  be a function such that

$$U := \{ (x, y) \in S^2 : \phi(x, y) > -\infty \} \subset \{ (x, y) \in S^2 : P(x, y) > 0 \}$$

is irreducible. Let U be as in (3.3), define A as in (3.14), let  $\alpha = e^r$  be its Perron-Frobenius eigenvalue and let  $\eta, h > 0$  be the associated left and right eigenvectors, normalized such that  $\sum_{x \in \overline{U}} h(x)\eta(x) = 1$ . Let  $A_h$  be the probability kernel defined in (3.15) and let  $\pi := h\eta$  be its unique invariant law. Let  $\mathcal{V} := \{\nu \in \mathcal{M}_1(S^2) : \nu^1 = \nu^2\}$ . Then the function

$$G_{\phi}(\nu) := \nu \phi - H(\nu | \nu^1 \otimes P)$$

satisfies  $G_{\phi}(\nu) \leq r \ (\nu \in \mathcal{V})$ , with equality if and only if  $\nu = \pi \otimes A_h$ .

**Proof** We have  $G_{\phi}(\nu) = -\infty$  if  $\nu(x_1, x_2) > 0$  for some  $(x_1, x_2) \notin U$ . On the other hand, for  $\nu \in \mathcal{V}(U)$ , we observe that

$$\begin{split} \nu\phi &- H(\nu|\nu^1 \otimes P) \\ &= \sum_{(x_1, x_2) \in U} \nu(x_1, x_2) \phi(x_1, x_2) - \sum_{(x_1, x_2) \in U} \nu(x_1, x_2) \log \frac{\nu(x_1, x_2)}{\nu^1(x_1) P(x_1, x_2)} \\ &= \sum_{(x_1, x_2) \in U} \nu(x_1, x_2) \Big( \phi(x_1, x_2) - \log \nu(x_1, x_2) + \log \nu^1(x_1) + \log P(x_1, x_2) \Big) \\ &= \sum_{(x_1, x_2) \in U} \nu(x_1, x_2) \Big( -\log \nu(x_1, x_2) + \log \nu^1(x_1) + \log A(x_1, x_2) \Big) \\ &= \sum_{(x_1, x_2) \in U} \nu(x_1, x_2) \Big( -\log \nu(x_1, x_2) + \log \nu^1(x_1) + \log A(x_1, x_2) \Big) \\ &+ \log \alpha + \log h(x_1) - \log h(x_2) \Big) \end{split}$$

 $= \log \alpha - H(\nu | \nu^1 \otimes A_h),$ 

where in the last step we have used that  $\nu^1 = \nu^2$ . Now the statement follows from Excercise 3.11.

**Proof of Theorem 3.23** We will deduce the claim from our basic large deviations results for Markov chains (Theorem 3.2 and Corollary 3.9). A direct proof (using a bit of matrix theory) is also possible, but our aim is to exhibit the links with our earlier results. In fact, the calculations below can be reversed, i.e., a direct proof of Theorem 3.23 can be used as the basis for an alternative proof of Theorem 3.2; see [Hol00, Section V.4].

Let  $M_n^{(2)}$  be the pair empirical distributions associated with X, defined in (3.1). Let

$$\mathcal{M}_1(U) \ni \nu \mapsto F(\nu) \in [-\infty, \infty)$$

be the continuous and bounded from above. Then, by Varadhan's lemma (Lemma 1.12) and Corollary 3.9,

$$\lim_{n \to \infty} \frac{1}{n} \log \int \mathbb{P}[M_n^{(2)} \in \mathrm{d}\nu] \big|_{\mathcal{M}_1(U)} e^{nF(\nu)} = \sup_{\nu \in \mathcal{M}_1(U)} \left[ F(\nu) - \tilde{I}^{(2)}(\nu) \right],$$

where  $\tilde{I}^{(2)}$  is the rate function from Corollary 3.9. A simpler way of writing this formula is

$$\lim_{n \to \infty} \frac{1}{n} \log \int \mathbb{E} \left[ e^{n F(M_n^{(2)})} \right] = \sup_{\nu \in \mathcal{M}_1(S^2)} \left[ F(\nu) - I^{(2)}(\nu) \right], \tag{3.16}$$

where  $I^{(2)}$  is the rate function from Theorem 3.2 and we have extended F to a function on  $\mathcal{M}_1(S^2)$  by setting  $F(\nu) := -\infty$  if  $\nu \in \mathcal{M}_1(S^2) \setminus \mathcal{M}_1(U)$ .

Applying this to the 'linear' function F defined by

$$F(\nu) := \nu \phi = \sum_{x \in S} \nu(x) \phi(x) \qquad \left(\nu \in \mathcal{M}_1(S^2)\right),$$

formula (3.16) tells us that

$$\lim_{n \to \infty} \frac{1}{n} \log \mathbb{E} \left[ e^{\sum_{k=1}^{n} \phi(X_{k-1}, X_k)} \right] = \sup_{\nu \in \mathcal{M}_1(S^2)} \left[ \nu \phi - I_n^{(2)}(\nu) \right]$$
$$= \sup_{\nu \in \mathcal{V}} \left[ \nu \phi - I_n^{(2)}(\nu) \right] = r,$$

where we have used that  $I^{(2)}(\nu) = H(\nu|\nu^1 \otimes P)$  for  $\nu \in \mathcal{V}$  and  $I^{(2)}(\nu) = \infty$  otherwise, and the final equality follows from Lemma 3.25.

Exercise 3.26 (First occurrence of a pattern: part 2) Let  $(X_k)_{k\geq 0}$  be i.i.d. random variables with  $\mathbb{P}[X_k = 0] = \mathbb{P}[X_k = 1] = \frac{1}{2}$ . Let  $\lambda_{001}$  be defined as in Excercise 3.21 and let

$$\lambda_{00} := \lim_{n \to \infty} \frac{1}{n} \log \mathbb{P} \big[ (X_k, X_{k+1}) \neq (0, 0) \ \forall k = 1, \dots, n \big]$$

Prove that  $\lambda_{001} = \lambda_{00}$ .

**Exercise 3.27 (First occurrence of a pattern: part 3)** Consider a Markov chain  $Z = (Z_k)_{k\geq 0}$  taking values in the space

$$S := \{ \underline{1}, \underline{10}, \underline{10}, \underline{100}, \underline{100}, \underline{100}, \underline{100}, \dagger \},\$$

that evolves according to the following rules:

 $\frac{\underline{10} \mapsto \underline{10}}{\underline{100} \mapsto \underline{100} \mapsto \underline{100}} \ \right\} \text{ with probability one,}$ 

$$\begin{array}{c} \underline{1}\\ \underline{10}\\ \underline{100} \end{array} \right\} \mapsto \begin{cases} \underline{1} & \text{with probability } 2^{-1}, \\ \underline{10} & \text{with probability } 2^{-2}, \\ \underline{100} & \text{with probability } 2^{-3}, \\ \frac{1}{7} & \text{with probability } 2^{-3}, \end{cases}$$

and

i.e., from each of the states  $\underline{1}, \underline{10}, \underline{100}$ , we jump with probability  $\frac{1}{2}$  to  $\underline{1}$ , with probability  $\frac{1}{4}$  to  $\underline{10}$ , with probability  $\frac{1}{8}$  to  $\underline{100}$ , and with probability  $\frac{1}{8}$  to  $\frac{1}{100}$ . The state  $\frac{1}{100}$ , finally, is a trap:

 $\dagger \mapsto \dagger$  with probability one.

Define  $\phi: S \times S \to [-\infty, \infty)$  by

$$\phi(x,y) := \begin{cases} 0 & \text{if } P(x,y) > 0 \text{ and } y \neq \dagger, \\ -\infty & \text{otherwise.} \end{cases}$$

Let  $\theta$  be the unique solution in the interval [0, 1] of the equation

$$\theta + \theta^2 + \theta^3 = 1.$$

and let  $\tilde{Z} = (\tilde{Z}_k)_{k \geq 0}$  be a Markov chain with state space  $S \setminus \{\dagger\}$  that evolves in the same way as Z, except that

$$\frac{1}{10} \\ 100 \\ 100 \\ \end{bmatrix} \mapsto \begin{cases} \frac{1}{10} & \text{with probability } \theta, \\ \frac{1}{10} & \text{with probability } \theta^2, \\ \frac{1}{100} & \text{with probability } \theta^3. \end{cases}$$

Let P and Q be the transition kernels of Z and  $\tilde{Z}$ , respectively. Set  $U := \{(x, y) \in S^2 : \phi(x, y) > -\infty\}$ . Prove that for any  $\nu \in \mathcal{V}(U)$ 

$$\nu\phi - H(\nu|\nu^1 \otimes P) = \log(\frac{1}{2}) - \log\theta - H(\nu|\nu^1 \otimes Q).$$
(3.17)

Hint: Do a calculation as in the proof of Lemma 3.25, and observe that for any  $\nu \in \mathcal{V}(U)$ 

$$\nu^{1}(\underline{1}1) = \nu^{1}(\underline{1}\underline{1})$$
 and  $\nu^{1}(\underline{1}11) = \nu^{1}(\underline{1}\underline{1}1) = \nu^{1}(\underline{1}1\underline{1}),$ 

hence  $\nu^{1}(\underline{1}) + 2\nu^{1}(\underline{1}1) + 3\nu^{1}(\underline{1}11) = 1.$ 

**Exercise 3.28 (First occurrence of a pattern: part 4)** Let  $(X_k)_{k\geq 0}$  be i.i.d. random variables with  $\mathbb{P}[X_k = 0] = \mathbb{P}[X_k = 1] = \frac{1}{2}$  and let  $\lambda_{000}$  be defined as in Excercise 3.21. Prove that  $\lambda_{000} = \log(\frac{1}{2}) - \log(\theta)$ , where  $\theta$  is the unique root of the equation  $\theta + \theta^2 + \theta^3 = 1$  in the interval [0, 1]. Hint: use formula (3.17).

**Exercise 3.29 (Percolation on a ladder)** Let  $0 and let <math>(Y_{i,k})_{i=1,2,3, k \ge 1}$  be i.i.d. Bernoulli random variables with  $\mathbb{P}[Y_{i,k} = 1] = p$  and  $\mathbb{P}[Y_{i,k} = 0] = 1 - p$ .

Let  $S := \{0,1\}^2$  and let  $x \in S \setminus \{(0,0)\}$  be fixed. Define inductively a Markov chain  $(X_k)_{k\geq 0}$  with state space S by first setting

$$\tilde{X}_k(1) := Y_{k,1}X_{k-1}(1)$$
 and  $\tilde{X}_k(2) := Y_{k,2}X_{k-1}(2),$ 

and then

$$X_k(1) := \tilde{X}_k(1) \lor Y_{3,k} \tilde{X}_k(2)$$
 and  $X_k(2) := \tilde{X}_k(2) \lor Y_{3,k} \tilde{X}_k(1).$ 

Calculate the limit

$$r := \lim_{n \to \infty} \frac{1}{n} \log \mathbb{P} \big[ X_n \neq (0, 0) \big].$$

Hint: find the transition kernel of X and calculate the relevant Perron-Frobenius eigenvalue. You can reduce the dimensionality of the problem by exploiting the symmetry between (1,0) and (0,1). Don't worry if the formula for r looks somewhat complicated.

## 3.5 Continuous time

Recall from Section 0.4 the definition of a continuous-time Markov process  $X = (X_t)_{t\geq 0}$  with finite state space S, initial law  $\mu$ , transition probabilities  $P_t(x, y)$ , semigroup  $(P_t)_{t\geq 0}$ , generator G, and transition rates r(x, y)  $(x \neq y)$ . To simplify notation, we set r(x, x) := 0.

By definition, an *invariant law* is a probability measure  $\rho$  on S such that  $\rho P_t = \rho$  for all  $t \ge 0$ , or, equivalently,  $\rho G = 0$ . This latter formula can be written more explicitly in terms of the rates r(x, y) as

$$\sum_{y \in S} \rho(y) r(y, x) = \rho(x) \sum_{y \in S} r(x, y) \qquad (x \in S),$$

i.e., in equilibrium, the frequency of jumps to x equals the frequency of jumps from x. Basic results about Markov processes with finite state spaces tell us that if the transition rates r(x, y) are irreducible, then the corresponding Markov process has a unique invariant law  $\rho$ , and  $\mu P_t \Rightarrow \rho$  as  $t \to \infty$  for every initial law  $\mu$ . (For continuous-time processes, there is no such concept a (a)periodicity.)

We let

$$M_T(x) := \frac{1}{T} \int_0^T \mathbf{1}_{\{X_t = x\}} dt \qquad (T > 0)$$

denote the *empirical distribution* of X up to time T. We denote the set of times when X makes a jump up to time T by

$$\Delta_T := \{ t \in (0, T] : X_{t-} \neq X_t \}$$

and we set

$$W_T(x,y) := \frac{1}{T} \sum_{t \in \Delta_T} \mathbb{1}\{X_{t-} = x, \ X_t = y\} \qquad (T > 0),$$

i.e.,  $W_T(x, y)$  is the *empirical frequency* of jumps from x to y. If the transition rates r(x, y) are irreducible, then, for large T, we expect  $M_T$  to be close to the (unique) invariant law  $\rho$  of X and we expect  $W_T(x, y)$  to be close to  $\rho(x)r(x, y)$ . We observe that  $(M_T, W_T)$  is a random variable taking values in the space  $\mathcal{M}_1(S) \times [0, \infty)^{S^2}$ . For any  $w \in [0, \infty)^{S^2}$ , we let

$$w^1(x_1) := \sum_{x_2 \in S} w(x_1, x_2)$$
 and  $w^2(x_2) := \sum_{x_1 \in S} w(x_1, x_2)$ 

denote the first and second marginal of w, and we set

$$\mathcal{W} := \{ (\rho, w) : \rho \in \mathcal{M}_1(S), \ w \in [0, \infty)^{S^2}, \ w^1 = w^2, \\ w(x, y) = 0 \text{ whenever } \rho(x)r(x, y) = 0 \}.$$

The aim of the present section is to prove the following analogue of Theorem 3.2.

**Theorem 3.30 (LDP for Markov processes)** Let  $(X_t)_{t\geq 0}$  be a continuous-time Markov process with finite state space S, irreducible transition rates r(x, y), and arbitrary initial law. Let  $M_T$  and  $W_T$  (T > 0) denote its empirical distributions and empirical frequencies of jumps, respectively, as defined above. Then the laws  $\mathbb{P}[(M_T, W_T) \in \cdot]$  satisfy the large deviation principle on  $\mathcal{M}_1(S) \times [0, \infty)^{S^2}$  with speed T and good rate function I given by

$$I(\rho, w) := \begin{cases} \sum_{x,y \in S} \rho(x) r(x,y) \psi\left(\frac{w(x,y)}{\rho(x) r(x,y)}\right) & \text{if } (\rho, w) \in \mathcal{W}, \\ \infty & \text{otherwise,} \end{cases}$$

where  $\psi(z) := 1 - z + z \log z$  (z > 0) and  $\psi(0) := 1$  and we set  $0 \psi(a/b) := 0$ , regardless of the values of  $a, b \ge 0$ .

**Remark** So far, we have only considered large deviation principles for sequences of measures  $\mu_n$ . The theory for families of measures  $(\mu_T)_{T>0}$  depending on a continuous parameter is completely analogous. Indeed, if the  $\mu_T$  are finite measures on a Polish space E and I is a good rate function, then one has

$$\lim_{T \to \infty} \|f\|_{T,\mu_T} = \|f\|_{\infty,I} \qquad \left(f \in \mathcal{C}_{b,+}(E)\right)$$

if and only if for each  $T_n \to \infty$ ,

$$\lim_{n \to \infty} \|f\|_{T_n, \mu_{T_n}} = \|f\|_{\infty, I} \qquad (f \in \mathcal{C}_{b, +}(E)).$$

A similar statement holds for the two conditions in Proposition 1.7. In other words: measures  $\mu_T$  depending on a continuous parameter T > 0 satisfy a large deviation principle with speed T and good rate function I if and only if for each  $T_n \to \infty$ , the measures  $\mu_{T_n}$  satisfy the large deviation principle with speed  $T_n$  and rate function I.

**Exercise 3.31 (Properties of the rate function)** Show that the function I from Theorem 3.30 is a good rate function and that  $I(\rho, w) \ge 0$  with equality if and only if  $\rho$  is the unique invariant law of the Markov process X and  $w(x, y) = \rho(x)r(x, y)$   $(x, y \in S)$ .

Our strategy is to derive Theorem 3.30 from Theorem 3.2 using approximation. We start with an abstract lemma.

**Lemma 3.32 (Diagonal argument)** Let  $(\mu_{m,n})_{m,n\geq 1}$  be finite measures on a Polish space E, let  $s_n$  be positive constants, tending to infinity, and let  $I_m$ , I be good rate functions on E. Assume that for each fixed  $m \geq 1$ , the  $\mu_{m,n}$  satisfy the large deviation principle with speed  $s_n$  and rate function  $I_m$ . Assume moreover that

$$\lim_{m \to \infty} \|f\|_{\infty, I_m} = \|f\|_{\infty, I} \qquad \left(f \in \mathcal{C}_{b, +}(E)\right).$$

Then there exist  $n(m) \to \infty$  such that for all  $n'(m) \ge n(m)$ , the measures  $\mu_{m,n'(m)}$  satisfy the large deviation principle with speed  $s_{n'(m)}$  and rate function I.

**Proof** Let  $\overline{E}$  be a metrizable compactification of E. We view the  $\mu_{m,n}$  as measures on  $\overline{E}$  such that  $\mu_{m,n}(\overline{E} \setminus E) = 0$  and we extend the rate functions  $I_m, I$  to  $\overline{E}$  by setting  $I_m, I := \infty$  on  $\overline{E} \setminus E$ . Then

$$\lim_{m \to \infty} \|f\|_{\infty, I_m} = \|f\|_{\infty, I} \qquad \left(f \in \mathcal{C}(\overline{E})\right).$$

Let  $\{f_i : i \geq 1\}$  be a countable dense subset of the separable Banach space  $\mathcal{C}(\overline{E})$  of continuous real functions on E, equipped with the supremumnorm. Choose  $n(m) \to \infty$  such that

$$\left| \|f_i\|_{s_{n'},\mu_{m,n'}} - \|f_i\|_{\infty,I_m} \right| \le 1/m \qquad (n' \ge n(m), \ i \le m).$$

Then, for any  $n'(m) \ge n(m)$ , one has

$$\begin{split} \limsup_{m \to \infty} |\|f_i\|_{s_{n'(m)}, \mu_{m,n'(m)}} - \|f_i\|_{\infty, I}| \\ \leq \limsup_{m \to \infty} |\|f_i\|_{s_{n'(m)}, \mu_{m,n'(m)}} - \|f_i\|_{\infty, I_m}| + \limsup_{m \to \infty} |\|f_i\|_{\infty, I_m} - \|f_i\|_{\infty, I}| = 0 \end{split}$$

for all  $i \geq 1$ . It is easy to see that the function  $\Lambda(f) := ||f||_{\infty,I}$  satisfies  $\Lambda(f) = \Lambda(|f|)$  and  $f \mapsto \Lambda(f)$  is continuous w.r.t. the supremumnorm (compare Proposition 1.21), and the same is true for the functions  $f \mapsto ||f||_{\infty,I_m}$ . Using this, it is easy to see that the functions  $|f_i|$  are rate function determining, hence by Proposition 1.29, the measures  $\mu_{m,n'(m)}$  satisfy the large deviation principle on  $\overline{E}$  with speed  $s_{n'(m)}$  and rate function I. By the restriction principle (Lemma 1.27), they also satisfy the large deviation principle on E.

**Proposition 3.33 (Approximation of LDP's)** Let E be a Polish space and let  $X_n, X_{m,n}$   $(m, n \ge 1)$  be random variables taking values in E. Assume that for each fixed  $m \ge 1$ , the laws  $\mathbb{P}[X_{m,n} \in \cdot]$  satisfy a large deviation principle with speed  $s_n$  and good rate function  $I_m$ . Assume moreover that there exists a good rate function I such that

$$\lim_{m \to \infty} \|f\|_{\infty, I_m} = \|f\|_{\infty, I} \qquad (f \in \mathcal{C}_{b, +}(E)), \tag{3.18}$$

and that there exists a metric d generating the topology on E such that for each  $n(m) \to \infty$ ,

$$\lim_{m \to \infty} \frac{1}{s_{n(m)}} \log \mathbb{P}[d(X_{n(m)}, X_{m,n(m)}) \ge \varepsilon] = -\infty \qquad (\varepsilon > 0), \tag{3.19}$$

i.e.,  $X_{n(m)}$  and  $X_{m,n(m)}$  are exponentially close in the sense of (1.8). Then the laws  $\mathbb{P}[X_n \in \cdot]$  satisfy the large deviation principle with speed  $s_n$  and good rate function I.

**Proof** By the argument used in the proof of Proposition 1.29, it suffices to show that each subsequence  $n(m) \to \infty$  contains a further subsequence  $n'(m) \to \infty$  such that the laws  $\mathbb{P}[X_{n'(m)} \in \cdot]$  satisfy the large deviation principle with speed  $s_{n'(m)}$ and good rate function *I*. By (3.18) and Lemma 3.32, we can choose  $n'(m) \to \infty$  such that the laws  $\mathbb{P}[X_{m,n'(m)} \in \cdot]$  satisfy the large deviation principle with speed  $s_{n'(m)}$  and good rate function I. By (3.19), the random variables  $X_{n'(m)}$  and  $X_{m,n'(m)}$  are exponentially close in the sense of Proposition 1.17, hence the large deviation principle for the laws of the  $X_{m,n'(m)}$  implies the large deviation principle for the laws of the  $X_{m,n'(m)}$  implies the large deviation principle for the laws of the  $X_{n'(m)}$ .

The following lemma gives sufficient conditions for the type of convergence in (3.18).

**Lemma 3.34 (Convergence of rate functions)** Let E be a Polish space and let  $I, I_m$  be good rate functions on E such that

- (i) For each  $a \in \mathbb{R}$ , there exists a compact set  $K \subset E$  such that  $\{x \in E : I_m(x) \leq a\} \subset K$  for all  $m \geq 1$ .
- (ii)  $\forall x_m, x \in E \text{ with } x_m \to x, \text{ one has } \liminf_{m \to \infty} I_m(x_m) \ge I(x).$
- (iii)  $\forall x \in E \; \exists x_m \in E \; such \; that \; x_m \to x \; and \; \limsup_{m \to \infty} I_m(x_m) \leq I(x).$

Then the  $I_m$  converge to I in the sense of (3.18).

**Proof** Formula (3.18) is equivalent to the statement that

$$\inf_{x \in E} [I_m(x) - F(x)] \underset{m \to \infty}{\longrightarrow} \inf_{x \in E} [I(x) - F(x)]$$

for any continuous  $F: E \to [-\infty, \infty)$  that is bounded from below. If  $I_m, I$  satisfy conditions (i)–(iii), then the same is true for I' := I - F,  $I'_m := I_m - F$ , so it suffices to show that conditions (i)–(iii) imply that

$$\inf_{x \in E} I_m(x) \xrightarrow[m \to \infty]{} \inf_{x \in E} I(x).$$

Since I is a good rate function, it achieves its minimum, i.e., there exists some  $x_{\circ} \in E$  such that  $I(x_{\circ}) = \inf_{x \in E} I(x)$ . By condition (iii), there exist  $x_m \in E$  such that  $x_m \to x$  and

$$\limsup_{m \to \infty} \inf_{x \in E} I_m(x) \le \limsup_{m \to \infty} I_m(x_m) \le I(x_\circ) = \inf_{x \in E} I(x).$$

To prove the other inequality, assume that

$$\liminf_{m \to \infty} \inf_{x \in E} I_m(x) < \inf_{x \in E} I(x).$$

Then, by going to a subsequence if necessary, we can find  $x_m \in E$  such that

$$\lim_{m \to \infty} I_m(x_m) < \inf_{x \in E} I(x),$$

where the limit on the left-hand side exists and may be  $-\infty$ . By condition (i), there exists a compact set  $K \subset E$  such that  $x_m \in K$  for all m, hence by going to a further subsequence if necessary, we may assume that  $x_m \to x_*$  for some  $x_* \in E$ . Condition (ii) now tells us that

$$\lim_{m \to \infty} I_m(x_m) \ge I(x_*) \ge \inf_{x \in E} I(x),$$

which leads to a contradiction.

#### Proof of Theorem 3.30 We set

$$M_T^{\varepsilon}(x) := \frac{1}{\lfloor T/\varepsilon \rfloor} \sum_{k=1}^{\lfloor T/\varepsilon \rfloor} 1\{(X_{\varepsilon(k-1)}, X_{\varepsilon k}) = (x, x)\} \qquad (x \in S),$$
$$W_T^{\varepsilon}(x, y) := \frac{1}{\varepsilon \lfloor T/\varepsilon \rfloor} \sum_{k=1}^{\lfloor T/\varepsilon \rfloor} 1\{(X_{\varepsilon(k-1)}, X_{\varepsilon k}) = (x, y)\} \qquad (x, y \in S, \ x \neq y),$$

and we let  $W_T^{\varepsilon}(x, x) := 0$  ( $x \in S$ ). By Proposition 3.33, the statements of the theorem will follow provided we prove the following three claims:

- 1. For each  $\varepsilon > 0$ , the laws  $\mathbb{P}[(M_T^{\varepsilon}, W_T^{\varepsilon}) \in \cdot]$  satisfy a large deviation principle with speed T and good rate function  $I_{\varepsilon}$ .
- 2. The function I from Theorem 3.30 is a good rate function and the rate functions  $I_{\varepsilon}$  converge to I in the sense of (3.18) as  $\varepsilon \downarrow 0$ .
- 3. For each  $T_m \to \infty$  and  $\varepsilon_m \downarrow 0$ , the random variables  $(M_{T_m}^{\varepsilon_m}, W_{T_m}^{\varepsilon_m})$  and  $(M_{T_m}, W_{T_m})$  are exponentially close with speed  $T_m$ .

Proof of Claim 1. For each  $\varepsilon > 0$ , let  $(X_k^{\varepsilon})_{k\geq 0}$  be the Markov chain given by

$$X_k^{\varepsilon} := X_{\varepsilon k} \qquad (k \ge 0)$$

and let  $M_n^{(2)\varepsilon}$  be its empirical pair distributions. Then

$$\begin{split} M_T^{\varepsilon}(x) &= M_{\lfloor T/\varepsilon \rfloor}^{(2) \varepsilon}(x, x) \qquad (x \in S), \\ W_T^{\varepsilon}(x, y) &= \varepsilon^{-1} M_{\lfloor T/\varepsilon \rfloor}^{(2) \varepsilon}(x, y) \qquad (x, y \in S, \ x \neq y) \end{split}$$

For each  $\varepsilon > 0$  and  $\nu \in \mathcal{M}_1(S^2)$ , let us define  $\rho_{\varepsilon} \in [0,\infty)^S$  and  $w_{\varepsilon}(\nu) \in [0,\infty)^{S^2}$ by

$$\rho_{\varepsilon}(\nu)(x) := \nu(x, x) \qquad (x \in S),$$
  
$$w_{\varepsilon}(\nu)(x) := \mathbf{1}_{\{x \neq y\}} \varepsilon^{-1} \nu(x, y) \qquad (x, y \in S)$$

Then, by Theorem 3.2, for each  $\varepsilon > 0$  the laws  $\mathbb{P}[(M_T^{\varepsilon}, W_T^{\varepsilon}) \in \cdot]$  satisfy a large deviation principle on  $[0, \infty)^S \times [0, \infty)^{S^2}$  with speed T and good rate function  $I_{\varepsilon}$  given by

$$I_{\varepsilon}(\rho_{\varepsilon}(\nu), w_{\varepsilon}(\nu)) := \varepsilon^{-1} H(\nu | \nu^{1} \otimes P_{\varepsilon}) \qquad (\nu \in \mathcal{V}),$$
(3.20)

while  $I_{\varepsilon}(\rho, w) := \infty$  if there exists no  $\nu \in \mathcal{V}$  such that  $(\rho, w) = (\rho_{\varepsilon}(\nu), w_{\varepsilon}(\nu))$ . Note the overall factor  $\varepsilon^{-1}$  which is due to the fact that the speed T differs a factor  $\varepsilon^{-1}$ from the speed n of the embedded Markov chain.

*Proof of Claim 2.* By Lemma 3.34, it suffices to prove, for any  $\varepsilon_n \downarrow 0$ , the following three statements.

- (i) If  $\rho_n \in [0,\infty)^S$  and  $w_n \in [0,\infty)^{S^2}$  satisfy  $w_n(x,y) \to \infty$  for some  $x, y \in S$ , then  $I_{\varepsilon_n}(\rho_n, w_n) \to \infty$ .
- (ii) If  $\rho_n \in [0,\infty)^S$  and  $w_n \in [0,\infty)^{S^2}$  satisfy  $(\rho_n, w_n) \to (\rho, w)$  for some  $\rho \in [0,\infty)^S$  and  $w \in [0,\infty)^{S^2}$ , then  $\liminf_{n\to\infty} I_{\varepsilon_n}(\rho_n, w_n) \ge I(\rho, w)$ .
- (iii) For each  $\rho \in [0,\infty)^S$  and  $w \in [0,\infty)^{S^2}$  there exist  $\rho_n \in [0,\infty)^S$  and  $w_n \in [0,\infty)^{S^2}$  such that  $\limsup_{n\to\infty} I_{\varepsilon_n}(\rho_n,w_n) \leq I(\rho,w)$ .

Obviously, it suffices to check conditions (i), (ii) for  $(\rho_n, w_n)$  such that  $I_{\varepsilon_n}(\rho_n, w_n) < \infty$  and condition (iii) for  $(\rho, w)$  such that  $I(\rho, w) < \infty$ . Therefore, taking into account our definition of  $I_{\varepsilon}$ , Claim 2 will follow provided we prove the following three subclaims.

2.I. If  $\nu_n \in \mathcal{V}$  satisfy  $\varepsilon_n^{-1}\nu_n(x,y) \to \infty$  for some  $x \neq y$ , then

$$\varepsilon_n^{-1} H(\nu_n | \nu_n^1 \otimes P_{\varepsilon_n}) \underset{n \to \infty}{\longrightarrow} \infty.$$

2.II. If  $\nu_n \in \mathcal{V}$  satisfy

$$\nu_n(x,x) \underset{n \to \infty}{\longrightarrow} \rho(x) \qquad (x \in S),$$
  

$$\varepsilon_n^{-1} \mathbb{1}_{\{x \neq y\}} \nu_n(x,y) \underset{n \to \infty}{\longrightarrow} w(x,y) \qquad (x,y \in S^2),$$
(3.21)

for some  $(\rho, w) \in [0, \infty)^S \times [0, \infty)^{S^2}$ , then  $\liminf_{n \to \infty} \varepsilon_n^{-1} H(\nu_n | \nu_n^1 \otimes P_{\varepsilon_n}) \ge I(\rho, w).$ 

2.III. For each  $(\rho, w) \in \mathcal{W}$ , we can find  $\nu_n \in \mathcal{V}$  satisfying (3.21) such that

$$\lim_{n \to \infty} \varepsilon_n^{-1} H(\nu_n | \nu_n^1 \otimes P_{\varepsilon_n}) = I(\rho, w).$$

We start by writing  $H(\nu|\nu^1 \otimes P)$  in a suitable way. Let  $\psi$  be as defined in the theorem. We observe that if  $\nu, \mu$  are probability measures on a finite set S and  $\mu(x) > 0$  for all  $x \in S$ , then

$$\sum_{x \in S} \mu(x)\psi\left(\frac{\nu(x)}{\mu(x)}\right) = \sum_{x \in S} \mu(x)\left[1 - \frac{\nu(x)}{\mu(x)} + \frac{\nu(x)}{\mu(x)}\log\left(\frac{\nu(x)}{\mu(x)}\right)\right]$$
$$= \sum_{x \in S} [\mu(x) - \nu(x)] + \sum_{x \in S} \nu(x)\log\left(\frac{\nu(x)}{\mu(x)}\right) = H(\nu|\mu),$$

where we use the convention that  $0 \log 0 := 0$ . By Excercise 3.10, it follows that for any probability measure  $\rho$  on S and probability kernels P, Q on S such that  $\rho \otimes Q \ll \rho \otimes P$ ,

$$H(\rho \otimes Q|\rho \otimes P) = \sum_{x} \rho(x)H(Q_{x}|P_{x})$$
$$= \sum_{x} \rho(x)\sum_{y} P(x,y)\psi\Big(\frac{Q(x,y)}{P(x,y)}\Big) = \sum_{x,y} \rho(x)P(x,y)\psi\Big(\frac{\rho(x)Q(x,y)}{\rho(x)P(x,y)}\Big),$$

where the sum runs over all  $x, y \in S$  such that  $\rho(x)P(x, y) > 0$ . In particular, if  $\nu$  is a probability measure on  $S^2$  and P is a probability kernel on S, then

$$H(\nu|\nu^{1} \otimes P) = \begin{cases} \sum_{x,y \in S} \nu^{1}(x)P(x,y)\psi\Big(\frac{\nu(x,y)}{\nu^{1}(x)P(x,y)}\Big) & \text{if } \nu \ll \nu^{1} \otimes P, \\ \infty & \text{otherwise,} \end{cases}$$

where we define  $0 \psi(a/b) := 0$ , irrespective of the values of  $a, b \ge 0$ .

To prove Claim 2.I, now, we observe that if  $\varepsilon_n^{-1}\nu_n(x,y) \to \infty$  for some  $x \neq y$ , then

$$\varepsilon_n^{-1} H(\nu_n | \nu_n^1 \otimes P_{\varepsilon_n}) \ge \varepsilon_n^{-1} \nu_n^1(x) P_{\varepsilon_n}(x, y) \psi\left(\frac{\nu_n(x, y)}{\nu_n^1(x) P_{\varepsilon_n}(x, y)}\right)$$
$$\ge \varepsilon_n^{-1} \nu_n(x, y) \Big( \log\left(\frac{\nu_n(x, y)}{\nu_n^1(x) P_{\varepsilon_n}(x, y)}\right) - 1 \Big),$$

where

$$\frac{\nu_n(x,y)}{\nu^1(x)P_{\varepsilon_n}(x,y)} \geq \frac{\nu_n(x,y)}{P_{\varepsilon_n}(x,y)} = \frac{\nu_n(x,y)}{\varepsilon_n r(x,y) + O(\varepsilon_n^2)} \xrightarrow[n \to \infty]{} \infty.$$

To prove Claim 2.II, we observe that if  $\nu_n, \rho, w$  satisfy (3.21), then, as  $n \to \infty$ ,

$$\nu_n^1(x)P_{\varepsilon_n}(x,x) = \rho(x) + O(\varepsilon_n), \\
 \nu_n(x,x) = \rho(x) + O(\varepsilon_n), 
 \right\} \quad (x \in S),$$

while

$$\left\{ \nu_n^1(x) P_{\varepsilon_n}(x, y) = \varepsilon_n \rho(x) r(x, y) + O(\varepsilon_n^2), \\ \nu_n(x, y) = \varepsilon_n w(x, y) + O(\varepsilon_n^2), \end{array} \right\} \quad (x, y \in S, \ x \neq y).$$

It follows that

$$\varepsilon_n^{-1} H(\nu_n | \nu_n^1 \otimes P_{\varepsilon_n}) = \varepsilon_n^{-1} \sum_{x,y} \nu_n^1(x) P_{\varepsilon_n}(x,y) \psi\left(\frac{\nu_n(x,y)}{\nu_n^1(x) P_{\varepsilon_n}(x,y)}\right)$$

$$= \varepsilon_n^{-1} \sum_x \left(\rho(x) + O(\varepsilon_n)\right) \psi\left(\frac{\rho(x) + O(\varepsilon_n)}{\rho(x) + O(\varepsilon_n)}\right)$$

$$+ \sum_{x \neq y} \left(\rho(x) r(x,y) + O(\varepsilon_n)\right) \psi\left(\frac{\varepsilon_n w(x,y) + O(\varepsilon_n^2)}{\varepsilon_n \rho(x) r(x,y) + O(\varepsilon_n^2)}\right)$$

$$\ge \sum_{x \neq y} \rho(x) r(x,y) \psi\left(\frac{w(x,y)}{\rho(x) r(x,y)}\right) + O(\varepsilon_n).$$
(3.22)

To prove Claim 2.III, finally, we observe that for each  $(\rho, w) \in \mathcal{W}$ , we can find  $\nu_n \in \mathcal{V}$  satisfying (3.21) such that moreover  $\nu_n(x, x) = 0$  whenever  $\rho(x) = 0$  and  $\nu_n(x, y) = 0$  whenever  $\rho(x)r(x, y) = 0$  for some  $x \neq y$ . It follows that  $\nu_n^1(x) = 0$  whenever  $\rho(x) = 0$ , so for each x, y such that  $\rho(x) = 0$ , we have

$$\varepsilon_n^{-1}\nu_n^1(x)P_{\varepsilon_n}(x,y)\psi\bigg(\frac{\nu_n(x,y)}{\nu_n^1(x)P_{\varepsilon_n}(x,y)}\bigg)=0,$$

while for  $x \neq y$  such that  $\rho(x) > 0$  but r(x, y) = 0, we have

$$\varepsilon_n^{-1}\nu_n^1(x)P_{\varepsilon_n}(x,y)\psi\Big(\frac{\nu_n(x,y)}{\nu_n^1(x)P_{\varepsilon_n}(x,y)}\Big)=O(\varepsilon_n)\psi(1).$$

It follows that in (3.22), only the terms where  $\rho(x)r(x,y) > 0$  contribute, and

$$\varepsilon_n^{-1} H(\nu_n | \nu_n^1 \otimes P_{\varepsilon_n}) = \sum_{x \neq y} \rho(x) r(x, y) \psi\left(\frac{w(x, y)}{\rho(x) r(x, y)}\right) + O(\varepsilon_n).$$

Proof of Claim 3. Set  $\varepsilon \mathbb{N} := \{\varepsilon k : k \in \mathbb{N}\}\$  and observe that  $\varepsilon \lfloor T/\varepsilon \rfloor = \sup\{T' \in \varepsilon \mathbb{N} : T' \leq T\}$ . It is not hard to show that for any  $T_m \to \infty$  and  $\varepsilon_m \downarrow 0$ , the random variables

$$(M_{T_m}, W_{T_m})$$
 and  $(M_{\varepsilon_m \lfloor T_m / \varepsilon_m \rfloor}, W_{\varepsilon_m \lfloor T_m / \varepsilon_m \rfloor})$  (3.23)

are exponentially close. Therefore, by Excercise 3.37 below and the fact that  $(M_{T_m}^{\varepsilon_m}, W_{T_m}^{\varepsilon_m})$  are functions of  $\varepsilon_m \lfloor T_m / \varepsilon_m \rfloor$  only, it suffices to prove the statement for times  $T_m \in \varepsilon_m \mathbb{N}$ .

Recall that  $\Delta_T := \{t \in (0, T] : X_{t-} \neq X_t\}$  is the set of times, up to time T, when X makes a jump. For any  $T \in \varepsilon \mathbb{N}$ , let

$$J_i(\varepsilon, T) := \sum_{k=1}^{T/\varepsilon} \mathbb{1}_{\left\{ \left| \Delta_T \cap (\varepsilon(k-1), \varepsilon k] \right| \ge i \right\}} \qquad (i = 1, 2)$$

denote the number of time intervals of the form  $(\varepsilon(k-1), \varepsilon k]$ , up to time T, during which X makes at least *i* jumps. We observe that for any  $T \in \varepsilon \mathbb{N}$ ,

$$\sum_{x \in S} \left| M_T^{\varepsilon}(x) - M_T(x) \right| \le \frac{\varepsilon}{T} J_1(\varepsilon, T),$$
$$\sum_{x, y \in S} \left| W_T^{\varepsilon}(x, y) - W_T(x, y) \right| \le \frac{1}{T} J_2(\varepsilon, T).$$

Thus, it suffices to show that for any  $\delta > 0$ ,  $\varepsilon_m \downarrow 0$  and  $T_m \in \varepsilon_m \mathbb{N}$  such that  $T_m \to \infty$ 

$$\lim_{m \to \infty} \frac{1}{T_m} \log \mathbb{P} \big[ \varepsilon_m J_1(\varepsilon_m, T_m) / T_m \ge \delta \big] = -\infty,$$
$$\lim_{m \to \infty} \frac{1}{T_m} \log \mathbb{P} \big[ J_2(\varepsilon_m, T_m) / T_m \ge \delta \big] = -\infty.$$

We observe that  $J_1(\varepsilon, T) \leq |\Delta_T|$ , which can in turn be estimated from above by a Poisson distributed random variable  $N_{RT}$  with mean

$$T \sup_{x \in S} \sum_{y \in S} r(x, y) =: RT.$$

By Excercise 3.35 below, it follows that for any  $0 < \varepsilon < \delta/R$ ,

$$\begin{split} \limsup_{m \to \infty} \frac{1}{T_m} \log \mathbb{P} \big[ \varepsilon_m J_1(\varepsilon_m, T_m) / T_m \ge \delta \big] \\ \le \limsup_{m \to \infty} \frac{1}{T_m} \log \mathbb{P} \big[ \varepsilon N_{RT_m} / T_m \ge \delta \big] \le \psi(\delta / R\varepsilon) \xrightarrow[\varepsilon \to 0]{} -\infty, \end{split}$$

where  $\psi(z) := 1 - z + z \log z$ . To also prove the statement for  $J_2$ , we observe that  $\Delta_T$  can be estimated from above by a Poisson point process with intensity R, hence

$$\mathbb{P}[\left|\Delta_T \cap (\varepsilon(k-1), \varepsilon k]\right| \ge 2] \le 1 - e^{-R\varepsilon} - R\varepsilon e^{-R\varepsilon}$$

and  $J_2(\varepsilon, T)$  can be estimated from above by a binomially distributed random variable with parameters  $(n, p) = (T/\varepsilon, 1 - e^{-R\varepsilon} - R\varepsilon e^{-R\varepsilon})$ . For small  $\varepsilon$ , this binomal distribution approximates a Poisson distribution. To turn this into a rigorous estimate, define  $\lambda_{\varepsilon}$  by

$$1 - e^{-\lambda_{\varepsilon}} := 1 - e^{-R\varepsilon} - R\varepsilon e^{-R\varepsilon}$$

In other words, if M and N are Poisson distributed random variable with mean  $\lambda_{\varepsilon}$  and  $R\varepsilon$ , respectively, then this says that  $\mathbb{P}[N \ge 1] = \mathbb{P}[M \ge 2]$ . Since the right-hand side of this equation is of order  $\frac{1}{2}R^2\varepsilon^2 + O(\varepsilon^3)$  as  $\varepsilon \downarrow 0$ , we see that

$$\lambda_{\varepsilon} = \frac{1}{2}R^2\varepsilon^2 + O(\varepsilon^3)$$
 as  $\varepsilon \downarrow 0$ .

Then  $J_2(\varepsilon, T)$  can be estimated from above by a Poisson disributed random variable with mean  $(T/\varepsilon)\lambda_{\varepsilon} = \frac{1}{2}R^2T\varepsilon + O(\varepsilon^2)$ . By the same argument as for  $J_1$ , we conclude that

$$\limsup_{m \to \infty} \frac{1}{T_m} \log \mathbb{P} \big[ \varepsilon_m J_2(\varepsilon_m, T_m) / T_m \ge \delta \big] = -\infty.$$

Exercise 3.35 (Large deviations for Poisson process) Let  $N = (N_t)_{t\geq 0}$  be a Poisson process with intensity one, i.e., N has independent increments where  $N_t - N_t$  is Poisson distributed with mean t - s. Show that the laws  $\mathbb{P}[N_T/T \in \cdot]$ satisfy the large deviation principle with speed T and good rate function

$$I(z) = \begin{cases} 1 - z + z \log z & \text{if } z \ge 0, \\ \infty & \text{otherwise.} \end{cases}$$

Hint: first consider the process at integer times and use that this is a sum of i.i.d. random variables. Then generalize to nontinteger times.

**Exercise 3.36 (Rounded times)** Prove that the random variables in (3.23) are exponentially close.

**Exercise 3.37 (Triangle inequality for exponential closeness)** Let  $(X_n)_{n\geq 1}$ ,  $(Y_n)_{n\geq 1}$  and  $(Z_n)_{n\geq 1}$  be random variables taking values in a Polish space E and let d be a metric generating the topology on E. Let  $s_n$  be positive constants, converging to infinity, and assume that

$$\lim_{n \to \infty} \frac{1}{s_n} \log \mathbb{P}[d(X_n, Y_n) \ge \varepsilon] = -\infty \qquad (\varepsilon > 0),$$
$$\lim_{n \to \infty} \frac{1}{s_n} \log \mathbb{P}[d(Y_n, Z_n) \ge \varepsilon] = -\infty \qquad (\varepsilon > 0).$$

Prove that

$$\lim_{n \to \infty} \frac{1}{s_n} \log \mathbb{P} \big[ d(X_n, Z_n) \ge \varepsilon \big] = -\infty \qquad (\varepsilon > 0).$$

### 3.6 Excercises

Exercise 3.38 (Testing the fairness of a dice) Imagine that we want to test if a dice is fair, i.e., if all sides come up with equal probabilities. To test this hypothesis, we throw the dice n times. General statistical theory tells us that any test on the distribution with which each side comes up can be based on the relative frequencies  $M_n(x)$  of the sides  $x = 1, \ldots, 6$  in these n throws. Let  $\mu_0$  be the uniform distribution on  $S := \{1, \ldots, 6\}$  and imagine that sides the dice come up according to some other, unknown distribution  $\mu_1$ . We are looking for a test function  $T : \mathcal{M}_1(S) \to \{0, 1\}$  such that if  $T(M_n) = 1$ , we reject the hypothesis that the dice is fair. Let  $\mathbb{P}_{\mu}$  denote the distribution of  $M_n$  when in a single throw, the sides of the dice come up with law  $\mu$ . Then

$$\alpha_n := \mathbb{P}_{\mu_0}[T(M_n) = 1] \text{ and } \beta_n := \mathbb{P}_{\mu_1}[T(M_n) = 0]$$

are the probability that we incorrectly reject the hypothesis that the dice is fair and the probability that we do not reckognize the non-fairness of the dice, respectively. A good test minimalizes  $\beta_n$  when  $\alpha_n$  is subject to a bound of the form  $\alpha_n \leq \varepsilon$ , with  $\varepsilon > 0$  small and fixed. Consider a test of the form

$$T(M_n) := 1\{H(M_n|\mu_0) \ge \lambda\},\$$

where  $\lambda > 0$  is fixed and small enough such that  $\{\mu \in \mathcal{M}_1(S) : H(\mu|\mu_0) \ge \lambda\} \neq \emptyset$ . Prove that

$$\lim_{n \to \infty} \frac{1}{n} \log \alpha_n = -\lambda_n$$

and, for any  $\mu_1 \neq \mu_0$ ,

$$\lim_{n \to \infty} \frac{1}{n} \log \beta_n = - \inf_{\mu: H(\mu|\mu_0) < \lambda} H(\mu|\mu_1).$$

#### 3.6. EXCERCISES

Let  $\tilde{T} : \mathcal{M}_1(S) \to \{0, 1\}$  be any other test such that  $\{\mu \in \mathcal{M}_1(S) : \tilde{T}(\mu) = 1\}$  is the closure of its interior and let  $\tilde{\alpha}_n, \tilde{\beta}_n$  be the corresponding error probabilities. Assume that

$$\limsup_{n \to \infty} \frac{1}{n} \log \tilde{\alpha}_n \le -\lambda.$$

Show that for any  $\mu_1 \neq \mu_0$ ,

$$\liminf_{n \to \infty} \frac{1}{n} \log \tilde{\beta}_n \ge -\inf_{\mu: \ H(\mu|\mu_0) < \lambda} H(\mu|\mu_0).$$

This shows that the test T is, in a sense, optimal.

**Exercise 3.39 (Reducible Markov chains)** Let  $X = (X_k)_{k\geq 0}$  be a Markov chain with finite state space S and transition kernel P. Assume that  $S = A \cup B \cup \{c\}$  where

- (i)  $\forall a, a' \in A \exists n \ge 0 \text{ s.t. } P^n(a, a') > 0,$
- (ii)  $\forall b, b' \in B \ \exists n \ge 0 \text{ s.t. } P^n(b, b') > 0,$
- (iii)  $\exists a \in A, b \in B \text{ s.t. } P(a,b) > 0,$
- (iv)  $\exists b \in B \text{ s.t. } P(b,c) > 0$ ,
- (v)  $P(a,c) = 0 \ \forall a \in A$ ,
- (vi)  $P(b,a) = 0 \ \forall a \in A, \ b \in B.$

Assume that  $X_0 \in A$  a.s. Give an expression for

$$\lim_{n \to \infty} \frac{1}{n} \log \mathbb{P}[X_n \neq c].$$

Hint: set  $\tau_B := \inf\{k \ge 0 : X_k \in B\}$  and consider the process before and after  $\tau_B$ .

**Exercise 3.40 (Sampling without replacement)** For each  $n \ge 1$ , consider an urn with n balls that have colors taken from some finite set S. Let  $c_n(x)$  be the number of balls of color  $x \in S$ . Imagine that we draw  $m_n$  balls from the urn without replacement. We assume that the numbers  $c_n(x)$  and  $m_n$  are deterministic (i.e., non-random), and that

$$\frac{1}{n}c_n(x) \xrightarrow[n \to \infty]{} \mu(x) \quad (x \in S) \quad \text{and} \quad \frac{m_n}{n} \xrightarrow[n \to \infty]{} \kappa_n$$

where  $\mu$  is a probability measure on S and  $0 < \kappa < 1$ . Let  $M_n(x)$  be the (random) number of balls of color x that we have drawn. Let  $k_n(x)$  satisfy

$$\frac{k_n(x)}{m_n} \xrightarrow[n \to \infty]{} \nu_1(x) \quad \text{and} \quad \frac{c_n(x) - k_n(x)}{n - m_n} \xrightarrow[n \to \infty]{} \nu_2(x) \qquad (x \in S),$$

where  $\nu_1, \nu_2$  are probability measures on S such that  $\nu_i(x) > 0$  for all  $x \in S$ , i = 1, 2. Prove that

$$\lim_{n \to \infty} \frac{1}{n} \log \mathbb{P}[M_n = k_n] = -\kappa H(\nu_1 | \mu) - (1 - \kappa) H(\nu_2 | \mu).$$
(3.24)

Sketch a proof, similar to the arguments following (3.9), that the laws  $\mathbb{P}[M_n \in \cdot]$  satisfy a large deviation principle with speed n and rate function given by the right-hand side of (3.24). Hint: use Stirling's formula to show that

$$\frac{1}{n}\log\binom{n}{m} \approx H\left(\frac{m}{n}\right),$$

where

$$H(z) := -z \log z - (1-z) \log(1-z).$$

**Exercise 3.41 (Conditioned Markov chain)** Let S be a finite set and let P be a probability kernel on S. Let

$$U \subset \{(x, y) \in S^2 : P(x, y) > 0\}$$

be irreducible, let  $\overline{U}$  be as in (3.3), and let A be the restriction of P to  $\overline{U}$ , i.e., A is the linear operator on  $\mathbb{R}^{\overline{U}}$  whose matrix is given by A(x, y) := P(x, y)  $(x, y \in \overline{U})$ . Let  $\alpha, \eta$  and h denote its Perron-Frobenius eigenvalue and associated left and right eigenvectors, respectively, normalized such that  $\sum_{x\in\overline{U}} h(x)\eta(x) = 1$ , and let  $A_h$  be the irreducible probability kernel on  $\overline{U}$  defined as in (3.15).

Fix  $x_0 \in \overline{U}$ , let  $X = (X_k)_{k\geq 0}$  be the Markov chain in S with transition kernel P started in  $X_0 = x_0$ , and let  $X^h = (X^h_k)_{k\geq 0}$  be the Markov chain in  $\overline{U}$  with transition kernel  $A_h$  started in  $X^h_0 = x_0$ . Show that

$$\mathbb{P}[X_1 = x_1, \dots, X_n = x_n \mid (X_{k-1}, X_k) \in U \; \forall k = 1, \dots, n] \\ \mathbb{E}[h^{-1}(X_n^h)]^{-1} \mathbb{E}[1_{\{X_1^h = x_1, \dots, X_n^h = x_n\}} h^{-1}(X_n^h)],$$

#### 3.6. EXCERCISES

where  $h^{-1}$  denotes the function  $h^{-1}(x) = 1/h(x)$ . Assuming moreover that  $A_h$  is aperiodic, prove that

$$\mathbb{P}[X_1 = x_1, \dots, X_m = x_m \mid (X_{k-1}, X_k) \in U \; \forall k = 1, \dots, n]$$
$$\xrightarrow[n \to \infty]{} \mathbb{P}[X_1^h = x_1, \dots, X_m^h = x_m]$$

for each fixed  $m \ge 1$  and  $x_1, \ldots, x_m \in \overline{U}$ . Hint:

$$\mathbb{P}[X_1 = x_1, \dots, X_m = x_m \mid (X_{k-1}, X_k) \in U \; \forall k = 1, \dots, n] (A_h^n h^{-1})(x_0)^{-1} \mathbb{E}[1_{\{X_1^h = x_1, \dots, X_m^h = x_m\}} (A_h^{n-m} h^{-1})(X_m^h)].$$

## Bibliography

- [Aco02] A. de Acosta. Moderate deviations and associated Laplace transformations for sums of independent random vectors. *Trans. Am. Math. Soc.* 329(1), 357–375, 2002.
- [Bil99] P. Billingsley. Convergence of Probability Measures. 2nd ed. Wiley, New York, 1999.
- [Bou58] N. Bourbaki. Éléments de Mathématique. VIII. Part. 1: Les Structures Fondamentales de l'Analyse. Livre III: Topologie Générale. Chap. 9: Utilisation des Nombres Réels en Topologie Générale. 2iéme éd. Actualités Scientifiques et Industrielles 1045. Hermann & Cie, Paris, 1958.
- [Bry90] W. Bryc. Large deviations by the asymptotic value method. Pages 447–472 in: Diffusion Processes and Related Problems in Analysis Vol. 1 (ed. M. Pinsky), Birkhäuser, Boston, 1990.
- [Cho69] G. Choquet. Lectures on Analysis. Volume I. Integration and Topological Vector Spaces. Benjamin, London, 1969.
- [Cra38] H. Cramér. Sur un nouveau théorème-limite de la théorie des probabilités. Actualités Scientifiques et Industrielles 736, 5-23, 1938.
- [Csi06] I. Csiszár. A simple proof of Sanov's theorem. Bull. Braz. Math. Soc. (N.S.) 37(4), 453–459, 2006.
- [DB81] C.M. Deo and G.J. Babu. Probabilities of moderate deviations in Banach spaces. Proc. Am. Math. Soc. 83(2), 392–397, 1981.
- [DE97] P. Dupuis and R.S. Ellis. A weak convergence approach to the theory of large deviations. Wiley Series in Probability and Statistics. Wiley, Chichester, 1997.

- [DS89] J.-D. Deuschel and D.W. Stroock. *Large deviations*. Academic Press, Boston, 1989.
- [Dud02] R.M. Dudley. *Real Analysis and Probability*. Reprint of the 1989 edition. Camebridge University Press, Camebridge, 2002.
- [DZ93] A. Dembo and O. Zeitouni. Large deviations techniques and applications. Jones and Bartlett Publishers, Boston, 1993.
- [DZ98] A. Dembo and O. Zeitouni. Large deviations techniques and applications 2nd edition. Applications of Mathematics 38. Springer, New York, 1998.
- [EL03] P. Eichelsbacher and M. Löwe. Moderate deviations for i.i.d. random variables. ESAIM, Probab. Stat. 7, 209–218, 2003.
- [Ell85] R.S. Ellis. Entropy, large deviations, and statistical mechanics. Grundlehren der Mathematischen Wissenschaften 271. Springer, New York, 1985.
- [Eng89] R. Engelking. General Topology. Heldermann, Berlin, 1989.
- [EK86] S.N. Ethier and T.G. Kurtz. Markov Processes; Characterization and Convergence. John Wiley & Sons, New York, 1986.
- [Gan00] F.R. Gantmacher. *The Theory of Matrices, Vol. 2.* AMS, Providence RI, 2000.
- [Hol00] F. den Hollander. *Large Deviations*. Fields Institute Monographs 14. AMS, Providence, 2000.
- [Kel75] J.L. Kelley. General Topology. Reprint of the 1955 edition printed by Van Nostrand. Springer, New York, 1975.
- [Led92] M. Ledoux. Sur les déviations modérés des sommes de variables aléatoires vectorielles indépendantes de même loi. Ann. Inst. Henri Poincaré, Probab. Stat., 28(2), 267–280, 1992.
- [OV91] G.L. O'Brien and W. Verwaat. Capacities, large deviations and loglog laws. Page 43–83 in: *Stable Processes and Related Topics* Progress in Probability 25, Birkhäuser, Boston, 1991.
- [Oxt80] J.C. Oxtoby. Measure and Category. Second Edition. Springer, New York, 1980.

- [Puk91] A.A. Pukhalski. On functional principle of large deviations. Pages 198– 218 in: New Trends in Probability and Statistics (eds. V. Sazonov and T. Shervashidze) VSP-Mokslas, 1991.
- [Puh01] A. Puhalskii. Large Deviations and Idempotent Probability. Monographs and Surveys in Pure and Applied Mathematics 119. Chapman & Hall, Boca Raton, 2001.
- [Roc70] R.T. Rockafellar. Convex Analysis. Princeton, New Jersey, 1970.
- [San61] I.N. Sanov. On the probability of large deviations of random variables. Mat. Sb. 42 (in russian). English translation in: Selected Translations in Mathematical Statistics and Probability I, 213–244, 1961.
- [Sen73] E. Seneta. Non-Negative Matrices: An Introduction to Theory and Applications. George Allen & Unwin, London, 1973.
- [Ste87] J. Štěpán. Teorie Pravěpodobnosti. Academia, Prague, 1987.
- [Var66] S.R.S. Varadhan. Asymptotic probabilities and differential equations. Comm. Pure Appl. Math. 19, 261–286, 1966.

# Index

 $A^{\rm c}, 42$ B(E), 19 $B_{+}(E), 19$  $B_b(E), 18, 19$  $B_r(x), 18, 42$  $B_{b,+}(E), 19$  $G_{\delta}$ -set, 44 *I*-continuous set, 24 int(A), 24 $\overline{A}$ , 24  $\overline{\mathbb{R}}, 19$  $\lor$ , 20  $\wedge, 20$  $\mathcal{C}_b(E), 18$  $\mathcal{B}(E), 18$  $\mathcal{C}(E), 19$  $C_{+}(E), 19$  $\mathcal{C}_b(E), 19$  $C_{b,+}(E), 19$  $\mathcal{L}(E), 19$  $\mathcal{L}_{+}(E), 19$  $\mathcal{L}_b(E), 19$  $\mathcal{L}_{b,+}(E), 19$  $\mathcal{U}(E), 19$  $\mathcal{U}_{+}(E), 19$  $\mathcal{U}_b(E), 19$  $\mathcal{U}_{b,+}(E), 19$  $\overline{A}$ , 18 int(A), 17aperiodic Markov chain, 92 boundary of a set, 67

central limit theorem, 10 closed convex hull, 70 closure, 18 compact level sets, 23 compactification, 44 contraction principle, 30 convex hull, 60 convex function, 56, 66 convex hull, 70 cumulant generating function, 8 dense set, 18diagonal argument, 96 distribution determining, 36 dual space, 74 eigenvector left or right, 112 empirical average, 7 empirical distribution finite space, 13 for pairs, 92 of Markov process, 118 empirical process, 110 epigraph, 56 exponential tightness, 42 exponentially close, 33 Fenchel-Legendre transform, 55 free energy function, 8 good rate function, 23

Hausdorff topological space, 17 image measure, 30 induced topology, 44, 75 initial law, 91 interior, 17 invariant law, 92 of Markov process, 117 inverse image, 30 irreducibility, 15, 92, 111 irreducible Markov chain, 92 Markov process, 15 set U, 95kernel probability, 91 Kullback-Leibler distance, 12 large deviation principle, 23 weak, 51 law of large numbers weak, 7 LDP, 23 Legendre transform, 55, 66 Legendre-Fenchel transform, 55 level set, 9 compact, 23logarithmic cumulant generating function, 8 logarithmic moment generating function, 8, 60 lower semi-continuous, 9 Markov chain, 91 moderate deviations, 11 moment generating function, 8

nonnegative matrix, 111 norm, 23 normalized rate function, 32 one-point compactification, 46 operator norm, 111 partial sum, 10 period of a state, 92 Perron-Frobenius theorem, 111 probability kernel, 91 projective limit, 52 rate, 23 rate function normalized, 32 rate function, 23 Cramér's theorem, 8 good, 23 rate function determining, 47, 49 relative compactness, 74 relative entropy, 72 finite space, 12 restriction principle, 45 Scott topology, 19 seminorm, 23 separable, 18 separation of points, 52

simple function, 20 spectral radius, 111 spectrum, 111 speed, 23 stationary process, 92 Stirling's formula, 97 supporting line, 58

tightness, 36 exponential, 42 tilted probability law, 32, 60, 69 total variation, 93 totally bounded, 43 transition kernel, 91 transition rate, 15

INDEX

transposed matrix, 112 trap, 104

uniform integrability, 74

138