

**The Theory of Large Deviations  
and Applications to Statistical Mechanics**

**Lectures for the International Seminar  
on Extreme Events in Complex Dynamics  
October 23–27, 2006**

**Max-Planck-Institut für  
Physik komplexer Systeme  
Dresden, Germany**

Richard S. Ellis  
Department of Mathematics and Statistics  
University of Massachusetts  
Amherst, MA 01003

rsellis@math.umass.edu  
<http://www.math.umass.edu/~rsellis>

Copyright © 2006 Richard S. Ellis

## Our Lives Are Large Deviations

Statistically, the probability of any one of us being here is so small that you'd think the mere fact of existing would keep us all in a contented dazzlement of surprise. We are alive against the stupendous odds of genetics, infinitely outnumbered by all the alternates who might, except for luck, be in our places.

Even more astounding is our statistical improbability in physical terms. The normal, predictable state of matter throughout the universe is randomness, a relaxed sort of equilibrium, with atoms and their particles scattered around in an amorphous muddle. We, in brilliant contrast, are completely organized structures, squirming with information at every covalent bond. We make our living by catching electrons at the moment of their excitement by solar photons, swiping the energy released at the instant of each jump and storing it up in intricate loops for ourselves. We violate probability, by our nature. To be able to do this systemically, and in such wild varieties of form, from viruses to whales, is extremely unlikely; to have sustained the effort successfully for the several billion years of our existence, without drifting back into randomness, was nearly a mathematical impossibility.

Lewis Thomas, *The Lives of a Cell*

## **Contents**

<b>1</b>	<b>Introduction</b>	<b>4</b>
<b>2</b>	<b>A Basic Probabilistic Model</b>	<b>9</b>
<b>3</b>	<b>Boltzmann's Discovery and Relative Entropy</b>	<b>11</b>
<b>4</b>	<b>The Most Likely Way for an Unlikely Event To Happen</b>	<b>19</b>
<b>5</b>	<b>Gibbs States for Models in Statistical Mechanics</b>	<b>26</b>
<b>6</b>	<b>Generalities: Large Deviation Principle and Laplace Principle</b>	<b>32</b>
<b>7</b>	<b>Cramér's Theorem</b>	<b>43</b>
<b>8</b>	<b>Gärtner-Ellis Theorem</b>	<b>58</b>
<b>9</b>	<b>The Curie-Weiss Model of Ferromagnetism</b>	<b>66</b>
<b>10</b>	<b>Equivalence of Ensembles for a General Class of Models in Statistical Mechanics</b>	<b>72</b>
<b>11</b>	<b>Maximum Entropy Principles in Two-Dimensional Turbulence</b>	<b>96</b>
	<b>References</b>	<b>108</b>

# 1 Introduction

The theory of large deviations studies the exponential decay of probabilities in certain random systems. It has been applied to a wide range of problems in which detailed information on rare events is required. One is often interested not only in the probability of rare events but also in the characteristic behavior of the system as the rare event occurs. For example, in applications to queueing theory and communication systems, the rare event could represent an overload or breakdown of the system. In this case, large deviation methodology can lead to an efficient redesign of the system so that the overload or breakdown does not occur. In applications to statistical mechanics the theory of large deviations gives precise, exponential-order estimates that are perfectly suited for asymptotic analysis.

These lectures will present a number of topics in the theory of large deviations and several applications to statistical mechanics, all united by the concept of relative entropy. This concept entered human culture through the first large deviation calculation in science, carried out by Ludwig Boltzmann. Stated in a modern terminology, his discovery was that the relative entropy expresses the asymptotic behavior of certain multinomial probabilities. This statistical interpretation of entropy has the following crucial physical implication [33, §1.1].

Entropy is a bridge between a microscopic level, on which physical systems are defined in terms of the complicated interactions among the individual constituent particles, and a macroscopic level, on which the laws describing the behavior of the system are formulated.

Building on the work of Boltzmann, Gibbs asked a fundamental question. How can one use probability theory to study equilibrium properties of physical systems such as an ideal gas, a ferromagnet, or a fluid? These properties include such phenomena as phase transitions; e.g., the liquid-gas transition or spontaneous magnetization in a ferromagnet. Another example arises in the study of freely evolving, inviscid fluids, for which one wants to describe coherent states. These are steady, stable mean flows comprised of one or more vortices that persist amidst the turbulent fluctuations of the vorticity field. Gibbs's answer,

which led to the development of classical equilibrium statistical mechanics, is that one studies equilibrium properties via probability measures on configuration space known today as Gibbs canonical ensembles or Gibbs states. For background in statistical mechanics, I recommend [33, 55, 79], which cover a number of topics relevant to these lectures.

One of my main purposes is to show the utility of the theory of large deviations by applying it to a number of statistical mechanical models. Our applications of the theory include the following.

- A derivation of the form of the Gibbs state for a discrete ideal gas (section 5).
- A probabilistic description of the phase transition in the Curie-Weiss model of a ferromagnet in terms of the breakdown of the law of large numbers for the spin per site (section 9).
- An analysis of equivalence and nonequivalence of ensembles for a general class of models, including spin models and models of coherent structures in turbulence (section 10).
- A derivation of variational formulas that describe the equilibrium macrostates in models of two-dimensional turbulence (section 11). In terms of these macrostates, coherent vortices of two-dimensional turbulence can be studied.

Like many areas of mathematics, the theory of large deviations has both a left hand and a right hand; the left hand provides heuristic insight while the right hand provides formal proofs. Although the theory is applicable in many diverse settings, the right-hand technicalities can be formidable. Recognizing this, I would like to supplement the rigorous, right-hand formulation of the theory with a number of basic results presented in a left-hand format useful to the applied researcher.

Boltzmann's calculation of the asymptotic behavior of multinomial probabilities in terms of relative entropy was carried out in 1877 as a key component of his paper that gave a probabilistic interpretation of the Second Law of Thermodynamics [4]. This momentous calculation represents a revolutionary moment

in human culture during which both statistical mechanics and the theory of large deviations were born. Boltzmann based his work on the hypothesis that atoms exist. Although this hypothesis is universally accepted today, one might be surprised to learn that it was highly controversial during Boltzmann's time [57, pp. vii–x].

Boltzmann's work is put in historical context by W. R. Everdell in his book *The First Moderns*, which traces the development of the modern consciousness in nineteenth and twentieth century thought [46]. Chapter 3 focuses on the mathematicians of Germany in the 1870's — namely, Cantor, Dedekind, and Frege — who “would become the first creative thinkers in any field to look at the world in a fully twentieth-century manner” [p. 31]. Boltzmann is then presented as the man whose investigations in stochastics and statistics made possible the work of the two other great founders of twentieth-century theoretical physics, Planck and Einstein. As Everdell writes, “he was at the center of the change” [p. 48].

Although the topic of these lectures is the theory of large deviations and not the history of science, it is important to appreciate the radical nature of Boltzmann's ideas. His belief in the existence of atoms and his use of probabilistic laws at the microscopic level of atoms and molecules to derive macroscopic properties of matter profoundly challenged the conventional wisdom of 19<sup>th</sup> century physics: physical laws express absolute truths based not on probabilistic assumptions, but on Newton's laws of motion and precise measurements of observable phenomena.

For his subversive attack on the temple of conventional wisdom, Boltzmann would eventually pay the ultimate price [8, p. 34].

Boltzmann had never worried about his health, but had sacrificed it to his scientific activity. When however even that vacation in Duino did not bring any relief from his illness, in a moment of deep depression he committed suicide by hanging on 5 September 1906. The next day he should have gone to Vienna to start his lectures.

The irony is that in 1905, the year before Boltzmann's suicide, Einstein applied Boltzmann's insights with great success [57, ch. 11]. In one paper he used

Boltzmann's idea to partition the energy of a gas into discrete units in order to explain a phenomenon known as the photoelectric effect. This work would mark the beginning of quantum mechanics and would eventually win him the Nobel Prize. In two other papers also written in 1905 Einstein gave a statistical mechanical explanation based directly on Boltzmann's theory to explain the random motion of a particle suspended in a fluid, a phenomenon known as Brownian motion. This work strongly corroborated the existence of atoms, putting statistical mechanics on a firm theoretical basis. From these papers and two additional 1905 papers on special relativity, the second of which contains the famous formula  $E = mc^2$ , modern physics was born.

Boltzmann's insights are now part of the canon, but he paid for this with his life. Without his insights, modern physics might never have been born, and unborn, it would not have become our civilization's main conceptual lens for interpreting the universe and our place in it.

Here is an overview of the contents of each section of these lectures.

- **Section 2.** A basic probabilistic model is introduced.
- **Section 3.** Boltzmann's discovery of the asymptotic behavior of multinomial probabilities in terms of relative entropy is described.
- **Section 4.** The probabilities of a loaded die are calculated as an illustration of a general principle expressed in the following question. What is the most likely way for an unlikely event to happen?
- **Section 5.** The probabilities of the energy states of a discrete ideal gas are calculated, generalizing the calculation in section 3.

The solutions of the problems in sections 4 and 5 motivate the form of the Gibbs canonical ensemble. This is a probability distribution used to determine the equilibrium properties of statistical mechanical systems; it is discussed in section 9 for a specific model and in section 10 for a general class of models.

- **Section 6.** We introduce the general concepts of a large deviation principle and a Laplace principle, together with related results.

- **Section 7.** We prove Cramér's Theorem, which is the large deviation principle for the sample means of i.i.d. random variables.
- **Section 8.** The generalization of Cramér's Theorem known as the Gärtner-Ellis Theorem is presented.

In the remainder of the sections the theory of large deviations is applied to a number of questions in statistical mechanics.

- **Section 9.** The theory of large deviations is used to study equilibrium properties of a basic model of ferromagnetism known as the Curie-Weiss model, which is a mean-field approximation to the much more complicated Ising model.
- **Section 10.** Our work in the preceding section leads to the formulation of a general procedure for applying the theory of large deviations to the analysis of an extensive class of statistical mechanical models, an analysis that will allow us to address the fundamental problem of equivalence and nonequivalence of ensembles.
- **Section 11.** The general procedure developed in the preceding section is used along with Sanov's Theorem to derive variational formulas that describe the equilibrium macrostates in two models of coherent states in two-dimensional turbulence; namely, the Miller-Robert theory and a modification of that theory proposed by Turkington.

Sanov's Theorem, which is used in section 11 to analyze two models of coherent states in two-dimensional turbulence, generalizes Boltzmann's 1877 calculation. Because this theorem plays a vital role in the derivation, this final application of the theory of large deviations brings our focus back home to Boltzmann, through whose research in the foundations of statistical mechanics the theory began to blossom.

**Acknowledgement.** The research of Richard S. Ellis is supported by a grant from the National Science Foundation (NSF-DMS-0604071).



## 2 A Basic Probabilistic Model

In later sections we will investigate a number of questions in the theory of large deviations in the context of a basic probabilistic model, which we now introduce. Let  $\alpha \geq 2$  be an integer,  $y_1 < y_2 < \dots < y_\alpha$  a set of  $\alpha$  real numbers, and  $\rho_1, \rho_2, \dots, \rho_\alpha$  a set of  $\alpha$  positive real numbers summing to 1. We think of  $\Lambda = \{y_1, y_2, \dots, y_\alpha\}$  as the set of possible outcomes of a random experiment in which each individual outcome  $y_k$  has the probability  $\rho_k$  of occurring. The vector  $\rho = (\rho_1, \rho_2, \dots, \rho_\alpha)$  is an element of the set of probability vectors

$$\mathcal{P}_\alpha = \left\{ \gamma = (\gamma_1, \gamma_2, \dots, \gamma_\alpha) \in \mathbb{R}^\alpha : \gamma_k \geq 0, \sum_{k=1}^{\alpha} \gamma_k = 1 \right\}.$$

Any vector  $\gamma \in \mathcal{P}_\alpha$  also defines a probability measure on the set of subsets of  $\Lambda$  via the formula

$$\gamma = \gamma(dy) = \sum_{k=1}^{\alpha} \gamma_k \delta_{y_k}(dy),$$

where for  $y \in \Lambda$ ,  $\delta_{y_k}\{y\} = 1$  if  $y = y_k$  and equals 0 otherwise. Thus for  $B \subset \Lambda$ ,  $\gamma\{B\} = \sum_{y_k \in B} \gamma_k$ .

For each positive integer  $n$ , the configuration space for  $n$  independent repetitions of the experiment is  $\Omega_n = \Lambda^n$ , a typical element of which is denoted by  $\omega = (\omega_1, \omega_2, \dots, \omega_n)$ . For each  $\omega \in \Omega_n$  we define

$$P_n\{\omega\} = \prod_{j=1}^n \rho\{\omega_j\}$$

and extend this to a probability measure on the set of subsets of  $\Omega_n$  by defining

$$P_n\{B\} = \sum_{\omega \in B} P_n\{\omega\} \text{ for } B \subset \Omega_n.$$

$P_n$  is called the product measure with one dimensional marginals  $\rho$ . With respect to  $P_n$  the coordinate functions  $X_j(\omega) = \omega_j, j = 1, 2, \dots, n$ , are independent, identically distributed (i.i.d.) random variables with common distribution

$\rho$ ; that is, for any subsets  $B_1, B_2, \dots, B_n$  of  $\Lambda$

$$\begin{aligned} P_n\{\omega \in \Omega_n : X_j(\omega) \in B_j \text{ for } j = 1, 2, \dots, n\} \\ = \prod_{j=1}^n P_n\{\omega \in \Omega_n : X_j(\omega) \in B_j\} = \prod_{j=1}^n \rho\{B_j\}. \end{aligned}$$

**Example 2.1.** Random phenomena that can be studied via this basic model include standard examples such as coin tossing and die tossing and also include a discrete ideal gas.

(a) *Coin tossing.* In this case  $\Lambda = \{1, 2\}$  and  $\rho_1 = \rho_2 = 1/2$ .

(b) *Die tossing.* In this case  $\Lambda = \{1, 2, \dots, 6\}$  and each  $\rho_k = 1/6$ .

(c) *Discrete ideal gas.* Consider a discrete ideal gas consisting of  $n$  identical, noninteracting particles, each having  $\alpha$  equally likely energy levels  $y_1, y_2, \dots, y_\alpha$ ; in this case each  $\rho_k$  equals  $1/\alpha$ . The coordinate functions  $X_j$  represent the random energy levels of the molecules of the gas. The statistical independence of these random variables reflects the fact that the molecules of the gas do not interact. ■

We will return to the discrete ideal gas in section 5 after introducing some basic concepts in theory of large deviations.

### 3 Boltzmann's Discovery and Relative Entropy

In its original form Boltzmann's discovery concerns the asymptotic behavior of certain multinomial coefficients. For the purpose of applications in these lectures, it is advantageous to formulate it in terms of a probabilistic quantity known as the empirical vector. We use the notation of the preceding section. Thus let  $\alpha \geq 2$  be an integer,  $y_1 < y_2 < \dots < y_\alpha$  a set of  $\alpha$  real numbers,  $\rho_1, \rho_2, \dots, \rho_\alpha$  a set of  $\alpha$  positive real numbers summing to 1,  $\Lambda$  the set  $\{y_1, y_2, \dots, y_\alpha\}$ , and  $P_n$  the product measure on  $\Omega_n = \Lambda^n$  with one dimensional marginals  $\rho = \sum_{k=1}^{\alpha} \rho_k \delta_{y_k}$ . For  $\omega = (\omega_1, \omega_2, \dots, \omega_n) \in \Omega_n$ , we let  $\{X_j, j = 1, \dots, n\}$  be the coordinate functions defined by  $X_j(\omega) = \omega_j$ . The  $X_j$  form a sequence of i.i.d. random variables with common distribution  $\rho$ .

We now turn to the object under study in the present section. For  $\omega \in \Omega_n$  and  $y \in \Lambda$  define

$$L_n(y) = L_n(\omega, y) = \frac{1}{n} \sum_{j=1}^n \delta_{X_j(\omega)}\{y\}.$$

Thus  $L_n(\omega, y)$  counts the relative frequency with which  $y$  appears in the configuration  $\omega$ ; in symbols,  $L_n(\omega, y) = n^{-1} \cdot \#\{j \in \{1, \dots, n\} : \omega_j = y\}$ . We then define the empirical vector

$$\begin{aligned} L_n &= L_n(\omega) = (L_n(\omega, y_1), \dots, L_n(\omega, y_\alpha)) \\ &= \frac{1}{n} \sum_{j=1}^n (\delta_{X_j(\omega)}\{y_1\}, \dots, \delta_{X_j(\omega)}\{y_\alpha\}). \end{aligned}$$

$L_n$  equals the sample mean of the i.i.d. random vectors  $(\delta_{X_j(\omega)}\{y_1\}, \dots, \delta_{X_j(\omega)}\{y_\alpha\})$ . It takes values in the set of probability vectors

$$\mathcal{P}_\alpha = \left\{ \gamma = (\gamma_1, \gamma_2, \dots, \gamma_\alpha) \in \mathbb{R}^\alpha : \gamma_k \geq 0, \sum_{k=1}^{\alpha} \gamma_k = 1 \right\}.$$

The limiting behavior of  $L_n$  is straightforward to determine. Let  $\|\cdot\|$  denote the Euclidean norm on  $\mathbb{R}^\alpha$ . For any  $\gamma \in \mathcal{P}_\alpha$  and  $\varepsilon > 0$ , we define the open ball

$$B(\gamma, \varepsilon) = \{\nu \in \mathcal{P}_\alpha : \|\gamma - \nu\| < \varepsilon\}.$$

Since the  $X_j$  have the common distribution  $\rho$ , for each  $y_k \in \Lambda$

$$E^{P_n}\{L_n(y_k)\} = E^{P_n}\left\{\frac{1}{n}\sum_{j=1}^n \delta_{X_j}\{y_k\}\right\} = \frac{1}{n}\sum_{j=1}^n P_n\{X_j = y_k\} = \rho_k,$$

where  $E^{P_n}$  denotes expectation with respect to  $P_n$ . Hence by the weak law of large numbers for the sample means of i.i.d. random variables, for any  $\varepsilon > 0$

$$\lim_{n \rightarrow \infty} P_n\{L_n \in B(\rho, \varepsilon)\} = 1. \quad (3.1)$$

It follows that for any  $\gamma \in \mathcal{P}_\alpha$  not equal to  $\rho$  and for any  $\varepsilon > 0$  satisfying  $0 < \varepsilon < \|\rho - \gamma\|$

$$\lim_{n \rightarrow \infty} P_n\{L_n \in B(\gamma, \varepsilon)\} = 0. \quad (3.2)$$

As we will see, Boltzmann's discovery implies that these probabilities converge to 0 exponentially fast in  $n$ . The exponential decay rate is given in terms of the relative entropy, which we now define.

**Definition 3.1 (Relative Entropy).** Let  $\rho = (\rho_1, \dots, \rho_\alpha)$  denote the probability vector in  $\mathcal{P}_\alpha$  in terms of which the basic probabilistic model is defined. The relative entropy of  $\gamma \in \mathcal{P}_\alpha$  with respect to  $\rho$  is defined by

$$I_\rho(\gamma) = \sum_{k=1}^{\alpha} \gamma_k \log \frac{\gamma_k}{\rho_k}.$$

Several properties of the relative entropy are given in the next lemma.

**Lemma 3.2.** For  $\gamma \in \mathcal{P}_\alpha$ ,  $I_\rho(\gamma)$  measures the discrepancy between  $\gamma$  and  $\rho$  in the sense that  $I_\rho(\gamma) \geq 0$  and  $I_\rho(\gamma) = 0$  if and only if  $\gamma = \rho$ . Thus  $I_\rho(\gamma)$  attains its infimum of 0 over  $\mathcal{P}_\alpha$  at the unique measure  $\gamma = \rho$ . In addition,  $I_\rho$  is strictly convex on  $\mathcal{P}_\alpha$ .

**Proof.** For  $x \geq 0$  the graph of the strictly convex function  $x \log x$  has the tangent line  $y = x - 1$  at  $x = 1$ . Hence  $x \log x \geq x - 1$  with equality if and only if  $x = 1$ . It follows that for any  $\gamma \in \mathcal{P}_\alpha$

$$\frac{\gamma_k}{\rho_k} \log \frac{\gamma_k}{\rho_k} \geq \frac{\gamma_k}{\rho_k} - 1 \quad (3.3)$$

with equality if and only if  $\gamma_k = \rho_k$ . Multiplying this inequality by  $\rho_k$  and summing over  $k$  yields

$$I_\rho(\gamma) = \sum_{k=1}^{\alpha} \gamma_k \log \frac{\gamma_k}{\rho_k} \geq \sum_{k=1}^{\alpha} (\gamma_k - \rho_k) = 0.$$

$I_\rho(\gamma) = 0$  if and only if equality holds in (3.3) for each  $k$ ; i.e., if and only if  $\gamma = \rho$ . This yields the first assertion in the proposition. This proof is typical of proofs of analogous results involving relative entropy [e.g., Prop. 4.2] in that we use a global convexity inequality — in this case,  $x \log x \geq x - 1$  with equality if and only if  $x = 1$  — rather than calculus to determine where  $I_\rho$  attains its infimum over  $\mathcal{P}_\alpha$ . Since

$$I_\rho(\gamma) = \sum_{k=1}^{\alpha} \rho_k \frac{\gamma_k}{\rho_k} \log \frac{\gamma_k}{\rho_k},$$

the strict convexity of  $I_\rho$  is a consequence of the strict convexity of  $x \log x$  for  $x \geq 0$ . ■

We are now ready to give the first formulation of Boltzmann’s discovery, which we state using a heuristic notation and which we label, in recognition of its formal status, as a “theorem.” However, the formal calculations used to motivate the “theorem” can easily be turned into a rigorous proof of an asymptotic theorem. That theorem is stated in Theorem 3.4. From Boltzmann’s momentous discovery both the theory of large deviations and the Gibbsian formulation of equilibrium statistical mechanics grew.

**“Theorem” 3.3 (Boltzmann’s Discovery–Formulation 1).** *For any  $\gamma \in \mathcal{P}_\alpha$  and all sufficiently small  $\varepsilon > 0$*

$$P_n\{L_n \in B(\gamma, \varepsilon)\} \approx \exp[-nI_\rho(\gamma)] \text{ as } n \rightarrow \infty.$$

**Heuristic Proof.** By elementary combinatorics

$$\begin{aligned}
P_n\{L_n \in B(\gamma, \varepsilon)\} &= P_n\left\{\omega \in \Omega_n : L_n(\omega) \sim \frac{1}{n}(n\gamma_1, n\gamma_2, \dots, n\gamma_\alpha)\right\} \\
&\approx P_n\{\#\{\omega_j\text{'s} = y_1\} \sim n\gamma_1, \dots, \#\{\omega_j\text{'s} = y_\alpha\} \sim n\gamma_\alpha\} \\
&\approx \frac{n!}{(n\gamma_1)!(n\gamma_2)! \cdots (n\gamma_\alpha)!} \rho_1^{n\gamma_1} \rho_2^{n\gamma_2} \cdots \rho_\alpha^{n\gamma_\alpha}.
\end{aligned}$$

Stirling's formula in the weak form  $\log(n!) = n \log n - n + \mathbf{O}(\log n)$  yields

$$\begin{aligned}
&\frac{1}{n} \log P_n\{L_n \in B(\gamma, \varepsilon)\} \\
&\approx \frac{1}{n} \log \left( \frac{n!}{(n\gamma_1)!(n\gamma_2)! \cdots (n\gamma_\alpha)!} \right) + \sum_{k=1}^{\alpha} \gamma_k \log \rho_k \\
&= - \sum_{k=1}^{\alpha} \gamma_k \log \gamma_k + \mathbf{O}\left(\frac{\log n}{n}\right) + \sum_{k=1}^{\alpha} \gamma_k \log \rho_k \\
&= - \sum_{k=1}^{\alpha} \gamma_k \log \frac{\gamma_k}{\rho_k} + \mathbf{O}\left(\frac{\log n}{n}\right) = -I_\rho(\gamma) + \mathbf{O}\left(\frac{\log n}{n}\right). \quad \blacksquare
\end{aligned}$$

“Theorem” 3.3 has the following interesting consequence. Let  $\gamma$  be any vector in  $\mathcal{P}_\alpha$  which differs from  $\rho$ . Since  $I_\rho(\gamma) > 0$  [Lemma 3.2], it follows that for all sufficiently small  $\varepsilon > 0$

$$P_n\{L_n \in B(\gamma, \varepsilon)\} \approx \exp[-nI_\rho(\gamma)] \rightarrow 0 \text{ as } n \rightarrow \infty,$$

a limit which, if rigorous, would imply (3.2).

Let  $A$  be a Borel subset of  $\mathcal{P}_\alpha$ ; the class of Borel subsets includes all closed sets and all open sets. If  $\rho$  is not contained in the closure of  $A$ , then by the weak law of large numbers

$$\lim_{n \rightarrow \infty} P_n\{L_n \in A\} = 0,$$

and by analogy with the heuristic asymptotic result given in “Theorem” 3.3 we expect that these probabilities converge to 0 exponentially fast with  $n$ . This is in fact the case. In order to express the exponential decay rate of such probabilities in terms of the relative entropy, we introduce the notation  $I_\rho(A) = \inf_{\gamma \in A} I_\rho(\gamma)$ .

The range of  $L_n(\omega)$  for  $\omega \in \Omega_n$  is the set of probability vectors having the form  $k/n$ , where  $k \in \mathbb{R}^\alpha$  has nonnegative integer coordinates summing to  $n$ ; hence the cardinality of the range does not exceed  $n^\alpha$ . Since

$$P_n\{L_n \in A\} = \sum_{\gamma \in A} P_n\{L_n \sim \gamma\} \approx \sum_{\gamma \in A} \exp[-nI_\rho(\gamma)]$$

and

$$\exp[-nI_\rho(A)] \leq \sum_{\gamma \in A} \exp[-nI_\rho(\gamma)] \leq n^\alpha \exp[-nI_\rho(A)],$$

one expects that at least to exponential order

$$P_n\{L_n \in A\} \approx \exp[-nI_\rho(A)] \text{ as } n \rightarrow \infty. \quad (3.4)$$

As formulated in Corollary 3.5, this asymptotic result is indeed valid. It is a consequence of the following rigorous reformulation of Boltzmann's discovery, known as Sanov's Theorem, which expresses the large deviation principle for the empirical vectors  $L_n$ . That concept is defined in general in Definition 6.1, and a general form of Sanov's Theorem is stated in Theorem 6.7.

**Theorem 3.4 (Boltzmann's Discovery–Formulation 2).** *The sequence of empirical vectors  $L_n$  satisfies the large deviation principle on  $\mathcal{P}_\alpha$  with rate function  $I_\rho$  in the following sense.*

(a) **Large deviation upper bound.** *For any closed subset  $F$  of  $\mathcal{P}_\alpha$*

$$\limsup_{n \rightarrow \infty} \frac{1}{n} \log P_n\{L_n \in F\} \leq -I_\rho(F).$$

(b) **Large deviation lower bound.** *For any open subset  $G$  of  $\mathcal{P}_\alpha$*

$$\liminf_{n \rightarrow \infty} \frac{1}{n} \log P_n\{L_n \in G\} \geq -I_\rho(G).$$

**Comments on the Proof.** For  $\gamma \in \mathcal{P}_\alpha$  and  $\varepsilon > 0$   $B(\gamma, \varepsilon)$  denotes the open ball with center  $\gamma$  and radius  $\varepsilon$  and  $\overline{B}(\gamma, \varepsilon)$  denotes the corresponding closed ball. Since  $\mathcal{P}_\alpha$  is a compact subset of  $\mathbb{R}^\alpha$ , any closed subset  $F$  of  $\mathcal{P}_\alpha$  is automatically compact. By a standard covering argument it is not hard to show that the large

deviation upper bound holds for any closed set  $F$  provided that one obtains the large deviation upper bound for any closed ball  $\overline{B}(\gamma, \varepsilon)$ :

$$\limsup_{n \rightarrow \infty} \frac{1}{n} \log P_n \{L_n \in \overline{B}(\gamma, \varepsilon)\} \leq -I_\rho(\overline{B}(\gamma, \varepsilon)).$$

Likewise, the large deviation lower bound holds for any open set  $G$  provided one obtains the large deviation lower bound for any open ball  $B(\gamma, \varepsilon)$ :

$$\liminf_{n \rightarrow \infty} \frac{1}{n} \log P_n \{L_n \in B(\gamma, \varepsilon)\} \geq -I_\rho(B(\gamma, \varepsilon)).$$

The bounds in the last two displays can be proved via combinatorics and Stirling's formula as in the heuristic proof of "Theorem" 3.3; one can easily adapt the calculations given in [33, §I.4]. The details are omitted. ■

Given  $A$  a Borel subset of  $\mathcal{P}_\alpha$ , we denote by  $A^\circ$  the interior of  $A$  relative to  $\mathcal{P}_\alpha$  and by  $\overline{A}$  the closure of  $A$ . For a class of Borel subsets  $A$  we can now derive a rigorous version of the asymptotic formula (3.4). This class consists of sets  $A$  such that  $\overline{A^\circ}$  equals  $\overline{A}$ . Any open ball  $B(\gamma, \varepsilon)$  or closed ball  $\overline{B}(\gamma, \varepsilon)$  satisfies this condition.

**Corollary 3.5.** *Let  $A$  be any Borel subset of  $\mathcal{P}_\alpha$  such that  $\overline{A^\circ} = \overline{A}$ . Then*

$$\lim_{n \rightarrow \infty} \frac{1}{n} \log P_n \{L_n \in A\} = -I_\rho(A).$$

**Proof.** We apply the large deviation upper bound to  $\overline{A}$  and the large deviation lower bound to  $A^\circ$ . Since  $\overline{A} \supset A \supset A^\circ$ , it follows that

$$\begin{aligned} -I_\rho(\overline{A}) &\geq \limsup_{n \rightarrow \infty} \frac{1}{n} \log P_n \{L_n \in \overline{A}\} \\ &\geq \limsup_{n \rightarrow \infty} \frac{1}{n} \log P_n \{L_n \in A\} \\ &\geq \liminf_{n \rightarrow \infty} \frac{1}{n} \log P_n \{L_n \in A\} \\ &\geq \liminf_{n \rightarrow \infty} \frac{1}{n} \log P_n \{L_n \in A^\circ\} \\ &\geq -I_\rho(A^\circ). \end{aligned}$$



The continuity of  $I_\rho$  on  $\mathcal{P}_\alpha$  implies that  $I_\rho(A^\circ) = I_\rho(\overline{A^\circ})$ . Since by hypothesis  $\overline{A^\circ} = \overline{A}$ , we conclude that the extreme terms in this display are equal. The desired limit follows. ■

The next corollary of Theorem 3.4 allows one to conclude that a large class of probabilities involving  $L_n$  converge to 0. The general version of this corollary given in Proposition 6.4 is extremely useful in applications. For example, we will use it in section 9 to analyze the Curie-Weiss model of ferromagnetism and in section 10 to motivate the definitions of the sets of equilibrium macrostates for the canonical ensemble and the microcanonical ensemble [Thms. 10.2(c), 10.5(c)].

**Corollary 3.6.** *Let  $A$  be any Borel subset of  $\mathcal{P}_\alpha$  such that  $\overline{A}$  does not contain  $\rho$ . Then  $I_\rho(\overline{A}) > 0$ , and for some  $C < \infty$*

$$P_n\{L_n \in A\} \leq C \exp[-nI_\rho(\overline{A})/2] \rightarrow 0 \text{ as } n \rightarrow \infty.$$

**Proof.** Since  $I_\rho(\gamma) > I_\rho(\rho) = 0$  for any  $\gamma \neq \rho$ , the positivity of  $I_\rho(\overline{A})$  follows from the continuity of  $I_\rho$  on  $\mathcal{P}_\alpha$ . The second assertion is an immediate consequence of the large deviation upper bound applied to  $\overline{A}$  and the positivity of  $I_\rho(\overline{A})$ . ■

Take any  $\varepsilon > 0$ . Applying Corollary 3.6 to the complement of the open ball  $B(\rho, \varepsilon)$  yields  $P_n\{L_n \notin B(\rho, \varepsilon)\} \rightarrow 0$  or equivalently

$$\lim_{n \rightarrow \infty} P_n\{L_n \in B(\rho, \varepsilon)\} = 1.$$

Although this rederives the weak law of large numbers for  $L_n$  as already expressed in (3.1), this second derivation relates the order-1 limit for  $L_n$  to the point  $\rho \in \mathcal{P}_\alpha$  at which the rate function  $I_\rho$  attains its infimum. In this context we call  $\rho$  the equilibrium value of  $L_n$  with respect to the measures  $P_n$ . This limit is the simplest example, and the first of several more complicated but related formulations to be encountered in this paper, of what is commonly called a maximum entropy principle. Following the usual convention in the physical literature, we will continue to use this terminology in referring to such principles

even though we are minimizing the relative entropy — equivalently, maximizing  $-I_\rho(\gamma)$  — rather than maximizing the physical entropy. When  $\rho_k = 1/\alpha$  for each  $k$ , the two quantities differ by a minus sign and an additive constant.

**Maximum Entropy Principle 3.7.**  $\gamma_0 \in \mathcal{P}_\alpha$  is an equilibrium value of  $L_n$  with respect to  $P_n$  if and only if  $\gamma_0$  minimizes  $I_\rho(\gamma)$  over  $\mathcal{P}_\alpha$ ; this occurs if and only if  $\gamma_0 = \rho$ .

In the next section we will present a limit theorem for  $L_n$  whose proof is based on the precise, exponential-order estimates given by the large deviation principle in Theorem 3.4.

## 4 The Most Likely Way for an Unlikely Event To Happen

You participate in a crooked gambling game being played with a loaded die. How can you determine the actual probabilities of each face  $1, 2, \dots, 6$ ? This question uncovers a basic issue in many areas of application. What is the most likely way for an unlikely event to happen?

We use the notation of the preceding section. Thus let  $\alpha \geq 2$  be an integer,  $y_1 < y_2 < \dots < y_\alpha$  a set of  $\alpha$  real numbers,  $\rho_1, \rho_2, \dots, \rho_\alpha$  a set of  $\alpha$  positive real numbers summing to 1,  $\Lambda$  the set  $\{y_1, y_2, \dots, y_\alpha\}$ , and  $P_n$  the product measure on  $\Omega_n = \Lambda^n$  with one dimensional marginals  $\rho = \sum_{k=1}^{\alpha} \rho_k \delta_{y_k}$ . For  $\omega = (\omega_1, \omega_2, \dots, \omega_n) \in \Omega_n$ , we let  $\{X_j, j = 1, \dots, n\}$  be the coordinate functions defined by  $X_j(\omega) = \omega_j$ . The  $X_j$  form a sequence of i.i.d. random variables with common distribution  $\rho$ . For  $\omega \in \Omega_n$  and  $y \in \Lambda$  we also define

$$L_n(y) = L_n(\omega, y) = \frac{1}{n} \sum_{j=1}^n \delta_{X_j(\omega)}\{y\}$$

and the empirical vector

$$L_n = L_n(\omega) = (L_n(\omega, y_1), \dots, L_n(\omega, y_\alpha)) = \frac{1}{n} \sum_{j=1}^n (\delta_{X_j(\omega)}\{y_1\}, \dots, \delta_{X_j(\omega)}\{y_\alpha\}).$$

$L_n$  takes values in the set of probability vectors

$$\mathcal{P}_\alpha = \left\{ \gamma = (\gamma_1, \gamma_2, \dots, \gamma_\alpha) \in \mathbb{R}^\alpha : \gamma_k \geq 0, \sum_{k=1}^{\alpha} \gamma_k = 1 \right\}.$$

In this section we prove a conditioned limit theorem for  $L_n$  that gives an answer to the apparently ambiguous question concerning a crooked gambling game posed in the first paragraph. This limit theorem has the added bonus of giving insight into a basic construction in statistical mechanics. As we will see in section 5, it motivates the form of the Gibbs state for the discrete ideal gas and, by extension, for any statistical mechanical system characterized by conservation of energy. These unexpected theorems are the first indication of the power of Boltzmann's discovery, which gives precise exponential-order estimates for probabilities of the form  $P_n\{L_n \in A\}$ .

The conditioned limit theorem that we will consider has the following form. Suppose that one is given a particular set  $A$  and wants to determine a set  $B$  belonging to a certain class (e.g., open balls) such that the conditioned limit

$$\lim_{n \rightarrow \infty} P_n\{L_n \in B \mid L_n \in A\} = \lim_{n \rightarrow \infty} P_n\{L_n \in B \cap A\} \cdot \frac{1}{P_n\{L_n \in A\}} = 1$$

is valid. Since to exponential order

$$P_n\{L_n \in B \cap A\} \cdot \frac{1}{P_n\{L_n \in A\}} \approx \exp[-n(I_\rho(B \cap A) - I_\rho(A))],$$

one should obtain the conditioned limit if  $B$  satisfies  $I_\rho(B \cap A) = I_\rho(A)$ . If one can determine the point in  $A$  where the infimum of  $I_\rho$  is attained, then one picks  $B$  to contain this point. In the examples involving the loaded die and the discrete ideal gas, such a minimizing point can be determined. It will lead to a second maximum entropy principle for  $L_n$  with respect to the conditional probabilities  $P_n\{\cdot \mid L_n \in A\}$ .

We return to the question concerning the loaded die, using the basic probabilistic model introduced in Example 2.1(b). Upon entering the crooked gambling game, one assigns the equal probabilities  $\rho_k = 1/6$  to each of the 6 faces because one has no additional information. One then observes the game for  $n$  tosses; probabilistically this corresponds to knowing a configuration  $\omega \in \{1, \dots, 6\}^n$ . Based on the value of

$$S_n(\omega) = \sum_{j=1}^n X_j(\omega) = \sum_{j=1}^n \omega_j,$$

one desires to recalculate the probabilities of the 6 faces. Being a mathematician rather than a professional gambler, I will carry this out in the limit  $n \rightarrow \infty$ .

If the die were fair, then the sample mean  $S_n(\omega)/n$  should equal approximately the theoretical mean

$$\bar{y} = \sum_{k=1}^6 k\rho_k = 3.5.$$

Hence let us assume that  $S_n/n \in [z - a, z]$ , where  $a$  is a small positive number and  $1 \leq z - a < z < \bar{y}$ ; a similar result would hold if we assumed that

$S_n/n \in [z, z+a]$ , where  $\bar{y} < z < z+a \leq 6$ . We can now formulate the question concerning the loaded die as the following conditioned limit: determine positive numbers  $\{\rho_k^*, k = 1, \dots, 6\}$  summing to 1 such that

$$\rho_k^* = \lim_{n \rightarrow \infty} P_n\{X_1 = k \mid S_n/n \in [z - a, z]\}.$$

This will be seen to follow from the following more easily answered question: conditioned on the event  $S_n/n \in [z - a, z]$ , determine the most likely configuration  $\rho^* = (\rho_1^*, \dots, \rho_6^*)$  of  $L_n$  in the limit  $n \rightarrow \infty$ . In other words, we want  $\rho^* \in \mathcal{P}_\alpha$  such that for any  $\varepsilon > 0$

$$\lim_{n \rightarrow \infty} P_n\{L_n \in B(\rho^*, \varepsilon) \mid S_n \in [z - a, z]\} = 1.$$

The form of  $\rho^*$  is given in the following theorem; it depends only on  $z$ , not on  $a$ .

We formulate the theorem for a general state space  $\Lambda = \{y_1, \dots, y_\alpha\}$  and a given positive vector  $\rho = (\rho_1, \dots, \rho_\alpha) \in \mathcal{P}_\alpha$ . As above, define

$$S_n = \sum_{j=1}^n X_j \quad \text{and} \quad \bar{y} = \sum_{k=1}^{\alpha} y_k \rho_k$$

and for  $a > 0$  fix a closed interval  $[z - a, z] \subset [y_1, \bar{y})$ . In the definition of  $\rho^{(\beta)}$  we write  $-\beta$  instead of  $\beta$  in order to be consistent with conventions in statistical mechanics.

**Theorem 4.1.** (a) *There exists  $\rho^{(\beta)} \in \mathcal{P}_\alpha$  such that for every  $\varepsilon > 0$*

$$\lim_{n \rightarrow \infty} P_n\{L_n \in B(\rho^{(\beta)}, \varepsilon) \mid S_n/n \in [z - a, z]\} = 1. \quad (4.1)$$

*The quantity  $\rho^{(\beta)} = (\rho_1^{(\beta)}, \dots, \rho_\alpha^{(\beta)})$  has the form*

$$\rho_k^{(\beta)} = \frac{1}{\sum_{j=1}^{\alpha} \exp[-\beta y_j] \rho_j} \cdot \exp[-\beta y_k] \rho_k,$$

*where  $\beta = \beta(z) \in \mathbb{R}$  is the unique value of  $\beta$  satisfying  $\sum_{k=1}^{\alpha} y_k \rho_k^{(\beta)} = z$ .*

(b) *For any continuous function  $f$  mapping  $\mathcal{P}_\alpha$  into  $\mathbb{R}$*

$$\lim_{n \rightarrow \infty} E^{P_n}\{f(L_n) \mid S_n/n \in [z - a, z]\} = f(\rho^{(\beta)}).$$

(c) For each  $j \in \{1, \dots, \alpha\}$

$$\lim_{n \rightarrow \infty} P_n \{X_1 = y_j \mid S_n/n \in [z - a, z]\} = \rho_j^{(\beta)}.$$

We first show that  $\rho^{(\beta)}$  is well defined. For  $r \in \mathbb{R}$  simple calculus gives the following properties of

$$c(r) = \log \left( \sum_{k=1}^{\alpha} \exp[ry_k] \rho_k \right) :$$

$c''(r) > 0$ ,  $c'(r) \rightarrow y_1$  as  $r \rightarrow -\infty$ ,  $c'(0) = \bar{y}$ , and  $c'(r) \rightarrow y_\alpha$  as  $r \rightarrow \infty$ . Hence there exists a unique  $\beta = \beta(z)$  satisfying

$$\begin{aligned} c'(-\beta) &= \frac{1}{\sum_{j=1}^{\alpha} \exp[-\beta y_j] \rho_j} \cdot \sum_{k=1}^{\alpha} y_k \exp[-\beta y_k] \rho_k & (4.2) \\ &= \sum_{k=1}^{\alpha} y_k \rho_k^{(\beta)} = z, \end{aligned}$$

as claimed. Since  $y_1 < z < \bar{y}$ ,  $\beta = \beta(z)$  is positive.

In order to prove the limit

$$\lim_{n \rightarrow \infty} P_n \{L_n \in B(\rho^{(\beta)}, \varepsilon) \mid S_n/n \in [z - a, z]\} = 1.$$

in part (a), we express the conditional probability in (4.1) in terms of the empirical vector  $L_n$ . Define the closed convex set

$$\Gamma(z) = \left\{ \gamma \in \mathcal{P}_\alpha : \sum_{k=1}^{\alpha} y_k \gamma_k \in [z - a, z] \right\},$$

which contains  $\rho^{(\beta)}$ . Since for each  $\omega \in \Omega_n$

$$\frac{1}{n} S_n(\omega) = \sum_{j=1}^n y_k L_n(\omega, y_k),$$

it follows that  $\{\omega \in \Omega_n : S_n(\omega)/n \in [z - a, z]\} = \{\omega \in \Omega_n : L_n(\omega) \in \Gamma(z)\}$ . Hence using the formal notation [see (3.4)]

$$P_n \{L_n \in A\} \approx \exp[-nI_\rho(A)] \text{ as } n \rightarrow \infty,$$

we have for large  $n$

$$\begin{aligned}
& P_n\{L_n \in B(\rho^{(\beta)}, \varepsilon) \mid S_n/n \in [z - a, z]\} \\
&= P_n\{L_n \in B(\rho^{(\beta)}, \varepsilon) \mid L_n \in \Gamma(z)\} \\
&= P_n\{L_n \in B(\rho^{(\beta)}, \varepsilon) \cap \Gamma(z)\} \cdot \frac{1}{P_n\{L_n \in \Gamma(z)\}} \\
&\approx \exp[-n(I_\rho(B(\rho^{(\beta)}, \varepsilon) \cap \Gamma(z)) - I_\rho(\Gamma(z)))].
\end{aligned}$$

The last expression, and thus the probability in the first line of the display, are of order 1 provided

$$I_\rho(B(\rho^{(\beta)}, \varepsilon) \cap \Gamma(z)) = I_\rho(\Gamma(z)). \quad (4.3)$$

The next proposition shows that  $I_\rho$  attains its infimum over  $\Gamma(z)$  at the unique point  $\rho^{(\beta)}$ . This gives (4.3) and motivates the fact that for large  $n$

$$P_n\{L_n \in B(\rho^{(\beta)}, \varepsilon) \mid S_n/n \in [z - a, z]\} \approx 1.$$

It is not difficult to convert these formal calculations into a proof of the limit

$$\lim_{n \rightarrow \infty} P_n\{L_n \in B(\rho^{(\beta)}, \varepsilon) \mid S_n/n \in [z - a, z]\} = 1.$$

The details are omitted.

**Proposition 4.2.**  $I_\rho$  attains its infimum over  $\Gamma(z) = \{\gamma \in \mathcal{P}_\alpha : \sum_{k=1}^\alpha y_k \gamma_k \in [z - a, z]\}$  at the unique point  $\rho^{(\beta)} = (\rho_1^{(\beta)}, \dots, \rho_\alpha^{(\beta)})$  defined in part (a) of Theorem 4.1: for each  $k = 1, \dots, \alpha$

$$\rho_k^{(\beta)} = \frac{1}{\sum_{j=1}^\alpha \exp[-\beta y_j] \rho_j} \cdot \exp[-\beta y_k] \rho_k,$$

where  $\beta = \beta(z) \in \mathbb{R}$  is the unique value of  $\beta$  satisfying  $\sum_{k=1}^\alpha y_k \rho_k^{(\beta)} = z$ .

**Proof.** We recall that  $\beta = \beta(z) > 0$  and that for each  $k \in \{1, \dots, \alpha\}$

$$\frac{\rho_k^{(\beta)}}{\rho_k} = \frac{1}{\sum_{j=1}^\alpha \exp[-\beta y_j] \rho_j} \cdot \exp[-\beta y_k] = \frac{1}{\exp[c(-\beta)]} \cdot \exp[-\beta y_k],$$

where for  $r \in \mathbb{R}$ ,  $c(r) = \log(\sum_{k=1}^{\alpha} \exp[ry_k] \rho_k)$ . Hence for any  $\gamma \in \Gamma(z)$

$$\begin{aligned} I_{\rho}(\gamma) &= \sum_{k=1}^{\alpha} \gamma_k \log \frac{\gamma_k}{\rho_k} = \sum_{k=1}^{\alpha} \gamma_k \log \frac{\gamma_k}{\rho_k^{(\beta)}} + \sum_{k=1}^{\alpha} \gamma_k \log \frac{\rho_k^{(\beta)}}{\rho_k} \\ &= I_{\rho^{(\beta)}}(\gamma) - \beta \sum_{k=1}^{\alpha} y_k \gamma_k - c(-\beta). \end{aligned}$$

Since  $I_{\rho^{(\beta)}}(\rho^{(\beta)}) = 0$  and by (4.2)  $\sum_{k=1}^{\alpha} y_k \rho_k^{(\beta)} = z$ , it follows that

$$I_{\rho}(\rho^{(\beta)}) = I_{\rho^{(\beta)}}(\rho^{(\beta)}) - \beta \sum_{k=1}^{\alpha} y_k \rho_k^{(\beta)} - c(-\beta) = -\beta z - c(-\beta).$$

Now consider any  $\gamma \in \Gamma(z)$ ,  $\gamma \neq \rho^{(\beta)}$ . Since  $I_{\rho^{(\beta)}}(\gamma) \geq 0$  with equality if and only if  $\gamma = \rho^{(\beta)}$  [Lemma 3.2], we obtain

$$\begin{aligned} I_{\rho}(\gamma) &= I_{\rho^{(\beta)}}(\gamma) - \beta \sum_{k=1}^{\alpha} y_k \gamma_k - c(-\beta) \\ &> -\beta \sum_{k=1}^{\alpha} y_k \gamma_k - c(-\beta) \geq -\beta z - c(-\beta) = I_{\rho}(\rho^{(\beta)}). \end{aligned}$$

We conclude that for any  $\gamma \in \Gamma(z)$ ,  $I_{\rho}(\gamma) \geq I_{\rho}(\rho^{(\beta)})$  with equality if and only if  $\gamma = \rho^{(\beta)}$ . Thus  $I_{\rho}$  attains its infimum over  $\Gamma(z)$  at the unique point  $\rho^{(\beta)}$ . ■

Combining this proposition with part (a) of Theorem 4.1 gives the second maximum entropy principle in these lectures.

**Maximum Entropy Principle 4.3.** *Conditioned on the event  $S_n/n \in [z-a, z]$ , the asymptotically most likely configuration of  $L_n$  is  $\rho^{(\beta)}$ , which is the unique  $\gamma \in \mathcal{P}_{\alpha}$  that minimizes  $I_{\rho}(\gamma)$  subject to the constraint that  $\gamma \in \Gamma(z)$ . In statistical mechanical terminology,  $\rho^{(\beta)}$  is the equilibrium macrostate of  $L_n$  with respect to the conditional probabilities  $P_n\{\cdot \mid S_n/n \in [z-a, z]\}$ .*

Part (b) of Theorem 4.1 states that for any continuous function  $f$  mapping  $\mathcal{P}_{\alpha}$  into  $\mathbb{R}$

$$\lim_{n \rightarrow \infty} E^{P_n}\{f(L_n) \mid S_n/n \in [z-a, z]\} = f(\rho^{(\beta)}).$$



This is an immediate consequence of part (a) and the continuity of  $f$ . Part (b) of Theorem 4.1 is another expression of the Maximum Entropy Principle 4.3.

Let  $y_k$  be any point in  $\Lambda$ . As in [20, p. 87], we prove part (c) of Theorem 4.1 by relating the conditional probability  $P_n\{X_1 = y_k \mid S_n/n \in [z - a, z]\}$  to the conditional expectation  $E^{P_n}\{f(L_n) \mid S_n/n \in [z - a, z]\}$  considered in part (b). Given  $\varphi$  any function mapping  $\Lambda$  into  $\mathbb{R}$ , we define a continuous function on  $\mathcal{P}_\alpha$  by

$$f(\gamma) = \sum_{k=1}^{\alpha} \varphi(y_k) \gamma_k.$$

Since  $f(L_n) = \sum_{k=1}^{\alpha} \varphi(y_k) L_n(y_k) = \frac{1}{n} \sum_{j=1}^n \varphi(X_j)$ , by symmetry and part (b)

$$\begin{aligned} & \lim_{n \rightarrow \infty} E^{P_n}\{\varphi(X_1) \mid S_n/n \in [z - a, z]\} \\ &= \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{j=1}^n E^{P_n}\{\varphi(X_j) \mid S_n/n \in [z - a, z]\} \\ &= \lim_{n \rightarrow \infty} E^{P_n}\{f(L_n) \mid S_n/n \in [z - a, z]\} \\ &= f(\rho^{(\beta)}) = \sum_{k=1}^{\alpha} \varphi(y_k) \rho_k^{(\beta)}. \end{aligned}$$

Setting  $\varphi = 1_{y_j}$  yields the limit in part (c) of Theorem 4.1:

$$\lim_{n \rightarrow \infty} P_n\{X_1 = y_j \mid S_n/n \in [z - a, z]\} = \rho_j^{(\beta)}.$$

With some additional work one can generalize part (a) of Theorem 4.1 by proving that with respect to the conditional probabilities  $P_n\{\cdot \mid S_n/n \in [z - a, a]\}$ ,  $L_n$  satisfies the large deviation principle on  $\mathcal{P}_\alpha$  with rate function

$$I(\gamma) = \begin{cases} I_\rho(\gamma) - I_\rho(\Gamma(z)) & \text{if } \gamma \in \Gamma(z) \\ \infty & \text{if } \gamma \in \mathcal{P}_\alpha \setminus \Gamma(z). \end{cases}$$

This large deviation principle is closely related to the large deviation principle for statistical mechanical models with respect to the microcanonical ensemble, which will be considered in Theorem 10.5.

In the next section we will show how calculations analogous to those used to motivate Theorem 4.1 can be used to derive the form of the Gibbs state for the discrete ideal gas.

## 5 Gibbs States for Models in Statistical Mechanics

The discussion in the preceding section concerning a loaded die applies with minor changes to the discrete ideal gas, introduced in part (c) of Examples 2.1. We continue to use the notation of the preceding sections. Thus let  $\alpha \geq 2$  be an integer,  $y_1 < y_2 < \dots < y_\alpha$  a set of  $\alpha$  real numbers,  $\rho_1, \rho_2, \dots, \rho_\alpha$  a set of  $\alpha$  positive real numbers summing to 1,  $\Lambda$  the set  $\{y_1, y_2, \dots, y_\alpha\}$ , and  $P_n$  the product measure on  $\Omega_n = \Lambda^n$  with one dimensional marginals  $\rho = \sum_{k=1}^{\alpha} \rho_k \delta_{y_k}$ . For  $\omega = (\omega_1, \omega_2, \dots, \omega_n) \in \Omega_n$ , we let  $\{X_j, j = 1, \dots, n\}$  be the coordinate functions defined by  $X_j(\omega) = \omega_j$ . The  $X_j$  form a sequence of i.i.d. random variables with common distribution  $\rho$ .

The discrete ideal gas consists of  $n$  identical, noninteracting particles, each having  $\alpha$  possible energy levels  $y_1, y_2, \dots, y_\alpha$ . For  $\omega \in \Omega_n$  we write  $H_n(\omega)$  in place of  $S_n(\omega) = \sum_{j=1}^n \omega_j$ ;  $H_n(\omega)$  denotes the total energy in the configuration  $\omega$ . In the absence of further information, one assigns the equal probabilities  $\rho_k = 1/\alpha$  to each of the  $y_k$ 's. Defining  $\bar{y} = \sum_{k=1}^{\alpha} y_k \rho_k$ , suppose that the energy per particle,  $H_n/n$ , is conditioned to lie in an interval  $[z - a, z]$ , where  $a$  is a small positive number and  $y_1 \leq z - a < z < \bar{y}$ . According to part (c) of Theorem 4.1, for each  $k \in \{1, \dots, \alpha\}$

$$\begin{aligned} & \lim_{n \rightarrow \infty} P_n \{X_1 = y_k \mid H_n/n \in [z - a, z]\} \\ &= \rho_k^{(\beta)} = \frac{1}{\sum_{j=1}^{\alpha} \exp[-\beta y_j] \rho_j} \cdot \exp[-\beta y_k] \rho_k, \end{aligned}$$

where  $\beta = \beta(z) \in \mathbb{R}$  is the unique value of  $\beta$  satisfying  $\sum_{k=1}^{\alpha} y_k \rho_k^{(\beta)} = z$ .

Let  $t \geq 2$  be a positive integer. The limit in the last display leads to a natural question. Conditioned on  $H_n/n \in [z - a, z]$ , as  $n \rightarrow \infty$  what is the limiting conditional distribution of the random variables  $X_1, \dots, X_t$ , which represent the energy levels of the first  $t$  particles? Although  $X_1, \dots, X_t$  are independent with respect to the original product measure  $P_n$ , this independence is lost when  $P_n$  is replaced by the conditional distribution  $P_n\{\cdot \mid H_n/n \in [z - a, z]\}$ . Hence the answer given in the next theorem is somewhat surprising: with respect to  $P_n\{\cdot \mid H_n/n \in [z - a, z]\}$ , the limiting distribution is the product measure on  $\Omega_t$  with one-dimensional marginals  $\rho^{(\beta)}$ . In other words, in the limit  $n \rightarrow \infty$  the

independence of  $X_1, \dots, X_t$  is regained. The theorem leads to, and in a sense motivates, the concept of the Gibbs state of the discrete ideal gas. We will end the section by discussing Gibbs states for this and other statistical mechanical models. As in Theorem 4.1, a theorem analogous to the following would hold if  $[z - a, z] \subset [y_1, \bar{y})$  were replaced by  $[z, z + a] \subset (\bar{y}, y_\alpha]$ .

**Theorem 5.1.** *Given  $t \in \mathbb{N}$ ,  $y_{k_1}, \dots, y_{k_t} \in \Lambda$ , and  $[z - a, z] \subset [y_1, \bar{y})$ ,*

$$\lim_{n \rightarrow \infty} P_n \{X_1 = y_{k_1}, \dots, X_t = y_{k_t} \mid H_n/n \in [z - a, z]\} = \prod_{j=1}^t \rho_{k_j}^{(\beta)}. \quad (5.1)$$

**Comments on the Proof.** We consider  $t = 2$ ; arbitrary  $t \in \mathbb{N}$  can be handled similarly. For  $\omega \in \Omega_n$  and  $i, j \in \{1, \dots, \alpha\}$  define

$$\begin{aligned} L_{n,2}(\{y_i, y_j\}) &= L_{n,2}(\omega, \{y_i, y_j\}) \\ &= \frac{1}{n} \left( \sum_{j=1}^{n-1} \delta_{X_j(\omega), X_{j+1}(\omega)} \{y_i, y_j\} + \delta_{X_n(\omega), X_1(\omega)} \{y_i, y_j\} \right). \end{aligned}$$

This counts the relative frequency with which the pair  $\{y_i, y_j\}$  appears in the configuration  $(\omega_1, \dots, \omega_n, \omega_1)$ . We then define the empirical pair vector

$$L_{n,2} = \{L_{n,2}(\{y_i, y_j\}), i, j = 1, \dots, \alpha\}.$$

This takes values in the set  $\mathcal{P}_{\alpha,2}$  consisting of all  $\tau = \{\tau_{i,j}, i, j = 1, \dots, \alpha\}$  satisfying  $\tau_{i,j} \geq 0$  and  $\sum_{i,j=1}^{\alpha} \tau_{i,j} = 1$ . Suppose one can show that  $\tau^* = \{\rho_i^{(\beta)} \rho_j^{(\beta)}, i, j = 1, \dots, \alpha\}$  has the property that for every  $\varepsilon > 0$

$$\lim_{n \rightarrow \infty} P_n \{L_{n,2} \in B(\tau^*, \varepsilon) \mid H_n/n \in [z - a, z]\} = 1. \quad (5.2)$$

Then as in Theorem 4.1, it will follow that

$$\lim_{n \rightarrow \infty} P_n \{X_1 = y_i, X_2 = y_j \mid H_n/n \in [z - a, z]\} = \rho_i^{(\beta)} \rho_j^{(\beta)}.$$

Like the analogous limit in part (a) of Theorem 4.1, (5.2) can be proved by showing that the sequence  $\{L_{n,2}, n \in \mathbb{N}\}$  satisfies the large deviation principle

on  $\mathcal{P}_{\alpha,2}$  [33, §I.5] and that the rate function attains its infimum over an appropriately defined, closed convex subset of  $\mathcal{P}_{\alpha,2}$  at the unique point  $\tau^*$  [cf. (4.3)]. The details are omitted. ■

The quantity appearing on the right side of (5.1) defines a probability measure  $P_{t,\beta}$  on  $\Omega_t$  that equals the product measure with one-dimensional marginals  $\rho^{(\beta)}$ . In the notation of Theorem 5.1,

$$P_{t,\beta}\{X_1 = y_{k_1}, \dots, X_t = y_{k_t}\} = \prod_{j=1}^t \rho_{k_j}^{(\beta)}.$$

$P_{t,\beta}$  can be written in terms of the total energy  $H_t(\omega) = \sum_{j=1}^t \omega_j$ : for  $\omega \in \Omega_t$

$$P_{t,\beta}\{\omega\} = \prod_{j=1}^t \rho^{(\beta)}\{\omega_j\} = \frac{1}{Z_t(\beta)} \cdot \exp[-\beta H_t(\omega)] P_t\{\omega\},$$

where  $P_t\{\omega\} = \prod_{j=1}^t \rho\{\omega_j\} = 1/\alpha^t$ ,

$$Z_t(\beta) = \sum_{\omega \in \Omega_t} \exp[-\beta H_t(\omega)] P_t\{\omega\} = \left( \sum_{k=1}^{\alpha} \exp[-\beta y_k] \rho_k \right)^t,$$

and  $\beta = \beta(z) \in \mathbb{R}$  is the unique value of  $\beta$  satisfying  $\sum_{k=1}^{\alpha} y_k \rho_k^{(\beta)} = z$ .

Theorem 5.1 can be motivated by a non-large deviation calculation that we present using a formal notation [64]. Since  $\bar{y} = \sum_{k=1}^{\alpha} y_k \rho_k = E^{P_n}\{X_1\}$ , by the weak law of large numbers  $P_n\{H_n/n \sim \bar{y}\} \approx 1$  for large  $n$ . Since the conditioning is on a set of probability close to 1, one expects that

$$\begin{aligned} & \lim_{n \rightarrow \infty} P_n\{X_1 = y_{k_1}, \dots, X_t = y_{k_t} \mid H_n/n \sim \bar{y}\} \\ &= \lim_{n \rightarrow \infty} P_n\{X_1 = y_{k_1}, \dots, X_t = y_{k_t}\} \\ &= \prod_{j=1}^t \rho_{k_j} = P_t\{X_1 = y_{k_1}, \dots, X_t = y_{k_t}\}. \end{aligned}$$

Now take  $z \neq \bar{y}$  and for any  $\beta > 0$  let  $P_{n,\beta}$  denote the product measure on  $\Omega_n$  with one-dimensional marginals  $\rho^{(\beta)}$ . A short calculation shows that for any

$\beta > 0$

$$\begin{aligned} P_n\{X_1 = y_{k_1}, \dots, X_t = y_{k_t} \mid H_n/n \sim z\} \\ = P_{n,\beta}\{X_1 = y_{k_1}, \dots, X_t = y_{k_t} \mid H_n/n \sim z\}. \end{aligned}$$

If one picks  $\beta = \beta(z)$  such that  $z = \sum_{k=1}^{\alpha} y_k \rho_k^{(\beta(z))} = E^{P_{n,\beta(z)}}\{X_1\}$ , then by the weak law of large numbers  $P_{n,\beta(z)}\{H_n/n \sim z\} \approx 1$ , and since the conditioning is on a set of probability close to 1, again one expects that

$$\begin{aligned} \lim_{n \rightarrow \infty} P_n\{X_1 = y_{k_1}, \dots, X_t = y_{k_t} \mid H_n/n \sim z\} \\ = \lim_{n \rightarrow \infty} P_{n,\beta(z)}\{X_1 = y_{k_1}, \dots, X_t = y_{k_t} \mid H_n/n \sim z\} \\ = \lim_{n \rightarrow \infty} P_{n,\beta(z)}\{X_1 = y_{k_1}, \dots, X_t = y_{k_t}\} \\ = \prod_{j=1}^t \rho_{k_j}^{(\beta(z))} = P_{t,\beta(z)}\{X_1 = y_{k_1}, \dots, X_t = y_{k_t}\}. \end{aligned}$$

This is consistent with Theorem 5.1.

For any subset  $B$  of  $\Omega_t$ , (5.1) implies that

$$\lim_{n \rightarrow \infty} P_n\{(X_1, \dots, X_t) \in B \mid H_n/n \in [z - a, z]\} = P_{t,\beta}\{B\}. \quad (5.3)$$

Since  $\sum_{\omega \in \Omega_t} [H_t(\omega)/t] P_{t,\beta}\{\omega\} = \sum_{k=1}^{\alpha} y_k \rho_k^{(\beta)}$ , the constraint on  $\beta = \beta(z)$  can be expressed as a constraint on  $P_{t,\beta}$ :

$$\text{choose } \beta = \beta(z) \text{ so that } \sum_{\omega \in \Omega_t} [H_t(\omega)/t] P_{t,\beta}\{\omega\} = z. \quad (5.4)$$

The conditional probability on the left side of (5.3) is known as the microcanonical ensemble, and the probability on the right side of (5.3) as the canonical ensemble or Gibbs state. This limit expresses the equivalence of the two ensembles provided  $\beta$  is chosen in accordance with (5.4). Since the canonical ensemble has a much simpler form than the microcanonical ensemble, one usually prefers to work with the former. One can interpret  $\beta$  as a parameter that is proportional to the inverse temperature. In section 10 we will discuss related issues involving the equivalence of ensembles in a much broader setting, showing that for models in which interactions are present, in general the microcanonical

formulation gives rise to a richer set of equilibrium properties than the canonical formulation.

This discussion motivates the definition of the Gibbs states for a wide class of statistical mechanical models that are defined in terms of an energy function. We will write the energy function, or Hamiltonian, and the corresponding Gibbs state as  $H_n$  and  $P_{n,\beta}$  rather than as  $H_t$  and  $P_{t,\beta}$ , as we did in the preceding paragraph. The notation of section 2 is used. Thus  $P_n$  is the product measure on the set of subsets of  $\Omega_n = \Lambda^n$  with one-dimensional marginals  $\rho$ . Noninteracting systems such as the discrete ideal gas have Hamiltonians of the form  $H_n(\omega) = \sum_{j=1}^n H_{n,j}(\omega_j)$ , where each  $H_{n,j}$  is a function only of  $\omega_j$ . In the next definition we do not restrict to this case.

**Definition 5.2.** *Let  $H_n$  be a function mapping  $\Omega_n$  into  $\mathbb{R}$ ;  $H_n(\omega)$  defines the energy of the configuration  $\omega$  and is known as a Hamiltonian. Let  $\beta$  be a parameter proportional to the inverse temperature. Then the canonical ensemble, or the Gibbs state, is the probability measure*

$$P_{n,\beta}\{\omega\} = \frac{1}{Z_n(\beta)} \cdot \exp[-\beta H_n(\omega)] P_n\{\omega\} \text{ for } \omega \in \Omega_n,$$

where  $Z_n(\beta)$  is the normalization factor that makes  $P_{n,\beta}$  a probability measure. That is,

$$Z_n(\beta) = \sum_{\omega \in \Omega_n} \exp[-\beta H_n(\omega)] P_n\{\omega\}.$$

We call  $Z_n(\beta)$  the partition function. For  $B \subset \Omega_n$  we define

$$P_{n,\beta}\{B\} = \sum_{\omega \in B} P_{n,\beta}\{\omega\}.$$

One can also characterize Gibbs states in terms of a maximum entropy principle [67, p. 6]. Given  $n \in \mathbb{N}$  and a Hamiltonian  $H_n$ , let  $B_n \subset \mathbb{R}$  denote the smallest closed interval containing the range of  $\{H_n(\omega)/n, \omega \in \Omega_n\}$ . For each  $z \in B_n^\circ$ , the interior of  $B_n$ , define  $C_n(z)$  to be the set of probability measures  $Q$  on  $\Omega_n$  satisfying the energy constraint  $\sum_{\omega \in \Omega_n} [H_n(\omega)/n] Q\{\omega\} = z$ .

**Maximum Entropy Principle 5.3.** *Let  $n \in \mathbb{N}$  and a Hamiltonian  $H_n : \Omega_n \mapsto \mathbb{R}$  be given. The following conclusions hold.*

(a) *For each  $z \in B_n^\circ$  there exists a unique  $\beta = \beta(z) \in \mathbb{R}$  such that  $P_{n,\beta} \in C_n(z)$ .*

(b) *The relative entropy  $I_{P_n}$  attains its infimum over  $C_n(z)$  at the unique measure  $P_{n,\beta}$ , and  $I_{P_n}(P_{n,\beta}) = nI_\rho(\rho^\beta)$ .*

Part (a) can be proved by a calculation similar to that given after the statement of Theorem 4.1 while part (b) can be proved like Proposition 4.2. We leave the details to the reader.

In the next section we formulate the general concepts of a large deviation principle and a Laplace principle. Subsequent sections will apply the theory of large deviations to study interacting systems in statistical mechanics.

## 6 Generalities: Large Deviation Principle and Laplace Principle

In Theorem 3.4 we formulated Sanov's Theorem, which is the large deviation principle for the empirical vectors  $L_n$  on the space  $\mathcal{P}_\alpha$  of probability vectors in  $\mathbb{R}^\alpha$ . Applications of the theory of large deviations to models in statistical mechanics require large deviation principles in much more general settings. As we will see in section 9, analyzing the Curie-Weiss model of ferromagnetism involves a large deviation principle on the closed interval  $[-1, 1]$  for the sample means of i.i.d. random variables. Analyzing the Ising model in dimensions  $d \geq 2$  is much more complicated. It involves a large deviation principle on the space of translation invariant probability measures on  $\{-1, 1\}^{\mathbb{Z}^d}$  [35, §11]. In section 11, our analysis of models of two-dimensional turbulence involves a large deviation principle on the space of probability measures on  $T^2 \times \mathcal{Y}$ , where  $T^2$  is the unit torus in  $\mathbb{R}^2$  and  $\mathcal{Y}$  is a compact subset of  $\mathbb{R}$ .

In order to define the general concept of a large deviation principle, we need some notation. First, for each  $n \in \mathbb{N}$  let  $(\Omega_n, \mathcal{F}_n, P_n)$  be a probability space. Thus  $\Omega_n$  is a set of points,  $\mathcal{F}_n$  is a  $\sigma$ -algebra of subsets of  $\Omega_n$ , and  $P_n$  is a probability measure on  $\mathcal{F}_n$ . An example is given by the basic model in section 2, where  $\Omega_n = \Lambda^n = \{y_1, y_2, \dots, y_\alpha\}^n$ ,  $\mathcal{F}_n$  is the set of all subsets of  $\Omega_n$ , and  $P_n$  is the product measure with one-dimensional marginals  $\rho$ .

Second, let  $\mathcal{X}$  be a complete, separable metric space or, as it is often called, a Polish space. Elementary examples are  $\mathcal{X} = \mathbb{R}^d$  for  $d \in \mathbb{N}$ ;  $\mathcal{X} = \mathcal{P}_\alpha$ , the set of probability vectors in  $\mathbb{R}^\alpha$ ; and in the notation of the basic probabilistic model in section 2,  $\mathcal{X}$  equal to the closed bounded interval  $[y_1, y_\alpha]$ . A class of Polish spaces arising naturally in applications is obtained by taking a Polish space  $\mathcal{Y}$  and considering the space  $\mathcal{P}(\mathcal{Y})$  of probability measures on  $\mathcal{Y}$ . We say that a sequence  $\{\Pi_n, n \in \mathbb{N}\}$  in  $\mathcal{P}(\mathcal{Y})$  converges weakly to  $\Pi \in \mathcal{P}(\mathcal{Y})$ , and write  $\Pi_n \Rightarrow \Pi$ , if  $\int_{\mathcal{Y}} f d\Pi_n \rightarrow \int_{\mathcal{Y}} f d\Pi$  for all bounded, continuous functions  $f$  mapping  $\mathcal{Y}$  into  $\mathbb{R}$ . A fundamental fact is that there exists a metric  $m$  on  $\mathcal{P}(\mathcal{Y})$  such that  $\Pi_n \Rightarrow \Pi$  if and only if  $m(\Pi, \Pi_n) \rightarrow 0$  and  $\mathcal{P}(\mathcal{Y})$  is a Polish space with respect to  $m$  [45, §3.1].



Third, for each  $n \in \mathbb{N}$  let  $Y_n$  be a random variable mapping  $\Omega_n$  into  $\mathcal{X}$ . For example, with  $\mathcal{X} = \mathcal{P}_\alpha$  let  $Y_n = L_n$ , or with  $\mathcal{X} = [y_1, y_\alpha]$  let  $Y_n = \sum_{j=1}^n X_j/n$ , where  $X_j(\omega) = \omega_j$  for  $\omega \in \Omega_n = \Lambda^n$ .

Finally, let  $I$  be a function mapping the complete, separable metric space  $\mathcal{X}$  into  $[0, \infty]$ .  $I$  is called a rate function if  $I$  has compact level sets; i.e., for all  $M < \infty$   $\{x \in \mathcal{X} : I(x) \leq M\}$  is compact. This technical regularity condition implies that  $I$  has closed level sets or equivalently that  $I$  is lower semicontinuous. Hence, if  $\mathcal{X}$  is compact, then the lower semicontinuity of  $I$  implies that  $I$  has compact level sets. For any subset  $A$  of  $\mathcal{X}$  we define  $I(A) = \inf_{x \in A} I(x)$ . When  $\mathcal{X} = \mathcal{P}_\alpha$ , an example of a rate function is the relative entropy  $I_\rho$  with respect to  $\rho$ ; when  $\mathcal{X} = [y_1, y_\alpha]$ , any continuous function  $I$  mapping  $[y_1, y_\alpha]$  into  $[0, \infty)$  is a rate function.

We next define the concept of a large deviation principle. If  $Y_n$  satisfies the large deviation principle with rate function  $I$ , then we summarize this by the formal notation

$$P_n\{Y_n \in dx\} \asymp \exp[-nI(x)] dx.$$

**Definition 6.1 (Large Deviation Principle).** *Let  $\{(\Omega_n, \mathcal{F}_n, P_n), n \in \mathbb{N}\}$  be a sequence of probability spaces,  $\mathcal{X}$  a complete, separable metric space,  $\{Y_n, n \in \mathbb{N}\}$  a sequence of random variables such that  $Y_n$  maps  $\Omega_n$  into  $\mathcal{X}$ , and  $I$  a rate function on  $\mathcal{X}$ . Then  $Y_n$  satisfies the large deviation principle on  $\mathcal{X}$  with rate function  $I$  if the following two limits hold.*

**Large deviation upper bound.** *For any closed subset  $F$  of  $\mathcal{X}$*

$$\limsup_{n \rightarrow \infty} \frac{1}{n} \log P_n\{Y_n \in F\} \leq -I(F).$$

**Large deviation lower bound.** *For any open subset  $G$  of  $\mathcal{X}$*

$$\liminf_{n \rightarrow \infty} \frac{1}{n} \log P_n\{Y_n \in G\} \geq -I(G).$$

We next explore several consequences of this definition. It is reassuring that a large deviation principle has a unique rate function. The following result is proved in Theorem II.3.2 in [33].

**Theorem 6.2.** *If  $Y_n$  satisfies the large deviation principle on  $\mathcal{X}$  with rate function  $I$  and with rate function  $J$ , then  $I(x) = J(x)$  for all  $x \in \mathcal{X}$ .*

The next theorem gives a condition that guarantees the existence of large deviation limits. The proof is analogous to the proof of Corollary 3.5.

**Theorem 6.3.** *Assume that  $Y_n$  satisfies the large deviation principle on  $\mathcal{X}$  with rate function  $I$ . Let  $A$  be a Borel subset of  $\mathcal{X}$  having closure  $\bar{A}$  and interior  $A^\circ$  and satisfying  $I(\bar{A}) = I(A^\circ)$ . Then*

$$\lim_{n \rightarrow \infty} \frac{1}{n} \log P\{Y_n \in A\} = -I(A).$$

**Proof.** We evaluate the large deviation upper bound for  $F = \bar{A}$  and the large deviation lower bound for  $G = A^\circ$ . Since  $\bar{A} \supset A \supset A^\circ$ , it follows that

$$\begin{aligned} I(\bar{A}) &\geq \limsup_{n \rightarrow \infty} \frac{1}{n} \log P\{Y_n \in \bar{A}\} \geq \limsup_{n \rightarrow \infty} \frac{1}{n} \log P\{Y_n \in A\} \\ &\geq \liminf_{n \rightarrow \infty} \frac{1}{n} \log P\{Y_n \in A\} \geq \liminf_{n \rightarrow \infty} \frac{1}{n} \log P\{Y_n \in A^\circ\} \geq I(A^\circ). \end{aligned}$$

By hypothesis the two extreme terms are equal, and so the theorem follows. ■

The next proposition states useful facts concerning the infimum of a rate function over the entire space and the use of the large deviation principle to show the convergence of a class of probabilities to 0. Part (b) generalizes Corollary 3.6.

**Proposition 6.4.** *Suppose that  $Y_n$  satisfies the large deviation principle on  $\mathcal{X}$  with rate function  $I$ . The following conclusions hold.*

(a) *The infimum of  $I$  over  $\mathcal{X}$  equals 0, and the set of  $x \in \mathcal{X}$  for which  $I(x) = 0$  is nonempty and compact.*

(b) *Define  $\mathcal{E}$  to be the nonempty, compact set of  $x \in \mathcal{X}$  for which  $I(x) = 0$  and let  $A$  be a Borel subset of  $\mathcal{X}$  such that  $\bar{A} \cap \mathcal{E} = \emptyset$ . Then  $I(\bar{A}) > 0$ , and for some  $C < \infty$*

$$P_n\{Y_n \in A\} \leq C \exp[-nI(\bar{A})/2] \rightarrow 0 \text{ as } n \rightarrow \infty.$$

**Proof.** (a) We evaluate the large deviation upper bound for  $F = \mathcal{X}$  and the large deviation lower bound for  $G = \mathcal{X}$ , obtaining  $I(\mathcal{X}) = 0$ . Since  $I$  has compact level sets, the set of  $x \in \mathcal{X}$  for which  $I(x) = 0$  is nonempty and compact. This gives part (a).

(b) If  $I(\overline{A}) > 0$ , then the desired upper bound follows immediately from the large deviation upper bound. We prove that  $I(\overline{A}) > 0$  by contradiction. If  $I(\overline{A}) = 0$ , then there exists a sequence  $x_n$  such that  $\lim_{n \rightarrow \infty} I(x_n) = 0$ . Since  $I$  has compact level sets and  $\overline{A}$  is closed, there exists a subsequence  $x_{n'}$  converging to an element  $x \in \overline{A}$ . Since  $I$  is lower semicontinuous, it follows that  $I(x) = 0$  and thus that  $x \in \mathcal{E}$ . This contradicts the assumption that  $\overline{A} \cap \mathcal{E} = \emptyset$ . The proof of the proposition is complete. ■

In the next section we will prove Cramér's Theorem, which is the large deviation principle for the sample means of i.i.d. random variables. Here is a statement of the theorem. The rate function is defined by a variational formula that in general cannot be evaluated explicitly. We denote by  $\langle \cdot, \cdot \rangle$  the inner product on  $\mathbb{R}^d$ .

**Theorem 6.5 (Cramér's Theorem).** *Let  $\{X_j, j \in \mathbb{N}\}$  be a sequence of i.i.d. random vectors taking values in  $\mathbb{R}^d$  and satisfying  $E\{\exp\langle t, X_1 \rangle\} < \infty$  for all  $t \in \mathbb{R}^d$ . We define the sample means  $S_n/n = \sum_{j=1}^n X_j/n$  and the cumulant generating function  $c(t) = \log E\{\exp\langle t, X_1 \rangle\}$ . The following conclusions hold.*

(a) *The sequence of sample means  $S_n/n$  satisfies the large deviation principle on  $\mathbb{R}^d$  with rate function  $I(x) = \sup_{t \in \mathbb{R}^d} \{\langle t, x \rangle - c(t)\}$ .*

(b)  *$I$  is a convex, lower semicontinuous function on  $\mathbb{R}^d$ , and it attains its infimum of 0 at the unique point  $x_0 = E\{X_1\}$ .*

For application in section 9, we next state a special case of Cramér's Theorem, for which the rate function can be given explicitly.

**Corollary 6.6.** *In the basic probability model of section 2, let  $\Lambda = \{-1, 1\}$  and  $\rho = (\frac{1}{2}, \frac{1}{2})$ , which corresponds to the probability measure  $\rho = \frac{1}{2}\delta_{-1} + \frac{1}{2}\delta_1$  on  $\Lambda$ . For  $\omega \in \Omega_n$  define  $S_n(\omega) = \sum_{j=1}^n \omega_j$ . Then the sequence of sample means*

$S_n/n$  satisfies the large deviation principle on the closed interval  $[-1, 1]$  with rate function

$$I(x) = \frac{1}{2}(1-x)\log(1-x) + \frac{1}{2}(1+x)\log(1+x). \quad (6.1)$$

**Proof.** In this case  $c(t) = \log(\frac{1}{2}[e^t + e^{-t}])$ . The function  $c(t)$  satisfies  $c''(t) > 0$  for all  $t$ , and the range of  $c'$  equals  $(-1, 1)$ . Hence for any  $x \in (-1, 1)$  the supremum in the definition of  $I$  is attained at the unique  $t = t(x)$  satisfying  $c'(t(x)) = x$ . One easily verifies that  $t(x) = \frac{1}{2}\log[(1+x)/(1-x)]$  and that  $I(x) = t(x) \cdot x - c(t(x))$  is given by (6.1). When  $x = 1$  or  $x = -1$ , the supremum in the definition of  $I(x)$  is not attained but equals  $\log 2$ . ■

Corollary 6.6 is easy to motivate using the formal notation of “Theorem” 3.3. For any  $x \in [-1, 1]$   $S_n(\omega)/n \sim x$  if and only if approximately  $\frac{n}{2}(1-x)$  of the  $\omega_j$ 's equal  $-1$  and approximately  $\frac{n}{2}(1+x)$  of the  $\omega_j$ 's equal  $1$ . Hence

$$\begin{aligned} P_n\{S_n/n \sim x\} &\approx P_n\{L_n(-1) = \frac{1}{2}(1-x), L_n(1) = \frac{1}{2}(1+x)\} \\ &\approx \exp[-nI_\rho(\frac{1}{2}(1-x), \frac{1}{2}(1+x))] = \exp[-nI(x)]. \end{aligned}$$

For application in section 11, we state a general version of Sanov's Theorem, which gives the large deviation principle for the sequence of empirical measures of i.i.d. random variables. Let  $(\Omega, \mathcal{F}, P)$  be a probability space,  $\mathcal{Y}$  a complete, separable metric space,  $\rho$  a probability measure on  $\mathcal{Y}$ , and  $\{X_j, j \in \mathbb{N}\}$  a sequence of i.i.d. random variables mapping  $\Omega$  into  $\mathcal{Y}$  and having the common distribution  $\rho$ . For  $n \in \mathbb{N}$ ,  $\omega \in \Omega$ , and  $A$  any Borel subset of  $\mathcal{Y}$  we define the empirical measure

$$L_n(A) = L_n(\omega, A) = \frac{1}{n} \sum_{j=1}^n \delta_{X_j(\omega)}\{A\},$$

where for  $y \in \mathcal{Y}$ ,  $\delta_y\{A\}$  equals 1 if  $y \in A$  and 0 if  $y \notin A$ . For each  $\omega$ ,  $L_n(\omega, \cdot)$  is a probability measure on  $\mathcal{Y}$ . Hence the sequence  $L_n$  takes values in the complete, separable metric space  $\mathcal{P}(\mathcal{Y})$ .

**Theorem 6.7 (Sanov's Theorem).** *The sequence  $L_n$  satisfies the large deviation principle on  $\mathcal{P}(\mathcal{Y})$  with rate function given by the relative entropy with*

respect to  $\rho$ . For  $\gamma \in \mathcal{P}(\mathcal{Y})$  this quantity is defined by

$$I_\rho(\gamma) = \begin{cases} \int_{\mathcal{Y}} \left( \log \frac{d\gamma}{d\rho} \right) d\gamma & \text{if } \gamma \ll \rho \\ \infty & \text{otherwise.} \end{cases}$$

This theorem is proved, for example, in [20, §6.2] and in [31, Ch. 2]. As we will see in the next section, if the support of  $\rho$  is a finite set  $\Lambda \subset \mathbb{R}$ , then Theorem 6.7 reduces to Theorem 3.4.

The concept of a Laplace principle will be useful in the analysis of statistical mechanical models.

**Definition 6.8 (Laplace Principle).** Let  $\{(\Omega_n, \mathcal{F}_n, P_n), n \in \mathbb{N}\}$  be a sequence of probability spaces,  $\mathcal{X}$  a complete, separable metric space,  $\{Y_n, n \in \mathbb{N}\}$  a sequence of random variables such that  $Y_n$  maps  $\Omega_n$  into  $\mathcal{X}$ , and  $I$  a rate function on  $\mathcal{X}$ . Then  $Y_n$  satisfies the Laplace principle on  $\mathcal{X}$  with rate function  $I$  if for all bounded, continuous functions  $f$  mapping  $\mathcal{X}$  into  $\mathbb{R}$

$$\lim_{n \rightarrow \infty} \frac{1}{n} \log \int_{\Omega_n} \exp[nf(Y_n)] dP_n = \sup_{x \in \mathcal{X}} \{f(x) - I(x)\}.$$

Suppose that  $Y_n$  satisfies the large deviation principle on  $\mathcal{X}$  with rate function  $I$ . Then substituting  $P_n\{Y_n \in dx\} \asymp \exp[-nI(x)] dx$  gives

$$\begin{aligned} \lim_{n \rightarrow \infty} \frac{1}{n} \log \int_{\Omega_n} \exp[nf(Y_n)] dP_n &= \frac{1}{n} \log \int_{\mathcal{X}} \exp[nf(x)] P_n\{Y_n \in dx\} \\ &\approx \frac{1}{n} \log \int_{\mathcal{X}} \exp[nf(x)] \exp[-nI(x)] dx. \end{aligned}$$

By analogy with Laplace's method on  $\mathbb{R}$ , the main contribution to the last integral should come from the largest value of the integrand, and thus the following limit should hold:

$$\lim_{n \rightarrow \infty} \frac{1}{n} \log \int_{\Omega_n} \exp[nf(Y_n)] dP_n = \sup_{x \in \mathcal{X}} \{f(x) - I(x)\}.$$

Hence it is plausible that  $Y_n$  satisfies the Laplace principle with rate function  $I$ . In fact, we have the following stronger result [31, Thms. 1.2.1, 1.2.3].

**Theorem 6.9.**  $Y_n$  satisfies the large deviation principle on  $\mathcal{X}$  with rate function  $I$  if and only if  $Y_n$  satisfies the Laplace principle on  $\mathcal{X}$  with rate function  $I$ .

As we will see in section 10, where a general class of statistical mechanical models are studied, the Laplace principle gives a variational formula for the canonical free energy [Thm. 10.2(a)].

We next introduce the concept of exponential tightness, which will be used in the proof of Theorem 6.11.

**Definition 6.10.** The sequence  $Y_n$  is said to be exponentially tight if for every  $M < \infty$  there exists a compact subset  $K_M$  such that

$$\limsup_{n \rightarrow \infty} \frac{1}{n} \log P_n \{Y_n \in K_M^c\} \leq -M. \quad (6.2)$$

The next theorem shows that if  $Y_n$  is exponentially tight, then the large deviation upper bound for all compact sets implies the bound for all closed sets. This is a useful observation because one can often prove the bound for compact sets by covering them with a finite class of sets such as balls or halfspaces for which the proof is easier to obtain. We will see an example of this in the proof of Cramér's Theorem in the next section.

**Theorem 6.11.** Assume that one can prove the large deviation upper bound for any compact subset of  $\mathcal{X}$ . Then the large deviation upper bound is valid for any closed subset of  $\mathcal{X}$ .

**Proof.** We give the proof under the assumption that  $F$  is a closed set for which  $I(F) < \infty$ , omitting the minor modifications necessary to handle the case in which  $I(F) = \infty$ . Choose  $M < \infty$  such that  $M > I(F)$  and let  $K_M$  be the compact set satisfying (6.2) in the definition of exponential tightness. Since

$F \subset (F \cap K_M) \cup K_M^c$  and  $F \cap K_M$  is compact, it follows that

$$\begin{aligned}
& \limsup_{n \rightarrow \infty} \frac{1}{n} \log P_n\{Y_n \in F\} \\
& \leq \limsup_{n \rightarrow \infty} \frac{1}{n} \log P_n\{Y_n \in (F \cap K_M) \cup K_M^c\} \\
& \leq \limsup_{n \rightarrow \infty} \frac{1}{n} \log(P_n\{Y_n \in F \cap K_M\} + P_n\{Y_n \in K_M^c\}) \\
& = \max \left\{ \limsup_{n \rightarrow \infty} \frac{1}{n} \log P_n\{Y_n \in F \cap K_M\}, \limsup_{n \rightarrow \infty} \frac{1}{n} \log P_n\{Y_n \in K_M^c\} \right\} \\
& \leq \max\{-I(F \cap K_M), -M\} \\
& \leq \max\{-I(F), -M\} = -I(F).
\end{aligned}$$

This completes the proof. ■

We end this section by presenting three ways to obtain large deviation principles from existing large deviation principles. In the first theorem we show that a large deviation principle is preserved under continuous mappings. An application involving the relative entropy is given after the statement of Theorem 6.14.

**Theorem 6.12 (Contraction Principle).** *Assume that  $Y_n$  satisfies the large deviation principle on  $\mathcal{X}$  with rate function  $I$  and that  $\psi$  is a continuous function mapping  $\mathcal{X}$  into a complete, separable metric space  $\mathcal{Y}$ . Then  $\psi(Y_n)$  satisfies the large deviation principle on  $\mathcal{Y}$  with rate function*

$$J(y) = \inf\{I(x) : x \in \mathcal{X}, \psi(x) = y\}.$$

**Proof.** Since  $I$  maps  $\mathcal{X}$  into  $[0, \infty]$ ,  $J$  maps  $\mathcal{Y}$  into  $[0, \infty]$ . It is left as an exercise to show that since  $I$  has compact level sets in  $\mathcal{X}$ ,  $J$  has compact level sets in  $\mathcal{Y}$ . If  $F$  is a closed subset of  $\mathcal{Y}$ , then since  $\psi$  is continuous,  $\psi^{-1}(F)$  is a closed

subset of  $\mathcal{X}$ . Hence by the large deviation upper bound for  $Y_n$

$$\begin{aligned} \limsup_{n \rightarrow \infty} \frac{1}{n} \log P_n \{ \psi(Y_n) \in F \} &= \limsup_{n \rightarrow \infty} \frac{1}{n} \log P_n \{ Y_n \in \psi^{-1}(F) \} \\ &\leq - \inf_{x \in \psi^{-1}(F)} I(x) \\ &= - \inf_{y \in F} \{ \inf \{ I(x) : x \in \mathcal{X}, \psi(x) = y \} \} \\ &= - \inf_{y \in F} J(y) = -J(F). \end{aligned}$$

Similarly, if  $G$  is an open subset of  $\mathcal{Y}$ , then

$$\liminf_{n \rightarrow \infty} \frac{1}{n} \log P_n \{ \psi(Y_n) \in G \} \geq -J(G).$$

This completes the proof. ■

In the next theorem we show that a large deviation principle is preserved if the probability measures  $P_n$  are multiplied by suitable exponential factors and then normalized. This result will be applied in section 10 when we prove the large deviation principle for statistical mechanical models with respect to the canonical ensemble [Thm. 10.2].

**Theorem 6.13.** *Assume that with respect to the probability measures  $P_n$ ,  $Y_n$  satisfies the large deviation principle on  $\mathcal{X}$  with rate function  $I$ . Let  $\psi$  be a bounded, continuous function mapping  $\mathcal{X}$  into  $\mathbb{R}$ . For  $A \in \mathcal{F}_n$  we define new probability measures*

$$P_{n,\psi} \{ A \} = \frac{1}{\int_{\mathcal{X}} \exp[-n \psi(Y_n)] dP_n} \cdot \int_A \exp[-n \psi(Y_n)] dP_n.$$

*Then with respect to  $P_{n,\psi}$ ,  $Y_n$  satisfies the large deviation principle on  $\mathcal{X}$  with rate function*

$$I_\psi(x) = I(x) + \psi(x) - \inf_{y \in \mathcal{X}} \{ I(y) + \psi(y) \}.$$

**Proof.** Clearly  $I_\psi$  maps  $\mathcal{X}$  into  $[0, \infty]$ , and it is easily checked that  $I_\psi$  has compact level sets. We prove the theorem by showing that with respect to  $P_{n,\psi}$ ,



$Y_n$  satisfies the Laplace principle with rate function  $I_\psi$  and then invoke Theorem 6.9. Let  $f$  be any bounded, continuous function mapping  $\mathcal{X}$  into  $\mathbb{R}$ . Since  $f + \psi$  is bounded and continuous and since with respect to  $P_n$ ,  $Y_n$  satisfies the Laplace principle with rate function  $I$ , it follows that

$$\begin{aligned}
& \lim_{n \rightarrow \infty} \frac{1}{n} \log \int_{\Omega_n} \exp[n f(Y_n)] dP_{n,\psi} \\
&= \lim_{n \rightarrow \infty} \frac{1}{n} \log \int_{\Omega_n} \exp[n(f(Y_n) - \psi(Y_n))] dP_n \\
&\quad - \lim_{n \rightarrow \infty} \frac{1}{n} \log \int_{\Omega_n} \exp[-n \psi(Y_n)] dP_n \\
&= \sup_{x \in \mathcal{X}} \{f(x) - \psi(x) - I(x)\} - \sup_{y \in \mathcal{X}} \{-\psi(y) - I(y)\} \\
&= \sup_{x \in \mathcal{X}} \{f(x) - I_\psi(x)\}.
\end{aligned}$$

Thus with respect to  $P_{n,\psi}$ ,  $Y_n$  satisfies the Laplace principle with rate function  $I_\psi$ , as claimed. This completes the proof. ■

According to our next result, if random variables  $X_n$  are superexponentially close to random variables  $Y_n$  that satisfy the large deviation principle, then  $X_n$  satisfies the large deviation principle with the same rate function. A proof based on the equivalent Laplace principle is given in Theorem 1.3.3 in [31].

**Theorem 6.14.** *Assume that  $Y_n$  satisfies the large deviation principle on  $\mathcal{X}$  with rate function  $I$  and denote by  $m(\cdot, \cdot)$  the metric on  $\mathcal{X}$  (e.g.,  $m(x, y) = |x - y|$  if  $\mathcal{X} = \mathbb{R}$ ). Assume also that  $X_n$  is superexponentially close to  $Y_n$  in the following sense: for each  $\delta > 0$*

$$\limsup_{n \rightarrow \infty} \frac{1}{n} \log P_n \{m(Y_n, X_n) > \delta\} = -\infty. \tag{6.3}$$

*Then  $X_n$  satisfies the large deviation principle on  $\mathcal{X}$  with the same rate function  $I$ . The condition (6.3) is satisfied if*

$$\lim_{n \rightarrow \infty} \sup_{\omega \in \Omega_n} m(X_n(\omega), Y_n(\omega)) = 0.$$

A standard application of the contraction principle stated in Theorem 6.12 is to relate the rate functions in Sanov's Theorem and in Cramér's Theorem. For simplicity, we restrict to the case of a nondegenerate probability measure on  $\mathbb{R}$ ; much more general versions are available. For example, in [27, Thm. 5.2] it is shown to hold in the case of random variables taking values in a Banach space. Let  $\rho$  be a nondegenerate probability on  $\mathbb{R}$  having compact support  $K$  and  $Y_n$  an i.i.d. sequence of random variables having common distribution  $\rho$ . Since  $K$  is compact, the function  $\psi$  mapping  $\gamma \in \mathcal{P}(K)$  to  $\int_K x \gamma(dx)$  is bounded and continuous, and

$$\psi(L_n) = \int_K x L_n(dx) = \frac{1}{n} \int_K x \delta_{X_j}(dx) = \frac{1}{n} \sum_{j=1}^n X_j = \frac{S_n}{n}.$$

Since  $L_n$  satisfies the large deviation principle on  $\mathcal{P}(K)$  with rate function given by the relative entropy  $I_\rho$  [Thm. 6.7], the contraction principle implies that  $S_n/n$  satisfies the large deviation on  $\mathcal{K}$  with rate function

$$J(y) = \inf \left\{ I_\rho(\gamma) : \gamma \in \mathcal{P}(K), \int_K x \gamma(dx) = y \right\}$$

Since a rate function in a large deviation principle is unique [Thm. 6.2],  $J$  must equal the rate function  $I$  in Cramér's Theorem. We conclude that for all  $y \in \mathbb{R}$

$$I(y) = \sup_{t \in \mathbb{R}} \{ty - c(t)\} = \inf \left\{ I_\rho(\gamma) : \gamma \in \mathcal{P}(K), \int_K x \gamma(dx) = x \right\}. \quad (6.4)$$

We emphasize that in order to apply the contraction principle, one needs the hypothesis that  $\rho$  has compact support. It is satisfying to know that (6.4) is valid without this extra hypothesis [33, Thm. VIII.3.1].

This completes our discussion of the large deviation principle, the Laplace principle, and related general results. In the next section we prove Cramér's Theorem.

## 7 Cramér's Theorem

Cramér's Theorem is the large deviation principle for sums of i.i.d. random vectors taking values in  $\mathbb{R}^d$ . In this section Cramér's Theorem will be proved and several applications will be given.

Let  $\{X_j, j \in \mathbb{N}\}$  be a sequence of i.i.d. random vectors defined on a probability space  $(\Omega, \mathcal{F}, P)$  and taking values in  $\mathbb{R}^d$ . We are interested in the large deviation principle for the sample means  $S_n/n$ , where  $S_n = \sum_{j=1}^n X_j$ . The basic assumption is that the moment generating function  $E\{\exp\langle t, X_1 \rangle\}$  is finite for all  $t \in \mathbb{R}^d$ . We define for  $t \in \mathbb{R}^d$  the cumulant generating function

$$c(t) = \log E\{\exp\langle t, X_1 \rangle\},$$

which is finite, convex, and differentiable, and for  $x \in \mathbb{R}^d$  we define the Legendre-Fenchel transform

$$I(x) = \sup_{t \in \mathbb{R}^d} \{\langle t, x \rangle - c(t)\}.$$

The basic theory of convex functions and the Legendre-Fenchel transform is developed in chapter VI of [33]. Here are some relevant definitions. A function  $f$  mapping  $\mathbb{R}$  into  $\mathbb{R} \cup \{\infty\}$  is said to be convex if for all  $x$  and  $y$  in  $\mathbb{R}^d$  and all  $\lambda \in (0, 1)$

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

Such a function is said to be lower semicontinuous if whenever  $x_n \rightarrow x \in \mathbb{R}$ , we have  $\liminf_{n \rightarrow \infty} f(x_n) \geq f(x)$ . The convexity of  $c(t)$  is an immediate consequence of Hölder's inequality with  $p = 1/\lambda$  and  $q = 1/(1 - \lambda)$  [33, Prop. VII.1.1].

We consider Cramér's Theorem first the case  $d = 1$ . Let  $\alpha$  be a real number exceeding the mean value  $E\{X_1\}$ . Assuming that  $\rho$  has an absolutely continuous component and that certain other conditions hold, Cramér obtained in his 1938 paper [17] an asymptotic expansion for the probability  $P\{S_n/n \in [\alpha, \infty)\}$ , which implies the large deviation limit

$$\lim_{n \rightarrow \infty} \frac{1}{n} \log P\{S_n/n \in [\alpha, \infty)\} = -I(\alpha) = -I([\alpha, \infty)).$$

In the modern theory of large deviations the following generalization of this limit is known as Cramér's Theorem.

**Theorem 7.1 (Cramér's Theorem).** *Let  $\{X_j, j \in \mathbb{N}\}$  be a sequence of i.i.d. random vectors taking values in  $\mathbb{R}^d$  and satisfying  $E\{\exp\langle t, X_1 \rangle\} < \infty$  for all  $t \in \mathbb{R}^d$ . The following conclusions hold.*

(a) *The sequence of sample means  $S_n/n$  satisfies the large deviation principle on  $\mathbb{R}^d$  with rate function  $I(x) = \sup_{t \in \mathbb{R}^d} \{\langle t, x \rangle - c(t)\}$ .*

(b)  *$I$  is a convex, lower semicontinuous function on  $\mathbb{R}^d$ , and it attains its infimum of 0 at the unique point  $x_0 = E\{X_1\}$ .*

Infinite-dimensional generalizations of Cramér's Theorem have been proved by many authors, including [1] and [27, §5]. The book [20] presents Cramér's Theorem first in the setting of  $\mathbb{R}^d$  and then in the setting of a complete, separable metric space. At the end of this section we will derive from Cramér's Theorem the large deviation principle for the empirical vectors stated in Theorem 3.4. This is a special case of Sanov's Theorem 6.7. We will also indicate how to prove a general version of Sanov's Theorem from an infinite-dimensional version of Cramér's Theorem.

The properties of  $I$  stated in part (b) of Theorem 7.1 as well as other properties of this function related to Legendre-Fenchel duality are proved in [33, Thm. VII.5.5]. Before proving Cramér's Theorem, it is worthwhile to motivate the form of the rate function  $I$ . Assuming that the sequence  $S_n/n$  satisfies the large deviation principle on  $\mathbb{R}^d$  with some convex, lower semicontinuous rate function  $J$ , we will prove that  $J = I$ .

Since for each  $t \in \mathbb{R}^d$

$$\begin{aligned} c(t) &= \log E\{\exp\langle t, X_1 \rangle\} = \frac{1}{n} \log E\{\exp\langle t, S_n/n \rangle\} \\ &= \frac{1}{n} \log \int_{\mathbb{R}^d} \exp[n\langle t, x \rangle] P\{S_n/n \in dx\}, \end{aligned}$$

it follows that

$$c(t) = \lim_{n \rightarrow \infty} \frac{1}{n} \log \int_{\mathbb{R}^d} \exp\langle t, x \rangle P\{S_n/n \in dx\}.$$

We now use the hypothesis that  $S_n/n$  satisfies the large deviation principle on  $\mathbb{R}^d$  with some convex, lower semicontinuous rate function  $J$ . Although the function mapping  $x \mapsto \langle t, x \rangle$  is not bounded, a straightforward extension of Theorem 6.9 allows us to apply the Laplace principle to evaluate the last limit, yielding

$$c(t) = \sup_{x \in \mathbb{R}^d} \{\langle t, x \rangle - J(x)\}.$$

The assumed convexity and lower semicontinuity of  $J$  combined with Legendre-Fenchel duality now yields the desired formula; namely, for each  $x \in \mathbb{R}^d$

$$J(x) = \sup_{t \in \mathbb{R}^d} \{\langle t, x \rangle - c(t)\} = I(x).$$

Legendre-Fenchel duality is explained, for example, in [33, §VI.5]. This completes the motivation of the form of the rate function in Cramer's Theorem.

We now turn to the proof of Cramér's Theorem. The main tool used in the proof of the large deviation upper bound is Chebyshev's inequality, introduced by Chernoff in [10], while the main tool used in the proof of the large deviation lower bound is a change of measure, introduced by Cramér in his 1938 paper [17]. These same tools for proving the large deviation bounds continue to be used in modern developments of the theory.

*Proof of Theorem 7.1.* We first show that  $I$  is a rate function, then prove part (b) followed by the proofs of the large deviation upper bound and lower bound.

*$I$  is a rate function.* Since  $I$  is defined as a Legendre-Fenchel transform, it is automatically convex and lower semicontinuous. By part (a) of Proposition 6.4, the infimum of  $I$  over  $\mathbb{R}^d$  equals 0, and so  $I$  maps  $\mathbb{R}^d$  into the extended nonnegative real numbers  $[0, \infty]$ . We now consider a level set  $K_L = \{x \in \mathbb{R}^d : I(x) \leq L\}$ , where  $L$  is any nonnegative real number. This set is closed since  $I$  is lower semicontinuous. If  $x$  is in  $K_L$ , then for any  $t \in \mathbb{R}^d$

$$\langle t, x \rangle \leq c(t) + I(x) \leq c(t) + L.$$

Fix any positive number  $R$ . The finite, convex, continuous function  $c$  is bounded on the ball of radius  $R$  with center 0, and so there exists a number

$\Gamma < \infty$  such that

$$\sup_{\|t\| \leq R} \langle t, x \rangle = R\|x\| \leq \sup_{\|t\| \leq R} c(t) + L \leq \Gamma < \infty.$$

This implies that  $K_L$  is bounded and thus that the level sets of  $I$  are compact. The sketch of the proof that  $I$  is a rate function is complete.

*Part (b).* We have already remarked that  $I$  is convex and lower semicontinuous. Since  $I$  is a rate function,  $I$  attains its infimum of 0 at some point  $x_0 \in \mathbb{R}^d$  [Prop. 6.4(a)]. It is easy to show that if  $x_0 = E\{X_1\}$ , then  $I(x_0) = 0$ . Indeed, since  $c(t)$  is convex and differentiable, the infimum in the definition of

$$I(x_0) = \sup_{t \in \mathbb{R}^d} \{\langle t, x_0 \rangle - c(t)\}$$

is attained when  $t$  satisfies  $\nabla c(t) = x_0 = E\{X_1\}$ . Since  $\nabla c(0) = E\{X_1\}$ , it follows  $t = 0$ . For this choice of  $t$ , we have  $I(x_0) = \langle 0, x_0 \rangle - c(0) = 0$ , as claimed. The proof that  $x_0 = E\{X_1\}$  is the unique minimum point of  $I$  requires some additional ideas from convex analysis, which we omit [33, Thm. VII.5.5].

*Large deviation upper bound.* We first prove this bound in the case  $d = 1$ . Our aim is to prove that for any nonempty closed subset  $F$  of  $\mathbb{R}$

$$\limsup_{n \rightarrow \infty} \frac{1}{n} \log P\{S_n/n \in F\} \leq -I(F).$$

Let  $x_0 = E\{X_1\}$ . We first show this for the closed intervals  $[\alpha, \infty)$ , where  $\alpha > x_0$ . For any  $t > 0$  Chebyshev's inequality implies

$$\begin{aligned} P\{S_n/n \in [\alpha, \infty)\} &= P\{tS_n > nt\alpha\} \\ &\leq \exp[-nt\alpha] E\{\exp[tS_n]\} \\ &= \exp[-nt\alpha] \prod_{i=1}^n E\{\exp[X_i]\} \\ &= \exp[-nt\alpha] (E\{\exp[X_1]\})^n \\ &= \exp[-nt\alpha + n \log E\{\exp[X_1]\}] \\ &= \exp[-n(t\alpha - c(t))]. \end{aligned}$$

It follows that for any  $t > 0$

$$\limsup_{n \rightarrow \infty} \frac{1}{n} \log P\{S_n/n \in [\alpha, \infty)\} \leq -(t\alpha - c(t)),$$

and thus that

$$\limsup_{n \rightarrow \infty} \frac{1}{n} \log P\{S_n/n \in [\alpha, \infty)\} \leq -\sup_{t>0} \{t\alpha - c(t)\}. \quad (7.1)$$

The next lemma will allow us to rewrite the right-hand side of this inequality as in the statement of Cramér's Theorem.

**Lemma 7.2.** *If  $\alpha > x_0$ , then*

$$\sup_{t>0} \{t\alpha - c(t)\} = I(\alpha) = I([\alpha, \infty).$$

**Proof.** Since  $c(t)$  is continuous at  $t = 0$ ,

$$I(\alpha) = \sup_{t \in \mathbb{R}} \{t\alpha - c(t)\} = \sup_{t \neq 0} \{t\alpha - c(t)\}.$$

Since  $c(t)$  is differentiable and convex, we have  $c'(0) \geq c(t)/t$  for any  $t < 0$ . Therefore, for any  $t < 0$

$$t\alpha - c(t) = t(\alpha - c(t)/t) \leq t(\alpha - c'(0)) < 0 = 0 \cdot \alpha - c(0).$$

The second inequality holds since  $\alpha > x_0 = E\{X_1\} = c'(0)$  and  $t < 0$ . From this display we see that the supremum in the formula for  $I(\alpha)$  cannot occur for  $t < 0$ . It follows that

$$I(\alpha) = \sup_{t>0} \{t\alpha - c(t)\}.$$

$I(x)$  is nonnegative, convex function satisfying  $I(x_0) = 0$ . Thus  $I(x)$  is non-increasing for  $x \leq x_0$  and is nonincreasing for  $x \geq x_0$ . This means that  $I(\alpha) = \inf\{I(x) : x \geq \alpha\} = I([\alpha, \infty)$ . The proof of the lemma is complete. ■

Inequality (7.1) and the lemma imply that if  $F$  is the closed interval  $[\alpha, \infty)$ , then the large deviation upper bound holds:

$$\limsup_{n \rightarrow \infty} \frac{1}{n} \log P\{S_n/n \in F\} \leq -I(\alpha) = -I(F).$$

A similar proof yields the large deviation upper bound if  $F = (-\infty, \alpha]$  and  $\alpha < x_0$ .

Now let  $F$  be an arbitrary nonempty closed set. If  $x_0 \in F$ , then  $I(F)$  equals 0 and the large deviation upper bound holds automatically since  $\log P\{S_n/n \in F\}$  is always nonpositive. If  $x_0 \notin F$ , then let  $(\alpha_1, \alpha_2)$  be the largest open interval containing  $x_0$  and having empty intersection with  $F$ .  $F$  is a subset of  $(-\infty, \alpha_1] \cup [\alpha_2, \infty)$ , and by the first part of the proof

$$\begin{aligned}
& \limsup_{n \rightarrow \infty} \frac{1}{n} \log P\{S_n/n \in F\} \\
& \leq \limsup_{n \rightarrow \infty} \frac{1}{n} \log P\{S_n/n \in (-\infty, \alpha_1] \cup [\alpha_2, \infty)\} \\
& \leq \limsup_{n \rightarrow \infty} \frac{1}{n} \log(P\{S_n/n \in [-\infty, \alpha_1]\} + P\{S_n/n \in [\alpha_2, \infty)\}) \\
& = \max \left\{ \limsup_{n \rightarrow \infty} \frac{1}{n} \log P\{S_n/n \in [-\infty, \alpha_1]\}, \limsup_{n \rightarrow \infty} \frac{1}{n} \log P\{S_n/n \in [\alpha_2, \infty)\} \right\} \\
& \leq \max\{-I(\alpha_1), -I(\alpha_2)\} \\
& = -\min\{I(\alpha_1), I(\alpha_2)\}.
\end{aligned}$$

If  $\alpha_1 = -\infty$  or  $\alpha_2 = \infty$ , then the corresponding term is missing. From the monotonicity properties of  $I(x)$  on  $(-\infty, x_0]$  and on  $[x_0, \infty)$ ,  $I(F) = \min\{I(\alpha_1), I(\alpha_2)\}$ . Hence from the last display we conclude that

$$\limsup_{n \rightarrow \infty} \frac{1}{n} \log P\{S_n/n \in F\} \leq -I(F).$$

This completes the proof of the large deviation upper bound for  $d = 1$ .

We now prove the large deviation upper bound for  $d > 1$ . Using the hypothesis that  $c(t) < \infty$  for all  $t \in \mathbb{R}^d$ , it is straightforward to prove that  $S_n/n$  is exponentially tight and that the compact set  $K_M$  appearing in the definition 6.10 of exponential tightness can be taken to be the hypercube  $K_m = [-M, M]^d$ . The details are left to the reader. By Theorem 6.11 the upper bound will follow for any closed set if we can prove it for any compact set  $K$ . If  $I(K) = 0$ , then the upper bound holds automatically since  $\log P_n\{S_n/n \in K\}$  is always nonpositive. Details will now be given under the assumption that  $I(K) < \infty$ . The



minor modifications necessary to prove the upper bound when  $I(K) = \infty$  are omitted.

The technique of the proof exhibits a remarkable interplay among analysis, geometry, and probability and readily extends to the much more general setting of the Gärtner-Ellis Theorem, which we will consider in the next section [Thm. 8.1]. We start by picking  $\varepsilon > 0$  to satisfy  $\varepsilon < I(K)$  and by defining for  $t \in \mathbb{R}^d$  the open halfspace

$$H_t = \{x \in \mathbb{R}^d : \langle t, x \rangle - c(t) > I(K) - \varepsilon\};$$

if  $t = 0$ , then  $H_0 = \emptyset$  since  $I(K) - \varepsilon > 0$ . Since for all  $x \in K$  we have  $I(x) > I(K) - \varepsilon$ , it follows that

$$\begin{aligned} K &\subset \{x \in \mathbb{R}^d : I(x) > I(K) - \varepsilon\} \\ &= \left\{ x \in \mathbb{R}^d : \sup_{t \in \mathbb{R}^d} \{\langle t, x \rangle - c(t)\} > I(K) - \varepsilon \right\} \\ &= \bigcup_{t \in \mathbb{R}^d} \{x \in \mathbb{R}^d : \langle t, x \rangle - c(t) > I(K) - \varepsilon\} \\ &= \bigcup_{t \in \mathbb{R}^d} H_t. \end{aligned}$$

Since  $K$  is compact, there exists  $r \in \mathbb{N}$  and nonzero  $t_1, \dots, t_r \in \mathbb{R}^d$  such that  $K \subset \cup_{i=1}^r H_{t_i}$ . Thus by Chebyshev's inequality

$$\begin{aligned} P\{S_n/n \in K\} &\leq \sum_{i=1}^r P\{S_n/n \in H_{t_i}\} && (7.2) \\ &= \sum_{i=1}^r P\{\langle t_i, S_n \rangle > n[c(t_i) + I(K) - \varepsilon]\} \\ &\leq \sum_{i=1}^r \exp[-n(c(t_i) + I(K) - \varepsilon)] E\{\exp[\langle t_i, S_n \rangle]\} \\ &= \sum_{i=1}^r \exp[-n(c(t_i) + I(K) - \varepsilon)] \exp[n c(t_i)] \\ &= r \exp[-n[I(K) - \varepsilon]], \end{aligned}$$

from which it follows that

$$\limsup_{n \rightarrow \infty} \frac{1}{n} \log P\{S_n/n \in K\} \leq -I(K) + \varepsilon.$$

Sending  $\varepsilon \rightarrow 0$  completes the proof of the large deviation upper bound for  $d > 1$ . This argument generalizes the proof for  $d = 1$ , in which we covered an arbitrary closed set  $F$  by the intervals  $(-\infty, \alpha_1] \cup [\alpha_2, \infty)$ .

*Proof of large deviation lower bound.* In contrast to the large deviation upper bound, which is proved by a global estimate involving Chebyshev's inequality, the large deviation lower bound is proved by a local estimate, the heart of which involves a change of measure. The proof is somewhat more technical than the proof of the large deviation upper bound. We denote the common distribution of the random vectors  $X_j$  by

$$\rho(dx) = P\{X_j \in dx\}.$$

In general

$$c(t) = \log E\{\exp\langle t, X_1 \rangle\} = \log \int_{\mathbb{R}^d} \exp\langle t, x \rangle \rho(dx)$$

is a finite, convex, differentiable function on  $\mathbb{R}^d$ . We first prove the lower bound under the highly restrictive assumption that the support of  $\rho$  is all of  $\mathbb{R}^d$  or more generally that the smallest convex set containing the support of  $\rho$  is all of  $\mathbb{R}^d$ . In this case, for each  $z \in \mathbb{R}^d$  there exists  $t \in \mathbb{R}^d$  such that  $\nabla c(t) = z$  [33, Thms. VIII.3.3, VIII.4.3].

For  $z \in \mathbb{R}^d$  and  $\varepsilon > 0$ , we denote by  $B(z, \varepsilon)$  the open ball with center  $z$  and radius  $\varepsilon$ . Let  $G$  be an open subset of  $\mathbb{R}^d$ . Then for any point  $z_0 \in G$  there exists  $\varepsilon > 0$  such that  $B(z_0, \varepsilon) \subset G$ , and so

$$P\{S_n/n \in G\} \geq P\{S_n/n \in B(z_0, \varepsilon)\} = \int_{\{\sum_{j=1}^n x_j/n \in B(z_0, \varepsilon)\}} \prod_{j=1}^n \rho(dx_j).$$

We first assume that  $G$  contains the point  $\int_{\mathbb{R}^d} x \rho(dx) = E\{X_1\}$  and let  $z_0 = \int_{\mathbb{R}^d} x \rho(dx)$ . In this case the weak law of large numbers implies that

$$\lim_{n \rightarrow \infty} \int_{\{\sum_{j=1}^n x_j/n \in B(z_0, \varepsilon)\}} \prod_{j=1}^n \rho(dx_j) = 1.$$

Since  $I(z_0) = 0$ , we obtain the large deviation lower bound:

$$\liminf_{n \rightarrow \infty} \frac{1}{n} \log P\{S_n/n \in G\} \geq 0 = -I(z_0) = -I(G).$$

Of course, in general  $G$  does not contain the point  $\int_{\mathbb{R}^d} x \rho(dx)$ , and the argument in the preceding paragraph breaks down. In this case we let  $z_0$  be an arbitrary point in  $G$  and introduce a change of measure, replacing  $\rho$  by a new measure  $\rho_{t_0}$  whose mean equals  $z_0$ . The exponential price that must be paid for introducing this new measure is of the order of  $\exp[-nI(z_0)]$ . Putting the various estimates together will yield the desired large deviation lower bound.

Given  $z_0 \in G$ , we choose  $t_0 \in \mathbb{R}^d$  such that  $\nabla c(t_0) = z_0$ . We then introduce the change of measure given by the exponential family

$$\rho_{t_0}(dx) = \frac{1}{\int_{\mathbb{R}^d} e^{\langle t_0, x \rangle} \rho(dx)} \cdot e^{\langle t_0, x \rangle} \rho(dx) = \frac{1}{e^{c(t_0)}} \cdot e^{\langle t_0, x \rangle} \rho(dx).$$

Similar exponential families arise in Theorems 4.1 and 5.1. By the definition of  $c(t_0)$ ,  $\rho_{t_0}$  is a probability measure, and the mean of  $\rho_{t_0}$  is  $z_0$ . Indeed

$$\int_{\mathbb{R}^d} x \rho_{t_0}(dx) = \frac{1}{e^{c(t_0)}} \cdot \int_{\mathbb{R}^d} x e^{\langle t_0, x \rangle} \rho(dx) = \nabla c(t_0) = z_0.$$

Furthermore, since  $c(t)$  is convex,

$$I(z_0) = \sup_{t \in \mathbb{R}^d} \{\langle t, z_0 \rangle - c(t)\} = \langle t_0, z_0 \rangle - c(t_0).$$

We thus obtain the lower bound

$$\begin{aligned}
& P_n\{S_n/n \in G\} \\
& \geq P_n\{S_n/n \in B(z_0, \varepsilon)\} \\
& = \int_{\{\sum_{j=1}^n x_j/n \in B(z_0, \varepsilon)\}} \prod_{j=1}^n \rho(dx_j) \\
& = \int_{\{\sum_{j=1}^n x_j/n \in B(z_0, \varepsilon)\}} \left( \prod_{j=1}^n \frac{d\rho}{d\rho_{t_0}}(x_j) \right) \prod_{j=1}^n \rho_{t_0}(dx_j) \\
& = \int_{\{\sum_{j=1}^n x_j/n \in B(z_0, \varepsilon)\}} \exp\left[-n \left( \langle t_0, \sum_{j=1}^n x_j/n \rangle - c(t_0) \right)\right] \prod_{j=1}^n \rho_{t_0}(dx_j) \\
& \geq \exp[-n(\langle t_0, z_0 \rangle - c(t_0)) - n\|t_0\|\varepsilon] \int_{\{\sum_{j=1}^n x_j/n \in B(z_0, \varepsilon)\}} \prod_{j=1}^n \rho_{t_0}(dx_j) \\
& = \exp[-nI(z_0) - n\|t_0\|\varepsilon] \int_{\{\sum_{j=1}^n x_j/n \in B(z_0, \varepsilon)\}} \prod_{j=1}^n \rho_{t_0}(dx_j).
\end{aligned}$$

Since the mean of the probability measure  $\rho_{t_0}$  equals  $z_0$ , the weak law of large numbers for i.i.d. random vectors with common distribution  $\rho_{t_0}$  implies that

$$\lim_{n \rightarrow \infty} \int_{\{\sum_{j=1}^n x_j/n \in B(z_0, \varepsilon)\}} \prod_{j=1}^n \rho_{t_0}(dx_j) = 1.$$

Hence it follows that

$$\liminf_{n \rightarrow \infty} \frac{1}{n} \log P\{S_n/n \in G\} \geq -I(z_0) - \|t_0\|\varepsilon.$$

We now send  $\varepsilon \rightarrow 0$ , and since  $z_0$  is an arbitrary point in  $G$ , we can replace  $-I(z_0)$  by  $-\inf_{z_0 \in G} I(z_0) = -I(G)$ . This completes the proof of the large deviation lower bound when the support of  $\rho$  is all of  $\mathbb{R}^d$  or more generally when the smallest convex set containing the support of  $\rho$  is all of  $\mathbb{R}^d$ .

When this hypothesis does not hold, then the range of  $\nabla c(t)$  is no longer all of  $\mathbb{R}^d$ , and the argument just given breaks down. To handle the case of general  $\rho$ , we find a set  $A$  with the properties that  $I(G) = I(G \cap A)$  and that  $A$  is a subset of the range of  $\nabla c(t)$ . If we can do this, then the proof just given, specialized

to arbitrary  $z_0 \in G \cap A$ , yields

$$\liminf_{n \rightarrow \infty} \frac{1}{n} \log P\{S_n/n \in G\} \geq -I(z_0).$$

Since  $z_0$  is an arbitrary point in  $G \cap A$ , we can replace  $-I(z_0)$  by  $-\inf_{z_0 \in G \cap A} I(z_0) = -I(G \cap A) = -I(G)$ , and we are done.

By definition, the domain of  $I$ ,  $\text{dom } I$ , is the set of  $x \in \mathbb{R}^d$  for which  $I(x) < \infty$ . The relative interior of  $\text{dom } I$ , denoted by  $\text{ri}(\text{dom } I)$ , is defined as the interior of  $\text{dom } I$  when considered as a subset of the smallest affine set that contains  $\text{dom } I$ . Clearly, if the smallest affine set that contains  $\text{dom } I$  is  $\mathbb{R}^d$ , then the relative interior of  $\text{dom } I$  equals the interior of  $\text{dom } I$ . This is the case if, for example,  $d = 1$  and  $\text{dom } I$  is a nonempty interval.

Using several properties of convex sets and convex functions, we will show that the desired set  $A$  equals  $\text{ri}(\text{dom } I)$ . In order to see this, we first note that since  $I(x)$  equals  $\infty$  for  $x \notin \text{dom } I$ ,  $I(G)$  equals  $\infty$  if  $G \cap \text{dom } I$  is empty. In this case the large deviation lower bound is valid. If  $G \cap \text{dom } I$  is nonempty, then  $I(G)$  equals  $I(G \cap \text{dom } I)$ . The set  $G \cap \text{ri}(\text{dom } I)$  is also nonempty [71, p. 46], and by the continuity property of  $I$  expressed in [33, Thm. VI.3.2]

$$I(G) = I(G \cap \text{dom } I) = I(G \cap \text{ri}(\text{dom } I)).$$

This is the first required property of  $A = \text{ri}(\text{dom } I)$ . The second desired property of this set — namely, that  $\text{ri}(\text{dom } I)$  is a subset of the range of  $\nabla c(t)$  — is a consequence of [33, Thm. VI.5.7], which is based on duality properties involving  $c(t)$  and its Legendre-Fenchel transform  $I(x)$  [71]. This completes the proof of the large deviation lower bound and thus the proof of Cramér's Theorem. ■

We now apply Cramér's Theorem to derive the special case of Sanov's Theorem given in Theorem 3.4; the latter states the large deviation principle for the empirical vectors of i.i.d. random variables having a finite state space. Let  $\alpha \geq 2$  be an integer,  $y_1 < y_2 < \dots < y_\alpha$  a set of  $\alpha$  real numbers,  $\rho_1, \rho_2, \dots, \rho_\alpha$  a set of  $\alpha$  positive real numbers summing to 1, and  $\{X_j, j \in \mathbb{N}\}$  a sequence of i.i.d. random variables defined on a probability space  $(\Omega, \mathcal{F}, P)$ , taking values in  $\Lambda = \{y_1, y_2, \dots, y_\alpha\}$ , and having distribution  $\rho = \sum_{k=1}^{\alpha} \rho_k \delta_{y_k}$ . In Theorem

3.4 we take  $\{X_j, j = 1, \dots, n\}$  to be the coordinate functions on  $\Lambda^n$  and impose on this space the product measure  $P_n$  with one dimensional marginals  $\rho$ , but there is no need to restrict to this case. For  $\omega \in \Omega$  and  $y \in \Lambda$  we consider

$$L_n(y) = L_n(\omega, y) = \frac{1}{n} \sum_{j=1}^n \delta_{X_j(\omega)}\{y\}$$

and the empirical vector

$$L_n = L_n(\omega) = (L_n(\omega, y_1), \dots, L_n(\omega, y_\alpha)) = \frac{1}{n} \sum_{j=1}^n (\delta_{X_j(\omega)}\{y_1\}, \dots, \delta_{X_j(\omega)}\{y_\alpha\}).$$

$L_n$  takes values in the set of probability vectors

$$\mathcal{P}_\alpha = \left\{ \gamma \in \mathbb{R}^\alpha : \gamma = (\gamma_1, \gamma_2, \dots, \gamma_\alpha) \geq 0, \sum_{k=1}^\alpha \gamma_k = 1 \right\}.$$

Since  $L_n$  equals the sample mean of the i.i.d. random variables

$$Y_j(\omega) = (\delta_{X_j(\omega)}\{y_1\}, \dots, \delta_{X_j(\omega)}\{y_\alpha\}),$$

the large deviation principle for  $L_n$  follows from Cramér's Theorem. For  $\gamma \in \mathbb{R}^\alpha$  the rate function is given by

$$I(\gamma) = \sup_{t \in \mathbb{R}^\alpha} \{ \langle \gamma, t \rangle - c(t) \}, \text{ where } c(t) = E\{\exp\langle t, Y_1 \rangle\} = \log \left( \sum_{k=1}^\alpha e^{t_k} \rho_k \right).$$

In the next proposition we show that for  $\gamma \in \mathcal{P}_\alpha$ ,  $I(\gamma)$  equals the relative entropy  $I_\rho(\gamma)$  and that for  $\gamma \in \mathbb{R}^\alpha \setminus \mathcal{P}_\alpha$ ,  $I_\rho(\gamma)$  equals  $\infty$ .

**Proposition 7.3.** *For  $\gamma \in \mathcal{P}_\alpha$  we define the relative entropy*

$$I_\rho(\gamma) = \sum_{k=1}^\alpha \gamma_k \log \frac{\gamma_k}{\rho_k}.$$

*Then*

$$I(\gamma) = \sup_{t \in \mathbb{R}^\alpha} \{ \langle \gamma, t \rangle - \log(\sum_{k=1}^\alpha e^{t_k} \rho_k) \} = \begin{cases} I_\rho(\gamma) & \text{for } \gamma \in \mathcal{P}_\alpha \\ \infty & \text{for } \gamma \in \mathbb{R}^\alpha \setminus \mathcal{P}_\alpha. \end{cases}$$

**Sketch of Proof.** Let  $G$  be any open set having empty intersection with  $\mathcal{P}_\alpha$ . Since  $P\{L_n \in G\} = 0$ , the large deviation lower bound implies that

$$-\infty = \liminf_{n \rightarrow \infty} \frac{1}{n} \log P\{L_n \in G\} \geq -I(G).$$

It follows that  $I(\gamma) = \infty$  for  $\gamma \in \mathbb{R}^\alpha \setminus \mathcal{P}_\alpha$ . The exercise of proving this directly from the definition of  $I(\gamma)$  as a Legendre-Fenchel transform is left to the reader.

Defining  $\mathcal{P}_\alpha^\circ$  to be the set of  $\gamma \in \mathcal{P}_\alpha$  having all positive components, we next prove that  $I(\gamma) = I_\rho(\gamma)$  for  $\gamma \in \mathcal{P}_\alpha^\circ$ . Let  $\mathbb{R}_+^\alpha$  denote the positive orthant of  $\mathbb{R}^\alpha$ . Since  $-\log$  is strictly convex on  $(0, \infty)$ , Jensen's inequality implies that for any  $\gamma \in \mathcal{P}_\alpha^\circ$

$$\sup_{s \in \mathbb{R}_+^\alpha} \left\{ \sum_{k=1}^{\alpha} \gamma_k \log s_k - \log \sum_{k=1}^{\alpha} \gamma_k s_k \right\} \leq 0$$

with equality if and only if  $s_k = \text{const.}$  For  $\gamma \in \mathcal{P}_\alpha^\circ$ , as  $t$  runs through  $\mathbb{R}^\alpha$ , the vector  $s$  having components  $s_k = e^{t_k} \rho_k / \gamma_k$  runs through  $\mathbb{R}_+^\alpha$ . Hence

$$\begin{aligned} I(\gamma) &= \sup_{t \in \mathbb{R}^\alpha} \left\{ \langle \gamma, t \rangle - \log \left( \sum_{k=1}^{\alpha} e^{t_k} \rho_k \right) \right\} \\ &= \sup_{s \in \mathbb{R}_+^\alpha} \left\{ \sum_{k=1}^{\alpha} \gamma_k \log \frac{\gamma_k s_k}{\rho_k} - \log \left( \sum_{k=1}^{\alpha} \gamma_k s_k \right) \right\} \\ &= \sum_{k=1}^{\alpha} \gamma_k \log \frac{\gamma_k}{\rho_k} + \sup_{s \in \mathbb{R}_+^\alpha} \left\{ \sum_{k=1}^{\alpha} \gamma_k \log s_k - \log \sum_{k=1}^{\alpha} \gamma_k s_k \right\} \\ &= \sum_{k=1}^{\alpha} \gamma_k \log \frac{\gamma_k}{\rho_k} \\ &= I_\rho(\gamma). \end{aligned}$$

This completes the proof that  $I(\gamma) = I_\rho(\gamma)$  for  $\gamma \in \mathcal{P}_\alpha^\circ$ . In order to prove this equality for all  $\gamma \in \mathcal{P}_\alpha$ , we use the continuity of  $I_\rho$  on  $\mathcal{P}_\alpha$  and the continuity property of  $I$  on  $\mathcal{P}_\alpha$  stated in [33, Thm. VI.3.2]. The proof of the proposition is complete. ■

In Theorem 5.2 in [27] the following infinite dimensional version of Cramér's Theorem is proved.

**Theorem 7.4.** *Let  $\mathcal{X}$  be a Banach space with dual space  $\mathcal{X}^*$  and  $\{X_j, j \in \mathbb{N}\}$  a sequence of i.i.d. random vectors taking values in  $\mathcal{X}$  and having common distribution  $\rho$ . Assume that  $E\{\exp(t\|X_1\|)\} < \infty$  for every  $t > 0$ . Then the sequence of sample means  $S_n/n$  satisfies the large deviation principle on  $\mathcal{X}$  with rate function*

$$I(x) = \sup_{\theta \in \mathcal{X}^*} \left\{ \langle \theta, x \rangle - \log \int_{\mathcal{X}} \exp \langle \theta, y \rangle \rho(dy) \right\}.$$

*The rate function  $I$  is convex and lower semicontinuous and attains its infimum of 0 at the unique point  $x_0 = E\{X_1\} = \int_{\mathcal{X}} x \rho(dx)$ .*

We now return to the setting of Sanov's Theorem, considering the empirical measures  $L_n$  of a sequence  $\{X_j, j \in \mathbb{N}\}$  of i.i.d. random variables taking values in  $\mathbb{R}^d$  [Thm. 6.7] and more generally in a complete, separable metric space  $\mathcal{X}$ . Let  $\rho$  denote the common distribution of  $X_j$ . Then

$$L_n = \frac{1}{n} \sum_{j=1}^n \delta_{X_j}$$

takes values in the complete, separable metric space  $\mathcal{P}(\mathcal{X})$  of probability measures on  $\mathcal{X}$ . Since  $L_n$  is the sample mean of the i.i.d. random variables  $\delta_{X_j}$ , it is reasonable to conjecture that Sanov's Theorem can be derived as a consequence of a suitable infinite-dimensional version of Cramér's Theorem. While Theorem 7.4 cannot be applied because  $\mathcal{P}(\mathcal{X})$  is not a Banach space, the derivation of Sanov's Theorem from a suitable infinite-dimensional version of Cramér's Theorem is carried out in [21, Thm. 3.2.17]. This reference first proves that  $L_n$  satisfies the large deviation principle on  $\mathcal{P}(\mathcal{X})$  with rate function  $I(\gamma)$  given by the Legendre-Fenchel transform

$$I(\gamma) = \sup_{f \in \mathcal{C}(\mathcal{X})} \left\{ \int_{\mathcal{X}} f d\gamma - \log \int_{\mathcal{X}} e^f d\rho \right\},$$

where  $\mathcal{C}(\mathcal{X})$  denotes the set of bounded, continuous functions mapping  $\mathcal{X}$  into  $\mathbb{R}$ . The proof of Sanov's Theorem is completed by showing that  $I(\gamma)$  equals the relative entropy  $I_\rho(\gamma)$ . A special case of this identification of the relative



entropy with a Legendre-Fenchel transform is given in Proposition 7.3. An independent derivation of Sanov's Theorem for i.i.d. random vectors taking values in a complete, separable space is given in [31, Ch. 2], which applies ideas from stochastic optimal control theory.

In the next section we present a generalization of Cramér's Theorem that does not require the underlying random variables to be independent. Both in Cramér's Theorem and in this generalization the rate functions are defined by Legendre-Fenchel transforms and so are always convex. This convexity is not a general feature. Indeed, at the end of the next section we present two examples of large deviation principles in which the rate function is not convex.

## 8 Gärtner-Ellis Theorem

For each  $n \in \mathbb{N}$  let  $(\Omega_n, \mathcal{F}_n, P_n)$  be a probability space and let  $Y_n$  be a random vector mapping  $\Omega_n$  into  $\mathbb{R}^d$ . In 1977 Gärtner proved an important generalization of Cramer's Theorem, assuming only that the limit

$$c(t) = \lim_{n \rightarrow \infty} \frac{1}{n} \log E^{P_n} \{ \exp[n \langle t, Y_n \rangle] \} \quad (8.1)$$

exists and is finite for every  $t \in \mathbb{R}^d$  and that  $c(t)$  is a differentiable function of  $t \in \mathbb{R}^d$  [50]. Gärtner's result is that  $Y_n$  satisfies the large deviation principle on  $\mathbb{R}^d$  with rate function equal to the Legendre-Fenchel transform

$$I(x) = \sup_{t \in \mathbb{R}^d} \{ \langle t, x \rangle - c(t) \}.$$

Using ideas from convex analysis, I generalized Gärtner's result by relaxing the condition that  $c(t)$  exist and be finite for every  $t \in \mathbb{R}^d$  [34]. The theorem is now known in the literature as the Gärtner-Ellis Theorem [20, §2.3, §2.5].

Gärtner's result contains Cramér's Theorem as a special case. In order to see this, let  $Y_n$  equal  $\sum_{j=1}^n X_j/n$ , where  $X_j$  is a sequence of i.i.d. random vectors satisfying  $E\{\exp\langle t, X_1 \rangle\} < \infty$  for every  $t \in \mathbb{R}^d$ . In this case the limit  $c(t)$  in (8.1) equals  $\log E\{\exp\langle t, X_1 \rangle\}$ , which is a differentiable function of  $t \in \mathbb{R}^d$ . The corresponding rate function is the same as in Cramer's Theorem.

We next state the Gärtner-Ellis Theorem under the hypotheses of [50] and in a form that is different from but equivalent to Gärtner's result in that paper. This is followed by comments on the generalization proved in [34]. In this theorem the differentiability of  $c(t)$  for all  $t \in \mathbb{R}^d$  is a sufficient condition for the large deviation lower bound; the large deviation upper bound is always valid. However, as we mention just before Example 8.3 in the context of the Ising model in statistical mechanics, the differentiability of  $c(t)$  is not a necessary condition for the validity of the lower bound.

**Theorem 8.1.** *For each  $n \in \mathbb{N}$  let  $(\Omega_n, \mathcal{F}_n, P_n)$  be a probability space and let  $Y_n$  be a random vector mapping  $\Omega_n$  into  $\mathbb{R}^d$ . We assume that the limit*

$$c(t) = \lim_{n \rightarrow \infty} \frac{1}{n} \log E^{P_n} \{ \exp[n \langle t, Y_n \rangle] \}$$

exists and is finite for every  $t \in \mathbb{R}^d$ . For  $x \in \mathbb{R}^d$  we define

$$I(x) = \sup_{t \in \mathbb{R}^d} \{\langle t, x \rangle - c(t)\}.$$

The following conclusions hold.

(a)  $I$  is a rate function. Furthermore,  $I$  is convex and lower semicontinuous.

(b) The large deviation upper bound is valid. Namely, for every closed subset  $F$  of  $\mathbb{R}^d$

$$\limsup_{n \rightarrow \infty} \frac{1}{a_n} \log P_n\{Y_n \in F\} \leq -I(F).$$

(c) Assume in addition that  $c(t)$  is differentiable for all  $t \in \mathbb{R}^d$ . Then the large deviation lower bound is valid. Namely, for every open subset  $G$  of  $\mathbb{R}^d$

$$\liminf_{n \rightarrow \infty} \frac{1}{a_n} \log P_n\{Y_n \in G\} \geq -I(G).$$

Hence, if  $c(t)$  is differentiable for all  $t \in \mathbb{R}^d$ , then  $Y_n$  satisfies the large deviation principle on  $\mathbb{R}^d$  with rate function  $I$ .

The theorem is proved by suitably generalizing the proof of Cramér's Theorem (see [33, Ch. 7]). In the case of the large deviation upper bound, the generalization is easy to see. As in the proof of Cramér's Theorem, the assumption that the limit function  $c(t)$  is finite for every  $t \in \mathbb{R}^d$  implies that  $Y_n$  is exponentially tight. Hence by Theorem 6.11, the upper bound will follow for any closed set if we can prove it for any compact set  $K$ . If  $I(K) = 0$ , then the upper bound holds automatically since  $\log P_n\{Y_n \in K\}$  is always nonpositive. In order to handle the case when  $I(K) < \infty$ , we argue as in the proof of Cramér's Theorem. Given  $\varepsilon > 0$  satisfying  $\varepsilon < I(K)$ , there exists  $r \in \mathbb{N}$  and nonzero  $t_1, \dots, t_r \in \mathbb{R}^d$  such that  $K \subset \cup_{i=1}^r H_{t_i}$ , where  $H_{t_i}$  denotes the open halfspace

$$H_{t_i} = \{x \in \mathbb{R}^d : \langle t_i, x \rangle - c(t_i) > I(K) - \varepsilon\}.$$

As in the display (7.2), Chebyshev's inequality yields

$$\begin{aligned}
P\{Y_n \in K\} &\leq \sum_{i=1}^r P\{Y_n \in H_{t_i}\} \\
&= \sum_{i=1}^r P\{n\langle t_i, Y_n \rangle > n[c(t_i) + I(K) - \varepsilon]\} \\
&\leq \sum_{i=1}^r \exp[-n(c(t_i) + I(K) - \varepsilon)] E\{\exp[n\langle t_i, Y_n \rangle]\} \\
&= \sum_{i=1}^r \exp[-n(c(t_i) + I(K) - \varepsilon)] \exp[n c_n(t_i)],
\end{aligned}$$

where

$$c_n(t) = \frac{1}{n} \log E^{P_n}\{\exp[n\langle t, Y_n \rangle]\}.$$

Since  $c_n(t_i) \rightarrow c(t_i)$ , there exists  $N \in \mathbb{N}$  such that  $c_n(t_i) < c(t_i) + \varepsilon$  for all  $n \geq N$  and all  $i = 1, \dots, r$ . Thus for all  $n \geq N$

$$P_n\{Y_n \in K\} \leq r \exp[-n(I(K) - 2\varepsilon)],$$

from which it follows that

$$\limsup_{n \rightarrow \infty} \frac{1}{n} \log P\{Y_n \in K\} \leq -I(K) + 2\varepsilon.$$

Sending  $\varepsilon \rightarrow 0$ , we complete the proof of the large deviation upper bound in the Gärtner-Ellis Theorem when  $I(K) < \infty$ . The minor modifications necessary to prove the upper bound when  $I(K) = \infty$  are omitted.

The proof of the large deviation lower bound requires a new idea. We recall that the proof of the lower bound in Cramér's Theorem invoked the weak law of large numbers with respect to the change of measure given by the product measure with one-dimensional marginals  $\rho_{t_0}$ . In the proof of the lower bound in the Gärtner-Ellis Theorem again one uses a change of measure, but the weak law of large numbers with respect to a product measure is not available. The innovation is to replace the weak law of large numbers by an order-1 estimate based on the large deviation upper bound in the Gärtner-Ellis Theorem.

In the extension of Gärtner's result proved in [34], it is assumed that for all  $t \in \mathbb{R}^d$ ,  $c(t)$  exists as an extended real number in  $(-\infty, \infty]$ . Then the large deviation upper bound as stated in part (b) of Theorem 8.1 is valid with rate function

$$I(x) = \sup_{t \in \mathbb{R}^d} \{ \langle t, x \rangle - c(t) \}.$$

We denote by  $\mathcal{D}$  the set of  $t \in \mathbb{R}^d$  for which  $c(t)$  is finite. If  $c$  is differentiable on the interior of  $\mathcal{D}$ , then  $c$  is called steep if  $\|\nabla c(t_n)\| \rightarrow \infty$  for any sequence  $t_n$  in the interior of  $\mathcal{D}$  that converges to a boundary point of  $\mathcal{D}$ . For example, if  $c$  is lower semicontinuous,  $\mathcal{D}$  is open, and  $c$  is differentiable on  $\mathcal{D}$ , then  $c$  is steep. In the extension of Gärtner's result, it is proved that if  $c$  is differentiable on the interior of  $\mathcal{D}$  and is steep, then the large deviation lower bound as stated in part (c) of Theorem 8.1 is valid with rate function

$$I(x) = \sup_{t \in \mathbb{R}^d} \{ \langle t, x \rangle - c(t) \}.$$

We next give an application of the Gärtner-Ellis Theorem to finite-state Markov chains. Let  $\alpha \geq 2$  be an integer,  $y_1 < y_2 < \dots < y_\alpha$  be a set of  $\alpha$  real numbers, and  $\{X_j, j \in \mathbb{N}\}$  a Markov chain taking values in  $\Lambda = \{y_1, y_2, \dots, y_\alpha\}$ . We denote by  $\rho \in \mathbb{R}^\alpha$  the initial distribution  $\rho_i = P\{X_1 = y_i\}$  and by  $\pi(i, j)$  the transition probabilities  $P\{X_{j+1} = y_j | X_j = y_i\}$  for  $1 \leq i, j \leq \alpha$ . Under the assumption that the matrix  $\pi = \{\pi(i, j)\}$  is irreducible and aperiodic, we have the following two-part theorem. Part (a) is the large deviation principle for the sample means  $Y_n = \sum_{j=1}^{\alpha} X_j/n$ , and part (b) is the large deviation principle for the empirical vectors  $L_n = \sum_{j=1}^{\alpha} \delta_{X_j}/n$ . The hypothesis that  $\pi$  is irreducible holds if, for example,  $\pi$  is a positive matrix. Part (b) is a special case of a result proved in [26].

**Theorem 8.2.** *We assume that the transition probability matrix  $\pi$  of the Markov chain  $X_j$  is aperiodic and irreducible. The following conclusions hold.*

(a) *For  $t \in \mathbb{R}$  let  $B(t)$  be the matrix with entries  $[B(t)]_{i,j} = \exp(tx_i)\pi(i, j)$ . Then for all  $t \in \mathbb{R}$ ,  $B(t)$  has a unique largest positive eigenvalue  $\lambda(t)$  which is differentiable for all  $t \in \mathbb{R}$ . Furthermore, for any choice of the initial distribution  $\rho$ , the sample means  $Y_n$  satisfy the large deviation principle on  $\mathbb{R}$  with rate*

function

$$I(x) = \sup_{t \in \mathbb{R}} \{tx - \log \lambda(t)\}.$$

(b) We denote by  $\mathcal{P}_\alpha$  the set of probability vectors in  $\mathbb{R}^\alpha$ . Then for any choice of the initial distribution  $\rho$ , the empirical vectors  $L_n$  satisfy the large deviation principle on  $\mathcal{P}_\alpha$  with rate function

$$I_\pi(\gamma) = - \inf_{u > 0} \sum_{i=1}^{\alpha} \gamma_i \log \frac{(\pi u)_i}{u_i}.$$

In this formula  $u$  is any positive vector in  $\mathbb{R}^\alpha$ , and  $(\pi u)_i = \sum_{j=1}^{\alpha} \pi(i,j)u_j$ .

**Sketch of Proof.** (a) That  $B(t)$  has a unique largest positive eigenvalue  $\lambda(t)$  is a consequence of the Perron-Frobenius Theorem [72, Thm. 1.1]. The differentiability of  $\lambda(t)$  is a consequence of the implicit function theorem and the fact that  $\lambda(t)$  is a simple root of the characteristic equation for  $B(t)$ . For  $t \in \mathbb{R}$  we calculate

$$\begin{aligned} c(t) &= \lim_{n \rightarrow \infty} \frac{1}{n} \log E\{\exp[ntY_n]\} \\ &= \lim_{n \rightarrow \infty} \frac{1}{n} \log E\{\exp\langle t \sum_{j=1}^n X_j \rangle\} \\ &= \lim_{n \rightarrow \infty} \frac{1}{n} \log \sum_{i_1, i_2, \dots, i_n=1}^{\alpha} \exp(t \sum_{j=1}^n x_{i_j}) \rho_{i_1} \pi(i_1, i_2) \pi(i_2, i_3) \pi(i_{n-1}, i_n) \\ &= \lim_{n \rightarrow \infty} \frac{1}{n} \log \sum_{i_1, i_2, \dots, i_n=1}^{\alpha} \rho_{i_1} (B(t))_{i_1, i_2} (B(t))_{i_2, i_3} \dots (B(t))_{i_{n-1}, i_n} e^{tx_{i_n}} \\ &= \lim_{n \rightarrow \infty} \frac{1}{n} \log \sum_{i_1, i_n=1}^{\alpha} \rho_{i_1} [B(t)^{n-1}]_{i_1, i_n} e^{tx_{i_n}}. \end{aligned}$$

Using a standard limit theorem for irreducible, aperiodic Markov chains [48, p. 356], one proves that for each  $1 \leq i, j \leq \alpha$

$$\lim_{n \rightarrow \infty} \frac{1}{n} \log [B(t)^n]_{i,j} = \log \lambda(t).$$

Details are given in [33, Lem. IX.4.1]. It follows that  $c(t) = \log \lambda(t)$ . Since  $\lambda(t)$  and thus  $c(t)$  are differentiable for all  $t$ , the Gärtner-Ellis Theorem implies that  $Y_n$  satisfies the large deviation principle on  $\mathbb{R}$  with the indicated rate function  $I$ .

(b) We refer the reader to [34, Thm. III.1], where this large deviation principle for the empirical vectors  $L_n$  is proved. The basic idea is that as in part (a) the rate function is given by a Legendre-Fenchel transform in  $\mathbb{R}^\alpha$ . One then shows that this Legendre-Fenchel transform equals  $I_\pi$  on  $\mathcal{P}_\alpha$  and equals  $\infty$  on  $\mathbb{R}^\alpha \setminus \mathcal{P}_\alpha$ . ■

We end this section by examining several features of the Gärtner-Ellis Theorem. Since in that theorem the rate function is always convex, a natural question is whether there exist large deviation principles having nonconvex rate functions. Two such examples are given next. Additional examples appear in [24].

One of the hypotheses of the Gärtner-Ellis Theorem is the differentiability of the limit function  $c(t)$ . An interesting problem is to investigate the existence of large deviation principles when this condition is violated. Unfortunately, the situation is complicated and a general theory does not exist. In Example 8.3, the differentiability of the limit function  $c(t)$  does not hold and the rate function is not given by a Legendre-Fenchel transform. In another example arising in the Ising model in statistical mechanics, the same hypothesis of the Gärtner-Ellis Theorem is not valid for all sufficiently large values of the inverse temperature defining the model. However, the rate function in the large deviation principle for the spin per site is defined by the identical Legendre-Fenchel transform appearing in the statement of the Gärtner-Ellis Theorem [35, Thm. 11.1].

The first example involves an extreme case of dependent random variables.

**Example 8.3.** We define a random variable  $X_1$  by the probability distribution  $P\{X_1 = 1\} = P\{X_1 = -1\} = \frac{1}{2}$ . For each integer  $j \geq 2$  we define random variables  $X_j = X_1$ , and for  $n \in \mathbb{N}$  we set

$$Y_n = \frac{1}{n} \sum_{j=1}^n X_j.$$

Let us first try to apply the Gärtner-Ellis Theorem to the sequence  $Y_n$ . For each  $t \in \mathbb{R}$  and  $x \in \mathbb{R}$  we calculate

$$c(t) = \lim_{n \rightarrow \infty} \frac{1}{n} \log E\{\exp(ntY_n)\} = \lim_{n \rightarrow \infty} \frac{1}{n} \log \left( \frac{1}{2} [e^{nt} + e^{-nt}] \right) = |t|$$

and

$$I(x) = \sup_{t \in \mathbb{R}} \{tx - c(t)\} = \begin{cases} 0 & \text{if } |x| \leq 1 \\ \infty & \text{if } |x| > 1. \end{cases}$$

Since  $c(t) = |t|$  is not differentiable at  $t = 0$ , the Gärtner-Ellis Theorem is not applicable. In fact,  $Y_n$  satisfies the large deviation principle on  $\mathbb{R}$  with the rate function

$$J(x) = \begin{cases} 0 & \text{if } x \in \{1, -1\} \\ \infty & \text{if } x \in \mathbb{R} \setminus \{1, -1\}. \end{cases}$$

This is easily checked since  $W_n$  has the distribution  $P\{W_n = 1\} = P\{W_n = -1\} = \frac{1}{2}$ . The function  $I$  is the largest convex function less than or equal to the rate function  $J$ . This completes the first example. ■

The second example generalizes Cramer's Theorem to the setting of a random walk with an interface.

**Example 8.4.** We define the sets

$$\Lambda^{(1)} = \{x \in \mathbb{R}^d : x_1 \leq 0\}, \Lambda^{(2)} = \{x \in \mathbb{R}^d : x_1 > 0\}, \partial = \{x \in \mathbb{R}^d : x_1 = 0\},$$

where  $x_1$  denotes the first component of  $x \in \mathbb{R}^d$ . We define a random walk model for which the distribution of the next step depends on the halfspace  $\Lambda^{(1)}$  or  $\Lambda^{(2)}$  in which the random walk is currently located. To this end let  $\rho^{(1)}$  and  $\rho^{(2)}$  be two distinct probability measures on  $\mathbb{R}^d$ . Although it is not necessary, for simplicity we assume that the support of each measure is all of  $\mathbb{R}^d$ . Let  $\{X_j^{(1)}, j \in \mathbb{N}\}$  and  $\{X_j^{(2)}, j \in \mathbb{N}\}$  be independent sequences of i.i.d. random vectors with probability distributions  $P\{X_j^{(1)} \in dx\} = \rho^{(1)}(dx)$  and  $P\{X_j^{(2)} \in dx\} = \rho^{(2)}(dx)$ . We consider the stochastic process  $\{S_n, n \in \mathbb{N} \cup \{0\}\}$ , where  $S_0 = 0$  and  $S_{n+1}$  is defined recursively from  $S_n$  by the formula

$$S_{n+1} = S_n + 1_{\{S_n \in \Lambda^{(1)}\}} \cdot X_n^{(1)} + 1_{\{S_n \in \Lambda^{(2)}\}} \cdot X_n^{(2)}.$$

For  $i = 1, 2$ ,  $1_{\{S_n \in \Lambda^{(i)}\}}$  denotes the indicator function of the set  $\{S_n \in \Lambda^{(i)}\}$ . Because of the abrupt change in distribution across the surface  $\partial$ , we call this random walk a model with discontinuous statistics. In [29] we show that  $S_n/n$  satisfies the large deviation principle on  $\mathbb{R}^d$ . The rate function is given by an



explicit formula that takes a complicated form along the interface  $\partial$ . We will not give the definition of the rate function here, but merely note that in general it is a nonconvex function on  $\mathbb{R}^d$  which is convex in each of the halfspaces  $\Lambda^{(1)}$  and  $\Lambda^{(2)}$ . If the measures  $\rho^{(1)}$  and  $\rho^{(2)}$  coincide, then the main theorem of [29] reduces to Cramér's Theorem.

The large deviation phenomena investigated in [29] are an example of the fascinating problems that arise in the study of other Markov processes with discontinuous statistics. The main theorem of [29] is generalized in [31, Ch. 6] to a large deviation principle for the entire path of the random walk. In [32] a large deviation upper bound is proved for a general class of Markov processes with discontinuous statistics. An important group of processes with discontinuous statistics arises in the study of queueing systems. The large deviation principle for a general class of such systems is proved in [30]. This completes the second example. ■

In the next section we begin our study of statistical mechanical models by considering the Curie-Weiss spin model.

## 9 The Curie-Weiss Model of Ferromagnetism

The Curie-Weiss model of ferromagnetism is one of the simplest examples of an interacting system in statistical mechanics. As we will see in the next section, using the theory of large deviations to analyze this model suggests how one can apply the theory to analyze much more complicated models.

The Curie-Weiss model is a spin system on the configuration spaces  $\Omega_n = \{-1, 1\}^n$ ; the value  $-1$  represents “spin-down” and the value  $1$  “spin-up.” Let  $\rho = \frac{1}{2}\delta_{-1} + \frac{1}{2}\delta_1$  and let  $P_n$  denote the product measure on  $\Omega_n$  with one-dimensional marginals  $\rho$ . Thus  $P_n\{\omega\} = 1/2^n$  for each configuration or microstate  $\omega = \{\omega_i, i = 1, \dots, n\} \in \Omega_n$ . The Hamiltonian, or energy, of  $\omega$  is defined by

$$H_n(\omega) = -\frac{1}{2n} \sum_{i,j=1}^n \omega_i \omega_j = -\frac{n}{2} \left( \frac{1}{n} \sum_{j=1}^n \omega_j \right)^2, \quad (9.1)$$

and the probability of  $\omega$  corresponding to inverse temperature  $\beta > 0$  is defined by the canonical ensemble

$$P_{n,\beta}\{\omega\} = \frac{1}{Z_n(\beta)} \exp[-\beta H_n(\omega)] P_n\{\omega\}, \quad (9.2)$$

where  $Z_n(\beta)$  is the partition function

$$Z_n(\beta) = \int_{\Omega_n} \exp[-\beta H_n(\omega)] P_n(d\omega) = \sum_{\omega \in \Omega_n} \exp[-\beta H_n(\omega)] \frac{1}{2^n}.$$

$P_{n,\beta}$  models a ferromagnet in the sense that the maximum of  $P_{n,\beta}\{\omega\}$  over  $\omega \in \Omega_n$  occurs at the two microstates having all coordinates  $\omega_i$  equal to  $-1$  or all coordinates equal to  $1$ . Furthermore, as  $\beta \rightarrow \infty$  all the mass of  $P_{n,\beta}$  concentrates on these two microstates. The Curie-Weiss model is used as a mean-field approximation to the much more complicated Ising model and related short-range, ferromagnetic models [33, §V.9].

A distinguishing feature of the Curie-Weiss model is its phase transition. Namely, the alignment effects incorporated in the Gibbs states  $P_{n,\beta}$  persist in the limit  $n \rightarrow \infty$ . This is most easily seen by evaluating the  $n \rightarrow \infty$  limit of the distributions  $P_{n,\beta}\{S_n/n \in dx\}$ , where  $S_n(\omega)/n$  equals the spin per site

$\sum_{j=1}^n \omega_j/n$ . We will see that for  $\beta \leq 1$  this limit acts like the classical weak law of large numbers, concentrating on the value 0. However, for  $\beta > 1$  the analogy with the classical law of large numbers breaks down; the alignment effects are so strong that the limiting  $P_{n,\beta}$ -distribution of  $S_n/n$  concentrates on the two points  $\pm m(\beta)$  for some  $m(\beta) \in (0, 1)$ . The analysis of the Curie-Weiss model to be presented below can be easily modified to handle an external magnetic field  $h$ . The resulting probabilistic description of the phase transition yields the predictions of mean field theory [33, §V.9], [67, §3.2].

We calculate the  $n \rightarrow \infty$  limit of  $P_{n,\beta}\{S_n/n \in dx\}$  by establishing a large deviation principle for the spin per site with respect to  $P_{n,\beta}$ . For each  $n$ ,  $S_n/n$  takes values in  $[-1, 1]$ . By the equivalence between the Laplace principle and the large deviation principle asserted in Theorem 6.9, it suffices to find a rate function  $I_\beta$  on  $[-1, 1]$  such that for any continuous function  $f$  mapping  $[-1, 1]$  into  $\mathbb{R}$

$$\lim_{n \rightarrow \infty} \frac{1}{n} \log \int_{\Omega_n} \{\exp[nf(S_n/n)]\} dP_{n,\beta} = \sup_{x \in [-1,1]} \{f(x) - I_\beta(x)\}.$$

In order to prove this Laplace principle, we define  $\psi(x) = -\frac{1}{2}\beta x^2$  for  $x \in [-1, 1]$  and appeal to a number of results established earlier in these lectures. Since

$$H_n(\omega) = -\frac{n}{2} \left( \frac{1}{n} \sum_{j=1}^n \omega_j \right)^2 = -\frac{n}{2} \left( \frac{S_n(\omega)}{n} \right)^2,$$

we can write

$$P_{n,\psi}\{\omega\} = \frac{1}{\int_{[-1,1]} \exp[-n\psi(S_n/n)] dP_n} \cdot \exp[-n\psi(S_n/n(\omega))] P_n(\omega).$$

In addition, by the version of Cramér's Theorem given in Corollary 6.6, with respect to  $P_n$ ,  $S_n/n$  satisfies the Laplace principle with rate function

$$I(x) = \frac{1}{2}(1-x) \log(1-x) + \frac{1}{2}(1+x) \log(1+x).$$

We can thus apply Theorem 6.13 with  $\mathcal{X} = [-1, 1]$ . We restate the theorem here for easy reference.

**Theorem 9.1.** *Assume that with respect to the probability measures  $P_n$ ,  $Y_n$  satisfies the large deviation principle on  $\mathcal{X}$  with rate function  $I$ . Let  $\psi$  be a bounded, continuous function mapping  $\mathcal{X}$  into  $\mathbb{R}$ . For  $A \in \mathcal{F}_n$  we define new probability measures*

$$P_{n,\psi}\{A\} = \frac{1}{\int_{\mathcal{X}} \exp[-n\psi(Y_n)] dP_n} \cdot \int_A \exp[-n\psi(Y_n)] dP_n.$$

*Then with respect to  $P_{n,\psi}$ ,  $Y_n$  satisfies the large deviation principle on  $\mathcal{X}$  with rate function*

$$I_\psi(x) = I(x) + \psi(x) - \inf_{y \in \mathcal{X}} \{I(y) + \psi(y)\}.$$

This gives the following large deviation principle for  $S_n/n$  with respect to the Curie-Weiss model. We write the rate function as  $I_\beta$  rather than as  $I_\psi$ .

**Theorem 9.2.** *With respect to the canonical ensemble  $P_{n,\beta}$  defined in (9.2), the spin per site  $S_n/n$  satisfies the large deviation principle on  $[-1, 1]$  with rate function*

$$I_\beta(x) = I(x) - \frac{1}{2}\beta x^2 - \inf_{y \in [-1,1]} \{I(y) - \frac{1}{2}\beta y^2\}.$$

The limiting behavior of the distributions  $P_{n,\beta}\{S_n/n \in dx\}$  is now determined by examining where  $I_\beta$  attains its infimum of 0 [33, §IV.4]. Infimizing points  $x^*$  satisfy

$$I'_\beta(x^*) = 0 \quad \text{or} \quad I'(x^*) = \beta x^*.$$

The second equation is equivalent to the mean field equation  $x^* = (I')^{-1}(\beta x^*) = \tanh(\beta x^*)$  [33, §V.9], [67, §3.2]. The next theorem is a consequence of the following easily verified properties of  $I$ :

- $I''(0) = 1$ .
- $I'$  is convex on  $[0, 1]$  and  $\lim_{x \rightarrow 1} I'(x) = \infty$ .
- $I'$  is concave on  $[-1, 0]$  and  $\lim_{x \rightarrow -1} I'(x) = -\infty$ .

**Theorem 9.3.** For each  $\beta > 0$  we define  $\mathcal{E}_\beta = \{x \in [-1, 1] : I_\beta(x) = 0\}$ . The following conclusions hold.

(a) For  $0 < \beta \leq 1$ ,  $\mathcal{E}_\beta = \{0\}$ .

(b) For  $\beta > 1$  there exists  $m(\beta) > 0$  such that  $\mathcal{E}_\beta = \{\pm m(\beta)\}$ . The function  $m(\beta)$  is monotonically increasing on  $(1, \infty)$  and satisfies  $m(\beta) \rightarrow 0$  as  $\beta \rightarrow 1^+$ ,  $m(\beta) \rightarrow 1$  as  $\beta \rightarrow \infty$ .

According to Proposition 6.4, if  $A$  is any closed subset of  $[-1, 1]$  such that  $A \cup \mathcal{E}_\beta = \emptyset$ , then  $I(A) > 0$  and for some  $C < \infty$

$$P_{n,\beta}\{S_n/n \in A\} \leq C \exp[-nI(A)/2] \rightarrow 0 \text{ as } n \rightarrow \infty.$$

In combination with Theorem 9.3, we are led to the following weak limits:

$$P_{n,\beta}\left\{\frac{1}{n}\sum_{i=1}^n \omega_i \in dx\right\} \Longrightarrow \begin{cases} \delta_0 & \text{if } 0 < \beta \leq 1 \\ \frac{1}{2}\delta_{m(\beta)} + \frac{1}{2}\delta_{-m(\beta)} & \text{if } \beta > 1. \end{cases} \quad (9.3)$$

We call  $m(\beta)$  the spontaneous magnetization for the Curie-Weiss model and  $\beta_c = 1$  the critical inverse temperature [33, §IV.4].

The limit (9.3) justifies calling  $\mathcal{E}_\beta$  the set of equilibrium macrostates for the spin per site  $S_n/n$  in the Curie-Weiss model. Because  $m(\beta) \rightarrow 0$  as  $\beta \rightarrow 1^+$  and 0 is the unique equilibrium macrostate for  $0 < \beta \leq 1$ , the phase transition at  $\beta_c$  is said to be continuous or second order. It is not difficult to show that points  $x^* \in \mathcal{E}_\beta$  have an equivalent characterization in terms of a maximum entropy principle. Because of the relatively simple nature of the model, this maximum entropy principle takes a rather trivial form. The details are omitted.

Before leaving the Curie-Weiss model, there are several points that should be emphasized. The first is to emphasize what makes possible the large deviation analysis of the phase transition in the model. In (9.1) we write the Hamiltonian as a quadratic function of the spin per site  $S_n/n$ , which by the version of Cramér's Theorem given in Corollary 6.6 satisfies the large deviation principle on  $[-1, 1]$  with respect to the product measures  $P_n$ . The equivalent Laplace principle allows us to convert this large deviation principle into a large deviation principle with respect to the canonical ensemble  $P_{n,\beta}$ . The form of the rate function  $I_\beta$  allows us to complete the analysis. In the next section we will

generalize these steps to formulate a large deviation approach to a wide class of models in statistical mechanics.

Our large deviation analysis of the phase transition in the Curie-Weiss model has the attractive feature that it directly motivates the physical importance of  $\mathcal{E}_\beta$ . This set is the support of the  $n \rightarrow \infty$  limit of the distributions  $P_{n,\beta}\{S_n/n \in dx\}$ . As we will see in the next section, an analogous fact is true for a large class of statistical mechanical models [Thm. 10.3].

The large deviation analysis of the Curie-Weiss model yields the limiting behavior of the  $P_{n,\beta}$ -distributions of  $S_n/n$ . For  $0 < \beta \leq 1$  this limit corresponds to the classical weak law of large numbers for the sample means of i.i.d. random variables and suggests examining the analogues of other classical limit results such as the central limit theorem. We end this section by summarizing these limit results for the Curie-Weiss, referring the reader to [33, §V.9] for proofs. If  $\theta \in (0, 1)$  and  $f$  is a nonnegative integrable function on  $\mathbb{R}$ , then the notation  $P_{n,\beta}\{S_n/n^\theta \in dx\} \implies f dx$  means that the distributions of  $S_n/n^\theta$  converge weakly to the probability measure on  $\mathbb{R}$  having a density proportional to  $f$  with respect to Lebesgue measure.

In the Curie-Weiss model for  $0 < \beta < 1$ , the interactions among the spins are relatively weak, and the analogue of the central limit theorem holds [33, Thm. V.9.4]:

$$P_{n,\beta}\{S_n/n^{1/2} \in dx\} \implies \exp[-\frac{1}{2}x^2/\sigma^2(\beta)] dx,$$

where  $\sigma^2(\beta) = 1/(1 - \beta)$ . However, when  $\beta = \beta_c = 1$ , the limiting variance  $\sigma^2(\beta)$  diverges, and the central limit scaling  $n^{1/2}$  must be replaced by  $n^{3/4}$ , which reflects the onset of long-range order at  $\beta_c$ . In this case we have [33, Thm. V.9.5]

$$P_{n,\beta_c}\{S_n/n^{3/4} \in dx\} \implies \exp[-\frac{1}{12}x^4] dx.$$

Finally, for  $\beta > \beta_c$ ,  $(S_n - n\tilde{z})/n^{1/2}$  satisfies a central-limit-type theorem when  $S_n/n$  is conditioned to lie in a sufficiently small neighborhood of  $\tilde{z} = m(\beta)$  or  $\tilde{z} = -m(\beta)$ ; see Theorem 2.4 in [41] with  $k = 1$ .

The results discussed in this section have been extensively generalized to a number of models, including the Curie-Weiss-Potts model [13, 44], the mean-field Blume-Emery-Griffiths model [12, 42], and the Ising and related models

[25, 49, 65]. For the latter models, refined large deviations at the surface level have been studied; see [20, p. 339] for references.

## 10 Equivalence of Ensembles for a General Class of Models in Statistical Mechanics

Equilibrium statistical mechanics specifies two ensembles that describe the probability distribution of microstates in statistical mechanical models. These are the microcanonical ensemble and the canonical ensemble. Particularly in the case of models of coherent structures in turbulence, the microcanonical ensemble is physically more fundamental because it expresses the fact that the Hamiltonian is a constant of the Euler dynamics underlying the model.

The introduction of two separate ensembles raises the basic problem of ensemble equivalence. As we will see in this section, the theory of large deviations and the theory of convex functions provide the perfect tools for analyzing this problem, which forces us to re-evaluate a number of deep questions that have often been dismissed in the past as being physically obvious. These questions include the following. Is the temperature of a statistical mechanical system always related to its energy in a one-to-one fashion? Are the microcanonical equilibrium properties of a system calculated as a function of the energy always equivalent to its canonical equilibrium properties calculated as a function of the temperature? Is the microcanonical entropy always a concave function of the energy? Is the heat capacity always a positive quantity? Surprisingly, the answer to each of these questions is in general no.

Starting with the work of Lynden-Bell and Wood [58] and the work of Thirring [75], physicists have come to realize in recent decades that systematic incompatibilities between the microcanonical and canonical ensembles can arise in the thermodynamic limit if the microcanonical entropy function of the system under study is nonconcave. The reason for this nonequivalence can be explained mathematically by the fact that when applied to a nonconcave function the Legendre-Fenchel transform is non-involutive; i.e., performing it twice does not give back the original function but gives back its concave envelope [42, 76]. As a consequence of this property, the Legendre-Fenchel structure of statistical mechanics, traditionally used to establish a one-to-one relationship between the entropy and the free energy and between the energy and the temperature, ceases



to be valid when the entropy is nonconcave.

From a more physical perspective, the explanation is even simpler. When the entropy is nonconcave, the microcanonical and canonical ensembles are nonequivalent because the nonconcavity of the entropy implies the existence of a nondifferentiable point of the free energy, and this, in turn, marks the presence of a first-order phase transition in the canonical ensemble [36, 51]. Accordingly, the ensembles are nonequivalent because the canonical ensemble jumps over a range of energy values at a critical value of the temperature and is therefore prevented from entering a subset of energy values that can always be accessed by the microcanonical ensemble [36, 51, 75]. This phenomenon lies at the root of ensemble nonequivalence, which is observed in systems as diverse as lattice spin models, including the Curie-Weiss-Potts model [13, 14], the mean-field Blume-Emery-Griffiths model [2, 3, 42, 43], mean-field versions of the Hamiltonian model [19, 56], and the XY model [18]; in gravitational systems [51, 52, 58, 75]; in models of coherent structures in turbulence [9, 36, 37, 47, 53, 70]; in models of plasmas [54, 73]; and in a model of the Lennard-Jones gas [5], to mention only a few. Many of these models can be analyzed by the methods to be introduced in this section, which summarize the results in [36]. Further developments in the theory are given in [15]. The reader is referred to these two paper for additional references to the large literature on ensemble equivalence for classical lattice systems and other models.

In the examples cited in the preceding paragraph as well as in other cases, the microcanonical formulation gives rise to a richer set of equilibrium macrostates than the canonical formulation, a phenomenon that occurs especially in the negative temperature regimes of the vorticity dynamics models [22, 23, 47, 53]. For example, it has been shown computationally that the strongly reversing zonal-jet structures on Jupiter as well as the Great Red Spot fall into the nonequivalent range of the microcanonical ensemble with respect to the energy and circulation invariants [78].

The general class of models to be considered include both spin models and models of coherent structures in turbulence, and for these two sets of models several of the definitions take slightly different forms. The models to be con-

sidered are defined in terms of the following quantities. After presenting the general setup, we will verify that it applies to the Curie-Weiss model. The large deviation analysis of that model, summarized in the preceding section, inspired the general approach presented here.

- A sequence of probability spaces  $(\Omega_n, \mathcal{F}_n, P_n)$  indexed by  $n \in \mathbb{N}$ , which typically represents a sequence of finite dimensional systems. The  $\Omega_n$  are the configuration spaces,  $\omega \in \Omega_n$  are the microstates, and the  $P_n$  are the prior measures.
- For each  $n \in \mathbb{N}$  the Hamiltonian  $H_n$ , a bounded, measurable function mapping  $\Omega_n$  into  $\mathbb{R}$ .
- A sequence of positive scaling constants  $a_n \rightarrow \infty$  as  $n \rightarrow \infty$ . In general  $a_n$  equals the total number of degrees of freedom in the model. In many cases  $a_n$  equals the number of particles.

Models of coherent structures in turbulence often incorporate other dynamical invariants besides the Hamiltonian; we will see such a model in the next section. In this case one replaces  $H_n$  in the second bullet by the vector of dynamical invariants and makes other corresponding changes in the theory, which are all purely notational. For simplicity we work only with the Hamiltonian in this section.

A large deviation analysis of the general model is possible provided that there exist, as specified in the next four items, a space of macrostates, a sequence of macroscopic variables, and an interaction representation function and provided that the macroscopic variables satisfy the large deviation principle on the space of macrostates. Item 3 takes one form for spin models and a different form for models of coherent structures in turbulence. Items 1, 2, and 4 are the same for these two sets of models.

1. **Space of macrostates.** This is a complete, separable metric space  $\mathcal{X}$ , which represents the set of all possible macrostates.

2. **Macroscopic variables.** These are a sequence of random variables  $Y_n$  mapping  $\Omega_n$  into  $\mathcal{X}$ . These functions associate a macrostate in  $\mathcal{X}$  with each microstate  $\omega \in \Omega_n$ .
3. **Hamiltonian representation function.** This is a bounded, continuous function  $\tilde{H}$  that maps  $\mathcal{X}$  into  $\mathbb{R}$  and enables us to write  $H_n$ , either exactly or asymptotically, as a function of the macrostate via the macroscopic variable  $Y_n$ . The precise description for the two sets of models is as follows.

*Spin models.* As  $n \rightarrow \infty$

$$H_n(\omega) = a_n \tilde{H}(Y_n(\omega)) + o(a_n) \quad \text{uniformly for } \omega \in \Omega_n;$$

i.e.,

$$\lim_{n \rightarrow \infty} \sup_{\omega \in \Omega_n} \left| \frac{1}{a_n} H_n(\omega) - \tilde{H}(Y_n(\omega)) \right| = 0. \quad (10.1)$$

*Models of coherent structures in turbulence.* As  $n \rightarrow \infty$

$$H_n(\omega) = \tilde{H}(Y_n(\omega)) + o(1) \quad \text{uniformly for } \omega \in \Omega_n;$$

i.e.,

$$\lim_{n \rightarrow \infty} \sup_{\omega \in \Omega_n} |H_n(\omega) - \tilde{H}(Y_n(\omega))| = 0. \quad (10.2)$$

4. **Large deviation principle for the macroscopic variables.** There exists a function  $I$  mapping  $\mathcal{X}$  into  $[0, \infty]$  and having compact level sets such that with respect to  $P_n$  the sequence  $Y_n$  satisfies the LDP on  $\mathcal{X}$  with rate function  $I$  and scaling constants  $a_n$ . In other words, for any closed subset  $F$  of  $\mathcal{X}$

$$\limsup_{n \rightarrow \infty} \frac{1}{a_n} \log P_n \{Y_n \in F\} \leq - \inf_{x \in F} I(x),$$

and for any open subset  $G$  of  $\mathcal{X}$

$$\liminf_{n \rightarrow \infty} \frac{1}{a_n} \log P_n \{Y_n \in G\} \geq - \inf_{x \in G} I(x).$$

We now verify that this general setup applies to the Curie-Weiss model.

**Example 10.1.**

- $n$  spins  $\omega_i \in \{-1, 1\}$ .
- Microstates:  $\omega = (\omega_1, \omega_2, \dots, \omega_n) \in \Omega_n = \{-1, 1\}^n$ .

- Prior measures:

$$P_n(\omega) = \frac{1}{2^n} \text{ for each } \omega \in \Omega_n.$$

- Scaling constants:  $a_n = n$ .
- Hamiltonians:

$$H_n(\omega) = H_n(\omega) = -\frac{1}{2n} \sum_{i,j=1}^n \omega_i \omega_j = -\frac{n}{2} \left( \frac{1}{n} \sum_{j=1}^n \omega_j \right)^2.$$

- Macroscopic variables:

$$Y_n(\omega) = \frac{1}{n} S_n(\omega) = \frac{1}{n} \sum_{j=1}^n \omega_j.$$

- $Y_n$  maps  $\Omega_n$  into  $[-1, 1]$ , which is the space of macrostates.
- Energy representation function:

$$H_n(\omega) = -\frac{1}{2}(Y_n(\omega))^2 = \tilde{H}(Y_n(\omega)), \text{ where } \tilde{H}(x) = -\frac{1}{2}x^2 \text{ for } x \in [-1, 1].$$

Thus (10.1) holds with equality for all  $\omega$  without the error term  $o(a_n)$ .

- Large deviation principle with respect to  $P_n$ :

$$P_n\{Y_n \in dx\} \asymp e^{-nI(x)}.$$

The version of Cramér's Theorem given in Corollary 6.6 gives the rate function

$$I(x) = \frac{1}{2}(1-x) \log(1-x) + \frac{1}{2}(1+x) \log(1+x).$$

This completes the example. ■

Here is a partial list of statistical mechanical models to which the large deviation formalism has been applied. Further details are given in [15, Ex. 2.1].

- The Miller-Robert model of fluid turbulence based on the two dimensional Euler equations [6]. This will be discussed in section 11.
- A model of geophysical flows based on equations describing barotropic, quasi-geostrophic turbulence [37].
- A model of soliton turbulence based on a class of generalized nonlinear Schrödinger equations [38]
- Lattice spin models including the Curie-Weiss model [33, §IV.4], the Curie-Weiss-Potts model [13], the mean-field Blume-Emery-Griffiths spin model [42], and the Ising model [49, 65]. The large deviation analysis of these models illustrate the three levels of the Donsker-Varadhan theory of large deviations, which are explained in Chapter 1 of [33].
  - Level 1. As we have seen, for the Curie-Weiss model the macroscopic variables are the sample means of i.i.d. random variables, and the large deviation principle with respect to the prior measures is the version of Cramér's Theorem given in Corollary 6.6.
  - Level 2. For the Curie-Weiss-Potts model [13] and the mean-field Blume-Emery-Griffiths spin model [42] the macroscopic variables are empirical vectors of i.i.d. random variables, and the large deviation principle with respect to the prior measures is the version of Sanov's Theorem given in Theorem 3.4.
  - Level 3. For the Ising model the macroscopic variables are an infinite-dimensional generalization of the empirical measure known as the empirical field, and the large deviation principle with respect to the prior measures is derived in [49, 65]. This is related to level 3 of the Donsker-Varadhan theory, which is formulated for a general class of Markov chains and Markov processes [28]. A special case is treated in [33, Ch. IX], which proves the large deviation principle for the empirical process of i.i.d. random variables taking values in a finite state

space. The complicated large deviation analysis of the Ising model is outlined in [35, §11].

Returning now to the general theory, we introduce the microcanonical ensemble, the canonical ensemble, and the basic thermodynamic functions associated with each ensemble: the microcanonical entropy and the canonical free energy. We then sketch the proofs of the large deviation principles for the macroscopic variables  $Y_n$  with respect to the two ensembles. As in the case of the Curie-Weiss model, the zeroes of the corresponding rate functions define the corresponding sets of equilibrium macrostates, one for the microcanonical ensemble and one for the canonical ensemble. The problem of ensemble equivalence investigates the relationship between these two sets of equilibrium macrostates.

In general terms, the main result is that a necessary and sufficient condition for equivalence of ensembles to hold at the level of equilibrium macrostates is that it holds at the level of thermodynamic functions, which is the case if and only if the microcanonical entropy is concave. The necessity of this condition has the following striking formulation. If the microcanonical entropy is not concave at some value of its argument, then the ensembles are nonequivalent in the sense that the corresponding set of microcanonical equilibrium macrostates is disjoint from any set of canonical equilibrium macrostates. The reader is referred to [36, §1.4] for a detailed discussion of models of coherent structures in turbulence in which nonconcave microcanonical entropies arise.

We start by introducing the function whose support and concavity properties completely determine all aspects of ensemble equivalence and nonequivalence. This function is the microcanonical entropy, defined for  $u \in \mathbb{R}$  by

$$s(u) = -\inf\{I(x) : x \in \mathcal{X}, \tilde{H}(x) = u\}. \quad (10.3)$$

Since  $I$  maps  $\mathcal{X}$  into  $[0, \infty]$ ,  $s$  maps  $\mathbb{R}$  into  $[-\infty, 0]$ . Moreover, since  $I$  is lower semicontinuous and  $\tilde{H}$  is continuous on  $\mathcal{X}$ ,  $s$  is upper semicontinuous on  $\mathbb{R}$ . We define  $\text{dom } s$  to be the set of  $u \in \mathbb{R}$  for which  $s(u) > -\infty$ . In general,  $\text{dom } s$  is nonempty since  $-s$  is a rate function [36, Prop. 3.1(a)]. The microcanonical ensemble takes two different forms depending on whether we consider spin

models or models of coherent structures in turbulence. For each  $u \in \text{dom } s$ ,  $r > 0$ ,  $n \in \mathbb{N}$ , and  $A \in \mathcal{F}_n$  the microcanonical ensemble for spin models is defined to be the conditioned measure

$$P_n^{u,r}\{A\} = P_n\{A \mid H_n/a_n \in [u - r, u + r]\}.$$

For models of coherent structures in turbulence we work with

$$P_n^{u,r}\{A\} = P_n\{A \mid H_n \in [u - r, u + r]\}.$$

As shown in [36, p. 1027], if  $u \in \text{dom } s$ , then for all sufficiently large  $n$  the conditioned measures  $P_n^{u,r}$  are well defined.

A mathematically more tractable probability measure is the canonical ensemble. For each  $n \in \mathbb{N}$ ,  $\beta \in \mathbb{R}$ , and  $A \in \mathcal{F}_n$  we define the partition function

$$Z_n(\beta) = \int_{\Omega_n} \exp[-\beta H_n] dP_n,$$

which is well defined and finite; the canonical free energy

$$\varphi(\beta) = - \lim_{n \rightarrow \infty} \frac{1}{a_n} \log Z_n(\beta);$$

and the probability measure

$$P_{n,\beta}\{A\} = \frac{1}{Z_n(\beta)} \cdot \int_A \exp[-\beta H_n] dP_n. \quad (10.4)$$

The measures  $P_{n,\beta}$  are Gibbs states that define the canonical ensemble for the given model. Although for spin models one usually takes  $\beta > 0$ , in general  $\beta \in \mathbb{R}$  is allowed; for example, negative values of  $\beta$  arise naturally in the study of coherent structures in two-dimensional turbulence.

Among other reasons, the canonical ensemble was introduced by Gibbs in the hope that in the limit  $n \rightarrow \infty$  the two ensembles are equivalent; all macroscopic properties of the model obtained via the microcanonical ensemble could be realized as macroscopic properties obtained via the canonical ensemble. However, as we will see, this in general is not the case.

The large deviation analysis of the canonical ensemble for spin models is summarized in the next theorem, Theorem 10.2. Additional information is given

in Theorem 10.3. The modifications in these two theorems necessary for analyzing the canonical ensemble for models of coherent structures in turbulence are indicated in Theorem 10.4.

Part (a) of Theorem 10.2 shows that the limit defining  $\varphi(\beta)$  exists and is given by the variational formula

$$\varphi(\beta) = \inf_{x \in \mathcal{X}} \{\beta \tilde{H}(x) + I(x)\}.$$

If in the definition of  $Z_n(\beta)$  one could replace  $H_n$  by  $a_n \tilde{H}$ , then this limit is a direct consequence of the Laplace principle for  $Y_n$  with respect to  $P_n$ , which is equivalent to the assumed large deviation principle for  $Y_n$  with respect to  $P_n$ . As we will see in the proof, the approximation property (10.1) of  $\tilde{H}$  allows us to make this replacement. Part (b) of Theorem 10.2 states the large deviation principle for the macroscopic variables with respect to canonical ensemble. This large deviation principle is easy to see. If in the definition of  $P_{n,\beta}$  one replaces  $H_n$  by  $a_n \tilde{H}$ , then it follows immediately from Theorem 6.13. Part (b) is the analogue of Theorem 9.2 for the Curie-Weiss model. In part (c) we consider the set  $\mathcal{E}_\beta$  consisting of points at which the rate function in part (b) attains its infimum of 0. The second property of  $\mathcal{E}_\beta$  given in part (c) justifies calling this the set of canonical equilibrium macrostates. Part (c) is a special case of Proposition 6.4.

**Theorem 10.2 (Canonical ensemble for spin models).** *For the general spin model we assume that there exists a space of macrostates  $\mathcal{X}$ , macroscopic variables  $Y_n$ , and a Hamiltonian representation function  $\tilde{H}$  satisfying*

$$\lim_{n \rightarrow \infty} \sup_{\omega \in \Omega_n} \left| \frac{1}{a_n} H_n(\omega) - \tilde{H}(Y_n(\omega)) \right| = 0, \quad (10.5)$$

where  $H_n$  is the Hamiltonian. We also assume that with respect to the prior measures  $P_n$ ,  $Y_n$  satisfies the large deviation principle on  $\mathcal{X}$  with some rate function  $I$  and scaling constants  $a_n$ . For each  $\beta \in \mathbb{R}$  the following conclusions hold.

(a) *The canonical free energy  $\varphi(\beta) = -\lim_{n \rightarrow \infty} \frac{1}{a_n} \log Z_n(\beta)$  exists and is given by*

$$\varphi(\beta) = \inf_{x \in \mathcal{X}} \{\beta \tilde{H}(x) + I(x)\}.$$



(b) With respect to the canonical ensemble  $P_{n,\beta}$  defined in (10.4),  $Y_n$  satisfies the large deviation principle on  $\mathcal{X}$  with scaling constants  $a_n$  and rate function

$$I_\beta(x) = I(x) + \beta\tilde{H}(x) - \varphi(\beta).$$

(c) We define the set of canonical equilibrium macrostates

$$\mathcal{E}_\beta = \{x \in \mathcal{X} : I_\beta(x) = 0\}.$$

Then  $\mathcal{E}_\beta$  is a nonempty, compact subset of  $\mathcal{X}$ . In addition, if  $A$  is a Borel subset of  $\mathcal{X}$  such that  $\bar{A} \cap \mathcal{E}_\beta = \emptyset$ , then  $I_\beta(\bar{A}) > 0$  and for some  $C < \infty$

$$P_{n,\beta}\{Y_n \in A\} \leq C \exp[-nI_\beta(\bar{A})/2] \rightarrow 0 \text{ as } n \rightarrow \infty.$$

**Proof.** Once we take into account the error between  $H_n$  and  $a_n\tilde{H}(Y_n)$  expressed in (10.5), the proofs of (a) and (b) follow from the Laplace principle. Here are the details.

(a) By (10.5)

$$\begin{aligned} & \left| \frac{1}{a_n} \log Z_n(\beta) - \frac{1}{a_n} \log \int_{\Omega_n} \exp[-\beta a_n \tilde{H}(Y_n)] dP_n \right| \\ &= \left| \frac{1}{a_n} \log \int_{\Omega_n} \exp[-\beta H_n] dP_n - \frac{1}{a_n} \log \int_{\Omega_n} \exp[-\beta a_n \tilde{H}(Y_n)] dP_n \right| \\ &\leq |\beta| \frac{1}{a_n} \sup_{\omega \in \Omega_n} |H_n(\omega) - a_n \tilde{H}(Y_n(\omega))| \rightarrow 0 \text{ as } n \rightarrow \infty. \end{aligned}$$

Since  $\tilde{H}$  is a bounded continuous function mapping  $\mathcal{X}$  into  $\mathbb{R}$ , the Laplace principle satisfied by  $Y_n$  with respect to  $P_n$  yields part (a):

$$\begin{aligned} \varphi(\beta) &= - \lim_{n \rightarrow \infty} \frac{1}{a_n} \log Z_n(\beta) \\ &= - \lim_{n \rightarrow \infty} \frac{1}{a_n} \log \int_{\Omega_n} \exp[-\beta a_n \tilde{H}(Y_n)] dP_n \\ &= - \sup_{x \in \mathcal{X}} \{-\beta \tilde{H}(x) - I(x)\} \\ &= \inf_{x \in \mathcal{X}} \{\beta \tilde{H}(x) + I(x)\}. \end{aligned}$$

(b) Rather than derive this from Theorem 6.13, we proceed as in the proof of part (a), but now with  $P_{n,\beta}$  replacing  $P_n$ . For any bounded continuous function  $f$  mapping  $\mathcal{X}$  into  $\mathbb{R}$ , again (10.5) and the Laplace principle satisfied by  $Y_n$  with respect to  $P_n$  yield

$$\begin{aligned}
& \lim_{n \rightarrow \infty} \frac{1}{a_n} \log \int_{\Omega_n} \exp[a_n f(Y_n)] dP_{n,\beta} \\
&= \lim_{n \rightarrow \infty} \frac{1}{a_n} \log \int_{\Omega_n} \exp[a_n f(Y_n) - \beta H_n] dP_n - \lim_{n \rightarrow \infty} \frac{1}{a_n} \log Z_n(\beta) \\
&= \lim_{n \rightarrow \infty} \frac{1}{a_n} \log \int_{\Omega_n} \exp[a_n (f(Y_n) - \beta \tilde{H}(Y_n))] dP_n - \lim_{n \rightarrow \infty} \frac{1}{a_n} \log Z_n(\beta) \\
&= \sup_{x \in \mathcal{X}} \{f(x) - \beta \tilde{H}(x) - I(x)\} + \varphi(\beta) \\
&= \sup_{x \in \mathcal{X}} \{f(x) - I_\beta(x)\}.
\end{aligned}$$

By hypothesis,  $I$  has compact level sets and  $\tilde{H}$  is bounded and continuous. Thus  $I_\beta$  has compact level sets. Since  $I_\beta$  maps  $\mathcal{X}$  into  $[0, \infty]$ ,  $I_\beta$  is a rate function. We conclude that with respect to  $P_{n,\beta}$ ,  $Y_n$  satisfies the Laplace principle, and thus the equivalent large deviation principle, with scaling constants  $a_n$  and rate function  $I_\beta$ .

(c) This is proved in Proposition 6.4. The display in part (c) is based on the large deviation upper bound for  $Y_n$  with respect to  $P_{n,\beta}$ , which was proved in part (b). The proof of the theorem is complete. ■

The second property of  $\mathcal{E}_\beta$  given in part (c) of the theorem can be regarded as a concentration property of the  $P_{n,\beta}$ -distributions of  $Y_n$  which justifies calling  $\mathcal{E}_\beta$  the set of canonical equilibrium macrostates. With respect to these distributions, the probability of any Borel set  $A$  whose closure has empty intersection with  $\mathcal{E}_\beta$  goes to 0 exponentially fast with  $a_n$ . This large deviation characterization of the equilibrium macrostates is an attractive feature of our approach.

The concentration property of the  $P_{n,\beta}$ -distributions of  $Y_n$  as expressed in part (c) of the theorem has a refinement that arises in our study of the Curie-Weiss model. From Theorem 9.3 we recall that  $\mathcal{E}_\beta = \{0\}$  for  $0 < \beta \leq 1$  and  $\mathcal{E}_\beta = \{\pm m(\beta)\}$  for  $\beta > 1$ , where  $m(\beta)$  is the spontaneous magnetiza-

tion. According to (9.3), for all  $\beta > 0$  the weak limit of  $P_{n,\beta}\{S_n/n \in dx\}$  is concentrated on  $\mathcal{E}_\beta$ . While in the case of the general model treated in the present section one should not expect such a precise formulation, the next theorem gives considerable information, relating weak limits of subsequences of  $P_{n,\beta}\{Y_n \in dx\}$  to the set of equilibrium macrostates  $\mathcal{E}_\beta$ . For example, if one knows that  $\mathcal{E}_\beta$  consists of a unique point  $\tilde{x}$ , then it follows that the entire sequence  $P_{n,\beta}\{Y_n \in dx\}$  converges weakly to  $\delta_{\tilde{x}}$ . This situation corresponds to the absence of a phase transition. The proof of the theorem is technical and is omitted.

**Theorem 10.3 (Canonical ensemble for spin systems).** *We fix  $\beta \in \mathbb{R}$  and use the notation of Theorem 10.2. If  $\mathcal{E}_\beta$  consists of a unique point  $\tilde{x}$ , then  $P_{n,\beta}\{Y_n \in dx\} \Rightarrow \delta_{\tilde{x}}$ . If  $\mathcal{E}_\beta$  does not consist of a unique point, then any subsequence of  $P_{n,\beta}\{Y_n \in dx\}$  has a subsubsequence converging weakly to a probability measure  $\Pi_\beta$  on  $\mathcal{X}$  that is concentrated on  $\mathcal{E}_\beta$ ; i.e.,  $\Pi_\beta\{(\mathcal{E}_\beta)^c\} = 0$ .*

In order to carry out the large deviation analysis of the canonical ensemble for models of coherent structures in turbulence, in Theorems 10.2 and 10.3 one must make two changes: replace the limit (10.5) by

$$\lim_{n \rightarrow \infty} \sup_{\omega \in \Omega_n} |H_n(\omega) - \tilde{H}(Y_n(\omega))| = 0, \quad (10.6)$$

where  $H_n$  is the Hamiltonian, and replace  $Z_n(\beta)$  and  $P_{n,\beta}$  by  $Z_n(a_n\beta)$  and  $P_{n,a_n\beta}$ . For easy reference, this is summarized in the next theorem.

**Theorem 10.4 (Canonical ensemble for models of coherent structures in turbulence).** *For the general model of coherent structures in turbulence we assume that there exists a space of macrostates  $\mathcal{X}$ , macroscopic variables  $Y_n$ , and a Hamiltonian representation function  $\tilde{H}$  satisfying (10.6). We also assume that with respect to the prior measures  $P_n$ ,  $Y_n$  satisfies the large deviation principle on  $\mathcal{X}$  with some rate function  $I$  and scaling constants  $a_n$ . Then for each  $\beta \in \mathbb{R}$  all the conclusions of Theorems 10.2 and 10.3 are valid provided that  $Z_n(\beta)$  and  $P_{n,\beta}$  are replaced by  $Z_n(a_n\beta)$  and  $P_{n,a_n\beta}$ .*

In order to carry out the large deviation analysis of the microcanonical ensemble, we recall the relevant definitions. For  $u \in \mathbb{R}$  the microcanonical entropy is defined by

$$s(u) = -\inf\{I(x) : x \in \mathcal{X}, \tilde{H}(x) = u\}.$$

For each  $u \in \text{dom } s$ ,  $r > 0$ ,  $n \in \mathbb{N}$ , and set  $A \in \mathcal{F}_n$  the microcanonical ensemble for spin models is defined by

$$P_n^{u,r}\{A\} = P_n\{A \mid H_n/a_n \in [u - r, u + r]\}, \quad (10.7)$$

while the microcanonical ensemble for models of coherent states in turbulence is defined by

$$P_n^{u,r}\{A\} = P_n\{A \mid H_n \in [u - r, u + r]\}, \quad (10.8)$$

In order to simplify the discussion we will work with the microcanonical ensemble for models of coherent states in turbulence. The treatment of the microcanonical ensemble for spin models is analogous. We start our analysis of the microcanonical ensemble by pointing out that  $-s$  is the rate function in the large deviation principles, with respect to the prior measures  $P_n$ , of both  $\tilde{H}(Y_n)$  and  $H_n$ . In order to see this, we recall that with respect to  $P_n$ ,  $Y_n$  satisfies the large deviation principle with rate function  $I$ . Since  $\tilde{H}$  is a continuous function mapping  $\mathcal{X}$  into  $\mathbb{R}$ , the large deviation principle for  $\tilde{H}(Y_n)$  is a consequence of the contraction principle [Thm. 6.12]. For  $u \in \mathbb{R}$  the rate function is given by

$$\inf\{I(x) : x \in \mathcal{X}, \tilde{H}(x) = u\} = -s(u).$$

In addition, since

$$\lim_{n \rightarrow \infty} \sup_{\omega \in \Omega_n} |H_n(\omega) - \tilde{H}(Y_n(\omega))| = 0,$$

$H_n$  inherits from  $\tilde{H}(Y_n)$  the large deviation principle with the same rate function. This follows from Theorem 6.14 or can be derived as in the proof of Theorem 10.2 by using the equivalent Laplace principle. We summarize this large deviation principle by the notation

$$P_n\{H_n \in du\} \asymp \exp[a_n s(u)]. \quad (10.9)$$

For  $x \in \mathcal{X}$  and  $\alpha > 0$ ,  $B(x, \alpha)$  denotes the open ball with center  $x$  and radius  $\alpha$ . We next motivate the large deviation principle for  $Y_n$  with respect to the microcanonical ensemble  $P_n^{u,r}$  by estimating the exponential order contribution to the probability  $P_n^{u,r}\{Y_n \in B(x, \alpha)\}$  as  $n \rightarrow \infty$ . Specifically we seek a function  $I^u$  such that for all  $u \in \text{dom } s$ , all  $x \in \mathcal{X}$ , and all  $\alpha > 0$  sufficiently small

$$P_n^{u,r}\{Y_n \in B(x, \alpha)\} \approx \exp[-a_n I^u(x)] \text{ as } n \rightarrow \infty, r \rightarrow 0, \alpha \rightarrow 0. \quad (10.10)$$

The calculation that we present shows both the interpretive power of the large deviation notation and the value of left-handed thinking. Although the calculation is a bit complicated, it is much more straightforward than the actual proof, which is given in [36, §3] (see Thm. 3.20).

We first work with  $x \in \mathcal{X}$  for which  $I(x) < \infty$  and  $\tilde{H}(x) = u$ . Such an  $x$  exists since  $u \in \text{dom } s$  and thus  $s(u) > -\infty$ . Because

$$\lim_{n \rightarrow \infty} \sup_{\omega \in \Omega_n} |H_n(\omega) - \tilde{H}(Y_n(\omega))| = 0,$$

for all sufficiently large  $n$  depending on  $r$  the set of  $\omega$  for which both  $Y_n(\omega) \in B(x, \alpha)$  and  $H_n(\omega) \in [u - r, u + r]$  is approximately equal to the set of  $\omega$  for which both  $Y_n(\omega) \in B(x, \alpha)$  and  $\tilde{H}(Y_n(\omega)) \in [u - r, u + r]$ . Since  $\tilde{H}$  is continuous and  $\tilde{H}(x) = u$ , for all sufficiently small  $\alpha$  compared to  $r$  this set reduces to  $\{\omega : Y_n(\omega) \in B(x, \alpha)\}$ . Hence for all sufficiently small  $r$ , all sufficiently large  $n$  depending on  $r$ , and all sufficiently small  $\alpha$  compared to  $r$ , the assumed large deviation principle for  $Y_n$  with respect to  $P_n$  and the large deviation principle for  $H_n$  summarized in (10.9) yield

$$\begin{aligned} P_n^{u,r}\{Y_n \in B(x, \alpha)\} &= \frac{P_n\{\{Y_n \in B(x, \alpha)\} \cap \{H_n \in [u - r, u + r]\}\}}{P_n\{H_n \in [u - r, u + r]\}} \\ &\approx \frac{P_n\{Y_n \in B(x, \alpha)\}}{P_n\{H_n \in [u - r, u + r]\}} \\ &\approx \exp[-a_n(I(x) + s(u))]. \end{aligned}$$

On the other hand, if  $\tilde{H}(x) \neq u$ , then a similar calculation shows that for all sufficiently small  $r$ , all sufficiently small  $\alpha$ , and all sufficiently large  $n$

$P_n^{u,r}\{Y_n \in B(x, \alpha)\} = 0$ . Comparing these approximate calculations with the desired asymptotic form (10.10) motivates the correct formula for the rate function [36, Thm. 3.2]:

$$I^u(x) = \begin{cases} I(x) + s(u) & \text{if } \tilde{H}(x) = u, \\ \infty & \text{if } \tilde{H}(x) \neq u. \end{cases} \quad (10.11)$$

We record the facts in the next theorem, which takes the same form both for spin models and for models of coherent structures in turbulence. An additional complication occurs in the statement of the large deviation principle in part (b) because it involves the double limit  $n \rightarrow \infty$  followed by  $r \rightarrow 0$ . In part (c) we introduce the set of microcanonical equilibrium macrostates  $\mathcal{E}^u$  and state a concentration property of this set with respect to the microcanonical ensemble that is analogous to the concentration satisfied by the set  $\mathcal{E}_\beta$  of canonical equilibrium macrostates with respect to the canonical ensemble. The proof is similar to the proof of the analogous property of  $\mathcal{E}_\beta$  given in part (c) of Theorem 10.2, and it is therefore omitted.

**Theorem 10.5 (Microcanonical ensemble).** *Both for the general spin model and for the general model of coherent structures in turbulence we assume that there exists a space of macrostates  $\mathcal{X}$ , macroscopic variables  $Y_n$ , and a Hamiltonian representation function  $\tilde{H}$  satisfying (10.1) in the case of spin models and (10.2) in the case of models of coherent structures in turbulence. We also assume that with respect to the prior measures  $P_n$ ,  $Y_n$  satisfies the large deviation principle on  $\mathcal{X}$  with scaling constants  $a_n$  and some rate function  $I$ . For each  $u \in \text{dom } s$  and any  $r \in (0, 1)$  the following conclusions hold.*

(a) *With respect to  $P_n$ ,  $\tilde{H}(Y_n)$  and  $H_n$  both satisfy the large deviation principle with scaling constants  $a_n$  and rate function  $-s$ .*

(b) *We consider the microcanonical ensemble  $P_n^{u,r}$  defined in (10.7) for spin models and defined in (10.8) for models of coherent structures in turbulence. With respect to  $P_n^{u,r}$  and in the double limit  $n \rightarrow \infty$  and  $r \rightarrow 0$ ,  $Y_n$  satisfies the large deviation principle on  $\mathcal{X}$  with scaling constants  $a_n$  and rate function  $I^u$  defined in (10.11). That is, for any closed subset  $F$  of  $\mathcal{X}$*

$$\lim_{r \rightarrow 0} \limsup_{n \rightarrow \infty} \frac{1}{a_n} \log P_n^{u,r}\{Y_n \in F\} \leq -I^u(F)$$

and for any open subset  $G$  of  $\mathcal{X}$

$$\lim_{r \rightarrow 0} \liminf_{n \rightarrow \infty} \frac{1}{a_n} \log P_n^{u,r} \{Y_n \in G\} \geq -I^u(G).$$

(c) We define the set of equilibrium macrostates

$$\mathcal{E}^u = \{x \in \mathcal{X} : I^u(x) = 0\}.$$

Then  $\mathcal{E}^u$  is a nonempty, compact subset of  $\mathcal{X}$ . In addition, if  $A$  is a Borel subset of  $\mathcal{X}$  such that  $\overline{A} \cap \mathcal{E}^u = \emptyset$ , then  $I^u(\overline{A}) > 0$  and there exists  $r_0 > 0$  and for all  $r \in (0, r_0]$  there exists  $C_r < \infty$

$$P_{n,\beta} \{Y_n \in A\} \leq C_r \exp[-n I_\beta(\overline{A})/2] \rightarrow 0 \text{ as } n \rightarrow \infty.$$

In the remainder of this section we investigate issues related to the equivalence and nonequivalence of the canonical and microcanonical ensembles, which involves studying the relationships between the two sets of equilibrium macrostates

$$\mathcal{E}_\beta = \{x \in \mathcal{X} : I_\beta(x) = 0\} \text{ and } \mathcal{E}^u = \{x \in \mathcal{X} : I^u(x) = 0\}.$$

The following questions will be considered.

1. Given  $\beta \in \mathbb{R}$  and  $x \in \mathcal{E}_\beta$ , does there exist  $u \in \mathbb{R}$  such that  $x \in \mathcal{E}^u$ ? In other words, is any canonical equilibrium macrostate realized microcanonically?
2. Given  $u \in \mathbb{R}$  and  $x \in \mathcal{E}^u$ , does there exist  $\beta \in \mathbb{R}$  such that  $x \in \mathcal{E}_\beta$ ? In other words, is any microcanonical equilibrium macrostate realized canonically?

As we will see in Theorem 10.6, the answer to question 1 is always yes, but the answer to question 2 is much more complicated, involving three possibilities.

2a. **Full equivalence.** There exists  $\beta \in \mathbb{R}$  such that  $\mathcal{E}^u = \mathcal{E}_\beta$ .

2b. **Partial equivalence.** There exists  $\beta \in \mathbb{R}$  such that  $\mathcal{E}^u \subset \mathcal{E}_\beta$  but  $\mathcal{E}^u \neq \mathcal{E}_\beta$ .

2c. **Nonequivalence.**  $\mathcal{E}^u$  is disjoint from  $\mathcal{E}_\beta$  for all  $\beta \in \mathbb{R}$ .

One of the big surprises of the theory to be presented here is that we are able to decide on which of these three possibilities occur by examining support and concavity properties of the microcanonical entropy  $s(u)$ . This is remarkable because the sets  $\mathcal{E}_\beta$  and  $\mathcal{E}^u$  are in general infinite dimensional while the microcanonical entropy is a function on  $\mathbb{R}$ .

In order to begin our study of ensemble equivalence and nonequivalence, we first recall the definitions of the corresponding rate functions:

$$I_\beta(x) = I(x) + \beta\tilde{H}(x) - \varphi(\beta),$$

where  $\varphi(\beta)$  denotes the canonical free energy

$$\varphi(\beta) = \inf_{x \in \mathcal{X}} \{\beta\tilde{H}(x) + I(x)\},$$

and

$$I^u(x) = \begin{cases} I(x) + s(u) & \text{if } \tilde{H}(x) = u, \\ \infty & \text{if } \tilde{H}(x) \neq u, \end{cases}$$

where  $s(u)$  denotes the microcanonical entropy

$$s(u) = - \inf \{I(x) : x \in \mathcal{X}, \tilde{H}(x) = u\}.$$

Using these definitions, we see that the two sets of equilibrium macrostates have the alternate characterizations

$$\mathcal{E}_\beta = \{x \in \mathcal{X} : I(x) + \beta\tilde{H}(x) \text{ is minimized}\}$$

and

$$\mathcal{E}^u = \{x \in \mathcal{X} : I(x) \text{ is minimized subject to } \tilde{H}(x) = u\}.$$

Thus  $\mathcal{E}^u$  is defined by the following constrained minimization problem for  $u \in \mathbb{R}$ :

$$\text{minimize } I(x) \text{ over } \mathcal{X} \text{ subject to the constraint } \tilde{H}(x) = u. \quad (10.12)$$

By contrast,  $\mathcal{E}_\beta$  is defined by the following related, unconstrained minimization problem for  $\beta \in \mathbb{R}$ :

$$\text{minimize } I(x) + \beta\tilde{H}(x) \text{ over } x \in \mathcal{X}. \quad (10.13)$$



In this formulation  $\beta$  is a Lagrange multiplier dual to the constraint  $\tilde{H}(x) = u$ . The theory of Lagrange multipliers outlines suitable conditions under which the solutions of the constrained problem (10.12) lie among the critical points of  $I + \beta\tilde{H}$ . However, it does not give, as we will do in Theorems 10.6, necessary and sufficient conditions for the solutions of (10.12) to coincide with the solutions of the unconstrained minimization problem (10.13). These necessary and sufficient conditions are expressed in terms of support and concavity properties of the microcanonical entropy  $s(u)$ .

Before we explain this, we reiterate a number of properties of  $\varphi(\beta)$  and  $s(u)$  that emphasize the fundamental nature of these two thermodynamic functions. Properties 1, 2, and 3 show a complete symmetry between the canonical and microcanonical ensembles, a state of affairs that is spoiled by property 4.

1. Both  $\varphi(\beta)$  and  $s(u)$  are given by limits and by variational formulas.

- $\varphi(\beta)$  expresses the asymptotics of the partition function

$$Z_n(\beta) = \int_{\Omega_n} \exp[-\beta H_n] dP_n$$

via the definition

$$\varphi(\beta) = - \lim_{n \rightarrow \infty} \frac{1}{a_n} \log Z_n(\beta).$$

In addition,  $\varphi(\beta)$  is given by the variational formula [Thm. 10.2(a)]

$$\varphi(\beta) = \inf_{x \in \mathcal{X}} \{\beta \tilde{H}(x) + I(x)\}.$$

- $s(u)$  is defined by the variational formula

$$s(u) = - \inf \{I(x) : x \in \mathcal{X}, \tilde{H}(x) = u\}.$$

In addition  $s(u)$  expresses the asymptotics of  $P_n\{H_n \in du\}$ , which satisfies the large deviation principle with rate function  $-s(u)$  [Thm. 10.5(a)]; i.e.,  $P_n\{H_n \in du\} \asymp \exp[a_n s(u)]$ . Furthermore, for  $u \in \text{dom } s$  we have the limit [36, Prop. 3.1(c)]

$$s(u) = \lim_{r \rightarrow 0} \lim_{n \rightarrow \infty} \frac{1}{a_n} \log P_n\{H_n \in [u - r, u + r]\}.$$

2. Both  $\varphi(\beta)$  and  $s(u)$  are respectively the normalization constants in the rate functions  $I_\beta$  and  $I^u$  in the large deviation principles for  $Y_n$  with respect to the canonical ensemble and with respect to the microcanonical ensemble:

$$I_\beta(x) = I(x) + \beta\tilde{H}(x) - \varphi(\beta)$$

and

$$I^u(x) = \begin{cases} I(x) + s(u) & \text{if } \tilde{H}(x) = u, \\ \infty & \text{if } \tilde{H}(x) \neq u, \end{cases}$$

3. The sets of equilibrium macrostates have the alternate characterizations

$$\mathcal{E}_\beta = \{x \in \mathcal{X} : I(x) + \beta\tilde{H}(x) \text{ is minimized}\}$$

and

$$\mathcal{E}^u = \{x \in \mathcal{X} : I(x) \text{ is minimized subject to } \tilde{H}(x) = u\}.$$

- Thus  $\mathcal{E}_\beta$  consists of all  $x \in \mathcal{X}$  at which the infimum is attained in

$$\varphi(\beta) = \inf_{x \in \mathcal{X}} \{\beta\tilde{H}(x) + I(x)\}.$$

- Thus  $\mathcal{E}^u$  consists of all  $x \in \mathcal{X}$  at which the infimum is attained in

$$s(u) = -\inf\{I(x) : x \in \mathcal{X}, \tilde{H}(x) = u\}.$$

4.  $\varphi(\beta)$  and  $s(u)$  are related via the Legendre-Fenchel transform

$$\varphi(\beta) = \inf_{u \in \mathbb{R}} \{\beta u - s(u)\}. \quad (10.14)$$

As do the two formulas for  $\varphi(\beta)$  in item 1, this Legendre-Fenchel transform shows that  $\varphi(\beta)$  is always concave, even if  $s(u)$  is not. Unless  $s(u)$  is concave on  $\mathbb{R}$ , the dual formula  $s(u) = \inf_{\beta \in \mathbb{R}} \{\beta u - \varphi(\beta)\}$  is not valid.

- Proof 1 of (10.14) using variational formulas:

$$\begin{aligned} \varphi(\beta) &= \inf_{x \in \mathcal{X}} \{\beta\tilde{H}(x) + I(x)\} \\ &= \inf_{u \in \mathbb{R}} \inf\{\beta\tilde{H}(x) + I(x) : x \in \mathcal{X}, \tilde{H}(x) = u\} \\ &= \inf_{u \in \mathbb{R}} \{\beta u + \inf\{I(x) : x \in \mathcal{X}, \tilde{H}(x) = u\}\} \\ &= \inf_{u \in \mathbb{R}} \{\beta u - s(u)\}. \end{aligned}$$

- Proof 2 of (10.14) using asymptotic properties:

$$\begin{aligned}
\varphi(\beta) &= - \lim_{n \rightarrow \infty} \frac{1}{a_n} \log Z_n(\beta) \\
&= - \lim_{n \rightarrow \infty} \frac{1}{a_n} \log \int_{\Omega_n} \exp[-\beta H_n] dP_n \\
&= - \lim_{n \rightarrow \infty} \frac{1}{a_n} \log \int_{\mathbb{R}} \exp[-\beta u] P_n\{H_n \in du\} \\
&= - \sup_{u \in \mathbb{R}} \{-\beta u + s(u)\} \\
&= \inf_{u \in \mathbb{R}} \{\beta u - s(u)\}.
\end{aligned}$$

To derive the next-to-last line we use the fact that with respect to  $P_n$ ,  $H_n$  satisfies the large deviation principle, and therefore the equivalent Laplace principle, with rate function  $-s(u)$  [Thm. 10.5(a)]. Invoking the Laplace principle is a bit of cheating since the identity function mapping  $u \in \mathbb{R} \mapsto u$  is not bounded.

The complete symmetry between the two ensembles as indicated by properties 1, 2, and 3 is spoiled by property 4. Although one can obtain  $\varphi(\beta)$  from  $s(u)$  via a Legendre-Fenchel transform, in general one cannot obtain  $s(u)$  from  $\varphi(\beta)$  via the dual formula unless  $s$  is concave on  $\mathbb{R}$ . The concavity of  $s$  on  $\mathbb{R}$  depends on the nature of  $I$  and  $\tilde{H}$ . For example, if  $I$  is convex on  $\mathcal{X}$  and  $\tilde{H}$  is affine, then  $s$  is concave on  $\mathbb{R}$ . Because of the local mean-field, long-range nature of the Hamiltonians arising in many models of coherent structures in turbulence, the associated microcanonical entropies are typically not concave on subsets of  $\mathbb{R}$  corresponding to a range of negative temperatures. This discussion indicates that of the two thermodynamic functions, the microcanonical entropy is the more fundamental, a state of affairs that is reinforced by the results on ensemble equivalence and nonequivalence to be presented in Theorem 10.6.

In order to state this theorem, we need several definitions. A function  $f$  on  $\mathbb{R}$  is said to be concave on  $\mathbb{R}$ , or concave, if  $-f$  is a proper convex function in the sense of [71, p. 24]; that is,  $f$  maps  $\mathbb{R}$  into  $\mathbb{R} \cup \{-\infty\}$ ,  $f \not\equiv -\infty$ , and for all  $u$  and  $v$  in  $\mathbb{R}$  and all  $\lambda \in (0, 1)$

$$f(\lambda u + (1 - \lambda)v) \geq \lambda f(u) + (1 - \lambda)f(v).$$

Given  $f \not\equiv -\infty$  a function mapping  $\mathbb{R}$  into  $\mathbb{R} \cup \{-\infty\}$ , we define  $\text{dom } f$  to be the set of  $u \in \mathbb{R}$  for which  $f(u) > -\infty$ . Let  $\beta$  be a point in  $\mathbb{R}$ . The function  $f$  is said to have a supporting line at  $u \in \text{dom } f$  with tangent  $\beta$  if

$$f(v) \leq f(u) + \beta(v - u) \text{ for all } v \in \mathbb{R}.$$

It follows from this inequality that  $u \in \text{dom } f$ . In addition,  $f$  is said to have a strictly supporting line at  $u \in \text{dom } f$  with tangent  $\beta$  if the inequality in the last display is strict for all  $v \neq u$ .

Let  $f \not\equiv -\infty$  be a function mapping  $\mathbb{R}$  into  $\mathbb{R} \cup \{-\infty\}$ . For  $\beta$  and  $u$  in  $\mathbb{R}$  the Legendre-Fenchel transforms  $f^*$  and  $f^{**}$  are defined by [71, p. 308]

$$f^*(\beta) = \inf_{u \in \mathbb{R}} \{\beta u - f(u)\} \text{ and } f^{**}(u) = \inf_{\beta \in \mathbb{R}} \{\beta u - f^*(\beta)\}.$$

As in the case of convex functions [33, Thm. VI.5.3],  $f^*$  is concave and upper semicontinuous on  $\mathbb{R}$ , and for all  $u \in \mathbb{R}$  we have  $f^{**}(u) = f(u)$  if and only if  $f$  is concave and upper semicontinuous on  $\mathbb{R}$ . If  $f$  is not concave and upper semicontinuous on  $\mathbb{R}$ , then  $f^{**}$  is the smallest concave, upper semicontinuous function on  $\mathbb{R}$  that satisfies  $f^{**}(u) \geq f(u)$  for all  $u \in \mathbb{R}$  [15, Prop. A.2]. In particular, if for some  $u$ ,  $f(u) \neq f^{**}(u)$ , then  $f(u) < f^{**}(u)$ .

Let  $f \not\equiv -\infty$  be a function mapping  $\mathbb{R}$  into  $\mathbb{R} \cup \{-\infty\}$ ,  $u$  a point in  $\text{dom } f$ , and  $K$  a convex subset of  $\text{dom } f$ . The first three of the following four definitions are reasonable because  $f^{**}$  is concave on  $\mathbb{R}$ .

- $f$  is concave at  $u$  if  $f(u) = f^{**}(u)$ .
- $f$  is not concave at  $u$  if  $f(u) < f^{**}(u)$ .
- $f$  is concave on  $K$  if  $f$  is concave at all  $u \in K$ .
- $f$  is strictly concave on  $K$  if for all  $u \neq v$  in  $K$  and all  $\lambda \in (0, 1)$

$$f(\lambda u + (1 - \lambda)v) > \lambda f(u) + (1 - \lambda)f(v).$$

We now state the main theorem concerning the equivalence and nonequivalence of the microcanonical and canonical ensembles. According to part (d),

canonical equilibrium macrostates are always realized microcanonically. However, according to parts (a)–(c), the converse in general is false. The three possibilities given in parts (a)–(c) depend on support and concavity properties of the microcanonical entropy  $s(u)$ .

**Theorem 10.6.** *In parts (a), (b), and (c),  $u$  denotes any point in  $\text{dom } s$ .*

(a) **Full equivalence.** *There exists  $\beta \in \mathbb{R}$  such that  $\mathcal{E}^u = \mathcal{E}_\beta$  if and only if  $s$  has a strictly supporting line at  $u$  with tangent  $\beta$ ; i.e.,*

$$s(v) < s(u) + \beta(v - u) \text{ for all } v \neq u.$$

(b) **Partial equivalence.** *There exists  $\beta \in \mathbb{R}$  such that  $\mathcal{E}^u \subset \mathcal{E}_\beta$  but  $\mathcal{E}^u \neq \mathcal{E}_\beta$  if and only if  $s$  has a nonstrictly supporting line at  $u$  with tangent  $\beta$ ; i.e.,*

$$s(v) \leq s(u) + \beta(v - u) \text{ for all } v \text{ with equality for some } v \neq u.$$

(c) **Nonequivalence.** *For all  $\beta \in \mathbb{R}$ ,  $\mathcal{E}^u \cap \mathcal{E}_\beta = \emptyset$  if and only if  $s$  has no supporting line at  $u$ ; i.e.,*

$$\text{for all } \beta \in \mathbb{R} \text{ there exists } v \text{ such that } s(v) > s(u) + \beta(v - u).$$

*Except possibly for boundary points of  $\text{dom } s$ , the latter condition is equivalent to the nonconcavity of  $s$  at  $u$  [Thm. A.5(c)].*

(d) **Canonical is always realized microcanonically.** *We define  $\tilde{H}(\mathcal{E}_\beta)$  to be the set of  $u \in \mathbb{R}$  having the form  $u = \tilde{H}(x)$  for some  $x \in \mathcal{E}_\beta$ . Then for any  $\beta \in \mathbb{R}$  we have  $\tilde{H}(\mathcal{E}_\beta) \subset \text{dom } s$  and*

$$\mathcal{E}_\beta = \bigcup_{u \in \tilde{H}(\mathcal{E}_\beta)} \mathcal{E}^u.$$

Here are two useful criteria for full or partial equivalence of ensembles.

- **Full or partial equivalence.** Except for boundary points of  $\text{dom } s$ ,  $s$  has a supporting line at  $u \in \text{dom } s$  if and only if  $s$  is concave at  $u$  [15, Thm. A.5(c)], and thus according to parts (a) and (b) of the next theorem, full or partial equivalence of ensembles holds.

- **Full equivalence.** Assume that  $\text{dom } s$  is a nonempty interval and that  $s$  is strictly concave on the interior of  $\text{dom } s$  and continuous on  $\text{dom } s$ . Then except for boundary points of  $\text{dom } s$ ,  $s$  has a strictly supporting line at all  $u \in \text{dom } s$ , and thus according to part (a) of the theorem, full equivalence of ensembles holds.

The reader is referred to [36, §4] for the proof of Theorem 10.6. A partial proof of the equality in part (d) is easily provided. Indeed, if  $x \in \mathcal{E}_\beta$ , then  $x$  minimizes  $I + \beta\tilde{H}$  over  $\mathcal{X}$ . Therefore  $x$  minimizes  $I + \beta\tilde{H}$  over the subset of  $\mathcal{X}$  consisting of all  $x$  satisfying  $\tilde{H}(x) = u$ . It follows that  $x$  minimizes  $I$  over  $\mathcal{X}$  subject to the constraint  $\tilde{H}(x) = u$  and thus that  $x \in \mathcal{E}^{\tilde{H}(x)}$ . We conclude that  $\mathcal{E}_\beta \subset \cup_{u \in \tilde{H}(\mathcal{E}_\beta)} \mathcal{E}^u$ , which is half of the assertion in part (d).

The various possibilities in parts (a), (b), and (c) are illustrated in [43] for the mean-field Blume-Emery-Griffiths spin model. In [37] the theory is applied to a model of coherent structures in two-dimensional turbulence. Numerical computations implemented for geostrophic turbulence over topography in a zonal channel demonstrate that nonequivalence of ensembles occurs over a wide range of the model parameters and that physically interesting equilibria seen microcanonically are often omitted by the canonical ensemble. The coherent structures observed in the model resemble the coherent structures observed in the mid-latitude, zone-belt domains on Jupiter.

In [15] we extend the theory developed in [36] and summarized in Theorem 10.6. In [15] it is shown that when the microcanonical ensemble is nonequivalent with the canonical ensemble on a subset of values of the energy, it is often possible to slightly modify the definition of the canonical ensemble so as to recover equivalence with the microcanonical ensemble. Specifically, we give natural conditions under which one can construct a so-called Gaussian ensemble that is equivalent with the microcanonical ensemble when the canonical ensemble is not. This is potentially useful if one wants to work out the equilibrium properties of a system in the microcanonical ensemble, a notoriously difficult problem because of the equality constraint appearing in the definition of this ensemble. An overview of [15] is given in [16], and in [14] it is applied to the Curie-Weiss-Potts model.

The general large deviation procedure presented in the first part of the present section is applied in the next section to the analysis of two models of coherent structures in two-dimensional turbulence, the Miller-Robert model [61, 62, 69, 70] and a related model due to Turkington [77].

## 11 Maximum Entropy Principles in Two-Dimensional Turbulence

This section presents an overview of work in which Gibbs states are used to predict the large-scale, long-lived order of coherent vortices that persist amid the turbulent fluctuations of the vorticity field in two dimensions [6]. This is done by applying a statistical equilibrium theory of the two-dimensional Euler equations, which govern the motion of an inviscid, incompressible fluid. As shown in [11, 59], these equations are reducible to the vorticity transport equations

$$\frac{\partial \omega}{\partial t} + \frac{\partial \omega}{\partial x_1} \frac{\partial \psi}{\partial x_2} - \frac{\partial \omega}{\partial x_2} \frac{\partial \psi}{\partial x_1} = 0 \quad \text{and} \quad -\Delta \psi = \omega, \quad (11.1)$$

in which  $\omega$  is the vorticity,  $\psi$  is the stream function, and  $\Delta = \partial^2/\partial x_1^2 + \partial^2/\partial x_2^2$  denotes the Laplacian operator on  $\mathbb{R}^2$ . The two-dimensionality of the flow means that these quantities are related to the velocity field  $v = (v_1, v_2, 0)$  according to  $(0, 0, \omega) = \text{curl } v$  and  $v = \text{curl}(0, 0, \psi)$ . All of these fields depend upon the time variable  $t \in [0, \infty)$  and the space variable  $x = (x_1, x_2)$ , which runs through a bounded domain in  $\mathbb{R}^2$ . Throughout this section we assume that this domain equals the unit torus  $T^2 = [0, 1) \times [0, 1)$ , and we impose doubly periodic boundary conditions on all the flow quantities.

The governing equations (11.1) can also be expressed as a single equation for the scalar vorticity field  $\omega = \omega(x, t)$ . The periodicity of the velocity field implies that  $\int_{T^2} \omega \, dx = 0$ . With this restriction on its domain, the Green's operator  $G = (-\Delta)^{-1}$  mapping  $\omega$  into  $\psi$  with  $\int_{\mathcal{X}} \psi \, dx = 0$  is well-defined. More explicitly,  $G$  is the integral operator

$$\psi(x) = G\omega(x) = \int_{\mathcal{X}} g(x - x') \omega(x') \, dx',$$

where  $g$  is the Green's function defined by the Fourier series

$$g(x - x') = \sum_{0 \neq z \in \mathbb{Z}^2} |2\pi z|^{-2} e^{2\pi i \langle z, (x-x') \rangle}.$$

Consequently, (11.1) can be considered as an equation in  $\omega$  alone.



Even though the initial value problem for the equation (11.1) is known to be well-posed for weak solutions whenever the initial data  $\omega^0 = \omega(\cdot, 0)$  belongs to  $L^\infty(\mathcal{X})$  [59], it is well known that this deterministic evolution does not provide a useful description of the system over long time intervals. When one seeks to quantify the long-time behavior of solutions, therefore, one is compelled to shift from the microscopic, or fine-grained, description inherent in  $\omega$  to some kind of macroscopic, or coarse-grained, description. We will make this shift by adopting the perspective of equilibrium statistical mechanics. That is, one views the underlying deterministic dynamics as a means of randomizing the microstate  $\omega$  subject to the conditioning inherent in the conserved quantities for the governing equations (11.1), and one takes the appropriate macrostates to be the canonical Gibbs measures built from these conserved quantities. In doing so, of course, one accepts an ergodic hypothesis that equates the time averages with canonical ensemble averages. Given this hypothesis, one hopes that these macrostates capture the long-lived, large-scale, coherent vortex structures that persist amid the small-scale vorticity fluctuations. The characterization of these self-organized macrostates, which are observed in simulations and physical experiments, is the ultimate goal of the theory.

The models that we will consider build on earlier and simpler theories, the first of which was due to Onsager [66]. Studying point vortices, he predicted that the equilibrium states with high enough energy have a negative temperature and represent large-scale, coherent vortices. This model was further developed in the 1970's, notably by Montgomery and Joyce [63]. However, the point vortex model fails to incorporate all the conserved quantities for two-dimensional ideal flow.

These conserved quantities are the energy, or Hamiltonian functional, and the family of generalized enstrophies, or Casimir functionals [59]. Expressed as a functional of  $\omega$ , the kinetic energy is

$$H(\omega) = \frac{1}{2} \int_{\mathcal{X} \times \mathcal{X}} g(x - x') \omega(x) \omega(x') dx dx'. \quad (11.2)$$

The so-called generalized enstrophies are the global vorticity integrals

$$A(\omega) = \int_{\mathcal{X}} a(\omega(x)) dx,$$

where  $a$  is an arbitrary continuous real function on the range of the vorticity. In terms of these conserved quantities, the canonical ensemble is defined by the formal Gibbs measure

$$P_{\beta,a}(d\omega) = Z(\beta, a)^{-1} \exp[-\beta H(\omega) - A(\omega)] \Pi(d\omega),$$

where  $Z(\beta, a)$  is the associated partition function and  $\Pi(d\omega)$  denotes some invariant product measure on some phase space of all admissible vorticity fields  $\omega$ . Of course, this formal construction is not meaningful as it stands due to the infinite dimensionality of such a phase space. We therefore proceed to define a sequence of lattice models on  $T^2$  in order to give a meaning to this formal construction.

One lattice model that respects conservation of energy and also the generalized enstrophy constraints was developed by Miller et. al. [61, 62] and Robert et. al. [69, 70]; we will refer to it as the Miller-Robert model. A related model, which discretizes the continuum dynamics in a different way, was developed by Turkington [77]. These authors use formal arguments to derive maximum entropy principles that are argued to be equivalent to variational formulas for the equilibrium macrostates. In terms of these macrostates, coherent vortices of two-dimensional turbulence can be studied. The purpose of this section is to outline how the large deviation analysis presented in section 10 can be applied to derive these variational formulas rigorously. References [6] and [77] discuss in detail the physical background.

The variational formulas will be derived for the following lattice model that includes both the Miller-Robert model and the Turkington model as special cases. Let  $T^2$  denote the unit torus  $[0, 1) \times [0, 1)$  with periodic boundary conditions and let  $\mathcal{L}$  be a uniform lattice of  $n = 2^{2m}$  sites  $s$  in  $T^2$ , where  $m$  is a positive integer. The intersite spacing in each coordinate direction is  $2^{-m}$ . We make this particular choice of  $n$  to ensure that the lattices are refined dyadically as  $m$  increases, a property that is needed later when we study the continuum

limit obtained by sending  $n \rightarrow \infty$  along the sequence  $n = 2^{2m}$ . In correspondence with this lattice we have a dyadic partition of  $T^2$  into  $n$  squares called microcells, each having area  $1/n$ . For each  $s \in \mathcal{L}$  we denote by  $M(s)$  the unique microcell having the site  $s$  in its lower left corner. Although  $\mathcal{L}$  and  $M(s)$  depend on  $n$ , this is not indicated in the notation.

The configuration spaces for the lattice model are the product spaces  $\Omega_n = \mathcal{Y}^n$ , where  $\mathcal{Y}$  is a compact set in  $\mathbb{R}$ . Configurations in  $\Omega_n$  are denoted by  $\zeta = \{\zeta(s), s \in \mathcal{L}\}$ , which represents the discretized vorticity field. Let  $\rho$  be a probability measure on  $\mathcal{Y}$  and let  $P_n$  denote the product measure on  $\Omega_n$  with one-dimensional marginals  $\rho$ . As discussed in [6], the Miller-Robert model and the Turkington model differ in their choices of the compact set  $\mathcal{Y}$  and the probability measure  $\rho$ .

For  $\zeta \in \Omega_n$  the Hamiltonian for the lattice model is defined by

$$H_n(\zeta) = \frac{1}{2n^2} \sum_{s, s' \in \mathcal{L}} g_n(s - s') \zeta(s) \zeta(s'),$$

where  $g_n$  is the lattice Green's function defined by the finite Fourier sum

$$g_n(s - s') = \sum_{0 \neq z \in \mathcal{L}^*} |2\pi z|^{-2} e^{2\pi i \langle z, s - s' \rangle}$$

over the finite set  $\mathcal{L}^* = \{z = (z_1, z_2) \in \mathbb{Z}^2 : -2^{m-1} < z_1, z_2 \leq 2^{m-1}\}$ . Let  $a$  be any continuous function mapping  $\mathcal{Y}$  into  $\mathbb{R}$ . For  $\zeta \in \Omega_n$  we also define functions known as the generalized enstrophies by

$$A_{n,a}(\zeta) = \frac{1}{n} \sum_{s \in \mathcal{L}} a(\zeta(s)),$$

In terms of these quantities we define the partition function

$$Z_n(\beta, a) = \int_{\Omega_n} \exp[-\beta H_n(\zeta) - A_{n,a}(\zeta)] P_n(d\zeta)$$

and the canonical ensemble  $P_{n,\beta,a}$ , which is the probability measure that assigns to a Borel subset  $B$  of  $\Omega_n$  the probability

$$P_{n,\beta,a}\{B\} = \frac{1}{Z_n(\beta, a)} \int_B \exp[-\beta H_n(\zeta) - A_{n,a}(\zeta)] P_n(d\zeta). \quad (11.3)$$

These probability measures are parametrized by the constant  $\beta \in \mathbb{R}$  and the function  $a \in \mathcal{C}(\mathcal{Y})$ . The dependence of Gibbs measures on the inverse temperature  $\beta$  is standard, while their dependence on the function  $a$  that determines the enstrophy functional is a novelty of this particular statistical equilibrium problem. The Miller-Robert model and the Turkington model also differ in their choices of the parameter  $\beta$  and the function  $a$ .

The main theorem in this section applies the theory of large deviations to derive the continuum limit  $n \rightarrow \infty$  of the lattice model just introduced. Because the interactions  $g_n(s - s')$  in the lattice model are long-range, one must replace  $\beta$  and  $a$  by  $n\beta$  and  $na$  in order to obtain a nontrivial continuum limit [6, 61, 62]. Replacing  $\beta$  and  $a$  by  $n\beta$  and  $na$  in the formulas for the partition function and the Gibbs state is equivalent to replacing  $H_n$  and  $A_n$  by  $nH_n$  and  $nA_n$  and leaving  $\beta$  and  $a$  unscaled. We carry out the large deviation analysis of the lattice model by applying the general procedure specified in the preceding section, making the straightforward modifications necessary to handle both the Hamiltonian and the generalized enstrophy. Thus we seek a space of macrostates, a sequence of macroscopic variables  $Y_n$ , representation functions  $\tilde{H}$  and  $\tilde{A}_a$  for the Hamiltonian and for the generalized enstrophy, and a large deviation principle for  $Y_n$  with respect to the product measures  $P_n$ . The first marginal of a probability measure  $\mu$  on  $T^2 \times \mathcal{Y}$  is defined to be the probability measure  $\mu_1\{A\} = \mu\{A \times \mathcal{Y}\}$  for Borel subsets  $A$  of  $T^2$ .

- **Space of macrostates.** This is the space  $\mathcal{P}_\theta(T^2 \times \mathcal{Y})$  of probability measures on  $T^2 \times \mathcal{Y}$  with first marginal  $\theta$ , where  $\theta(dx) = dx$  is Lebesgue measure on  $T^2$ .
- **Macroscopic variables.** For each  $n \in \mathbb{N}$ ,  $Y_n$  is the measure-valued function mapping  $\zeta \in \Omega_n$  to  $Y_n(\zeta, dx \times dy) \in \mathcal{P}_\theta(T^2 \times \mathcal{Y})$  defined by

$$Y_n(dx \times dy) = Y_n(\zeta, dx \times dy) = dx \otimes \sum_{s \in \mathcal{L}} 1_{M(s)}(x) \delta_{\zeta(s)}(dy).$$

Thus for Borel subsets  $A$  of  $T^2 \times \mathcal{Y}$

$$Y_n\{A\} = \sum_{s \in \mathcal{L}} \int_A 1_{M(s)}(x) dx \delta_{\zeta(s)}(dy).$$

Since  $\sum_{s \in \mathcal{L}} 1_{M(s)}(x) = 1$  for all  $x \in T^2$ , the first marginal of  $Y_n$  equals  $dx$ .

- **Hamiltonian representation function.**  $\tilde{H} : \mathcal{P}_\theta(T^2 \times \mathcal{Y}) \mapsto \mathbb{R}$  is defined by

$$\tilde{H}(\mu) = \frac{1}{2} \int_{(T^2 \times \mathcal{Y})^2} g(x - x') y y' \mu(dx \times dy) \mu(dx' \times dy'),$$

where

$$g(x - x') = \sum_{0 \neq z \in \mathbb{Z}^2} |2\pi z|^{-2} \exp[2\pi i \langle z, x - x' \rangle].$$

As proved in [6, Lem. 4.4],  $\tilde{H}$  is bounded and continuous and there exists  $C < \infty$  such that

$$\sup_{\zeta \in \Omega_n} |H_n(\zeta) - \tilde{H}(Y_n(\zeta, \cdot))| \leq C \left( \frac{\log n}{n} \right)^{1/2} \quad \text{for all } n \in \mathbb{N}. \quad (11.4)$$

- **Generalized entrophy representation function.**  $\tilde{A}_a : \mathcal{P}_\theta(T^2 \times \mathcal{Y}) \mapsto \mathbb{R}$  is defined by

$$\tilde{A}_a(\mu) = \int_{T^2 \times \mathcal{Y}} a(y) \mu(dx \times dy).$$

$\tilde{A}_a$  is bounded and continuous and

$$A_{n,a}(\zeta) = \tilde{A}_a(Y_n(\zeta, \cdot)) \quad \text{for all } \zeta \in \Omega_n. \quad (11.5)$$

- **Large deviation principle for  $Y_n$ .** With respect to the product measures  $P_n$ ,  $Y_n$  satisfies the large deviation principle on  $\mathcal{P}_\theta(T^2 \times \mathcal{Y})$  with rate function the relative entropy

$$I_{\theta \times \rho}(\mu) = \begin{cases} \int_{T^2 \times \mathcal{Y}} \left( \log \frac{d\mu}{d(\theta \times \rho)} \right) d\mu & \text{if } \mu \ll \theta \times \rho \\ \infty & \text{otherwise.} \end{cases}$$

We first comment on the last item. The large deviation principle for  $Y_n$  with respect to  $P_n$  is far from obvious and in fact is one of the main contributions of [6]. We will address this issue after specifying the large deviation behavior

of the model in Theorem 11.1. Concerning (11.4), since  $\theta\{M(s)\} = 1/n$ , it is plausible that

$$\tilde{H}(Y_n(\zeta, \cdot)) = \frac{1}{2} \sum_{s, s' \in \mathcal{L}} \int_{M(s) \times M(s')} g(x - x') dx dx' \zeta(s) \zeta(s')$$

is a good approximation to  $H_n(\zeta) = [1/(2n^2)] \sum_{s, s' \in \mathcal{L}} g_n(s - s') \zeta(s) \zeta(s')$ . Concerning (11.5), for  $\zeta \in \Omega_n$  we have

$$\tilde{A}_a(Y_n(\zeta, \cdot)) = \int_{T^2 \times \mathcal{Y}} a(y) Y_n(\zeta, dx \times dy) = \frac{1}{n} \sum_{s \in \mathcal{L}} a(\zeta(s)) = A_{n,a}(\zeta).$$

The proofs of the boundedness and continuity of  $\tilde{A}_a$  are straightforward.

Part (a) of Theorem 11.1 gives the asymptotic behavior of the scaled partition functions  $Z_n(n\beta, na)$ , and part (b) states the large deviation principle for  $Y_n$  with respect to the scaled canonical ensemble  $P_{n,n\beta,na}$ . The rate function has the familiar form

$$I_{\beta,a} = I_{\rho \times \theta} + \beta \tilde{H} + \tilde{A} - \varphi(\beta, a),$$

where  $\varphi(\beta, a)$  denotes the canonical free energy. In the formula for  $I_{\beta,a}$  the relative entropy  $I_{\rho \times \theta}$  arises from the large deviation principle for  $Y_n$  with respect to  $P_n$ , and the other terms arise from (11.4), (11.5), and the form of  $P_{n,n\beta,na}$ . Part (c) of the theorem gives properties of the set  $\mathcal{E}_{\beta,a}$  of equilibrium macrostates.  $\mathcal{E}_{\beta,a}$  consists of measures  $\mu$  at which the rate function  $I_{\beta,a}$  in part (b) attains its infimum of 0 over  $\mathcal{P}_\theta(T^2 \times \mathcal{Y})$ . The proof of the theorem is omitted since it is similar to the proof of Theorem 10.4, which adapts Theorems 10.2 and 10.3 to the setting of models of coherent structures in turbulence.

**Theorem 11.1.** *For each  $\beta \in \mathbb{R}$  and  $a \in \mathcal{C}(\mathcal{Y})$  the following conclusions hold.*

(a) *The canonical free energy  $\varphi(\beta, a) = -\lim_{n \rightarrow \infty} \frac{1}{n} \log Z_n(n\beta, na)$  exists and is given by the variational formula*

$$\varphi(\beta, a) = \inf_{\mu \in \mathcal{P}_\theta(T^2 \times \mathcal{Y})} \{\beta \tilde{H}(\mu) + \tilde{A}_a(\mu) + I_{\rho \times \theta}(\mu)\}.$$

(b) *With respect to the scaled canonical ensemble  $P_{n,n\beta,na}$  defined in (11.3),  $Y_n$  satisfies the large deviation principle on  $\mathcal{P}_\theta(T^2 \times \mathcal{Y})$  with scaling constants*

*n* and rate function

$$I_{\beta,a}(\mu) = I_{\rho \times \theta}(\mu) + \beta \tilde{H}(\mu) + \tilde{A}_a(\mu) - \varphi(\beta, a).$$

(c) We define the set of equilibrium macrostates

$$\mathcal{E}_{\beta,a} = \{\mu \in \mathcal{P}_\theta(T^2 \times \mathcal{Y}) : I_{\beta,a}(\mu) = 0\}.$$

Then  $\mathcal{E}_{\beta,a}$  is a nonempty, compact subset of  $\mathcal{P}_\theta(T^2 \times \mathcal{Y})$ . In addition, if  $A$  is a Borel subset of  $\mathcal{P}_\theta(T^2 \times \mathcal{Y})$  such that  $\bar{A} \cap \mathcal{E}_{\beta,a} = \emptyset$ , then  $I_{\beta,a}(\bar{A}) > 0$  and for some  $C < \infty$

$$P_{n,\beta,a}\{Y_n \in A\} \leq \exp[-nI_{\beta,a}(\bar{A})/2] \rightarrow 0 \text{ as } n \rightarrow \infty.$$

In section 3 of [6] we discuss the physical implications of the theorem and the relationship between the following concepts in the context of the Miller-Robert model and the Turkington model:  $\mu \in \mathcal{P}_\theta(T^2 \times \mathcal{Y})$  is a canonical equilibrium macrostate (i.e.,  $\mu \in \mathcal{E}_{\beta,a}$ ) and  $\mu$  satisfies a corresponding maximum entropy principle. In the Miller-Robert model, the maximum entropy principle takes the form of minimizing the relative entropy  $I_{\theta \times \rho}(\mu)$  over  $\mu \in \mathcal{P}_\theta(T^2 \times \mathcal{Y})$  subject to the constraints

$$\tilde{H}(\mu) = H(\omega^0) \quad \text{and} \quad \int_{T^2} \mu(dx \times \cdot) = \int_{T^2} \delta_{\omega^0(x)}(\cdot) dx,$$

where  $\omega^0$  is an initial vorticity field and  $H(\omega^0)$  is defined in (11.2). By analogy with our work in the preceding section, this constrained minimization problem defines the set of equilibrium macrostates with respect to the microcanonical ensemble for the Miller-Robert model. The fact that each  $\mu \in \mathcal{E}_{\beta,a}$  is also a microcanonical equilibrium macrostate is a consequence of part (d) of Theorem 10.6 adapted to handle both the Hamiltonian and the generalized enstrophy. In the Turkington model, the maximum entropy principle takes a somewhat related form in which the second constraint appearing in the Miller-Robert maximum entropy principle is relaxed to a family of convex inequalities parametrized by points in  $\mathcal{Y}$ . Understanding for each model the relationship between equilibrium macrostates  $\mu$  and the corresponding maximum entropy principle allows one to identify a steady vortex flow with a given equilibrium macrostate  $\mu$ . Through

this identification, which is described in [6], one demonstrates how the equilibrium macrostates capture the long-lived, large-scale, coherent structures that persist amid the small-scale vorticity fluctuations.

We spend the rest of this section outlining how the large deviation principle is proved for the macroscopic variables

$$Y_n(dx \times dy) = dx \otimes \sum_{s \in \mathcal{L}} 1_{M(s)}(x) \delta_{\zeta(s)}(dy)$$

with respect to the product measures  $P_n$ . The proof is based on the innovative technique of approximating  $Y_n$  by a doubly indexed sequence of random measures  $W_{n,r}$  for which the large deviation principle is, at least formally, almost obvious. This doubly indexed sequence, obtained from  $Y_n$  by averaging over an intermediate scale, clarifies the physical basis of the large deviation principle and reflects the multiscale nature of turbulence. A similar large deviation principle is derived in [60, 68] by an abstract approach that relies on a convex analysis argument. That approach obscures the role of spatial coarse-graining in the large deviation behavior.

In order to define  $W_{n,r}$ , we recall that  $\mathcal{L}$  contains  $n = 2^{2m}$  sites  $s$ . For even  $r < 2m$  we consider a regular dyadic partition of  $T^2$  into  $2^r$  macrocells  $\{D_{r,k}, k = 1, 2, \dots, 2^r\}$ . Each macrocell contains  $n/2^r$  lattice sites and is the union of  $n/2^r$  microcells  $M(s)$ , where  $M(s)$  contains the site  $s$  in its lower left corner. We now define

$$W_{n,r}(dx \times dy) = W_{n,r}(\zeta, dx \times dy) = dx \otimes \sum_{k=1}^{2^r} 1_{D_{r,k}}(x) \frac{1}{n/2^r} \sum_{s \in D_{r,k}} \delta_{\zeta(s)}(dy).$$

$W_{n,r}$  is obtained from  $Y_n$  by replacing, for each  $s \in D_{r,k}$ , the point mass  $\delta_{\zeta(s)}$  by the average  $(n/2^r)^{-1} \sum_{s \in D_{r,k}} \delta_{\zeta(s)}$  over the  $n/2^r$  sites contained in  $D_{r,k}$ .

We need the key fact that with respect to a suitable metric  $d$  on  $\mathcal{P}_\theta(T^2 \times \mathcal{Y})$ ,  $d(Y_n, W_{n,r}) \leq \sqrt{2}/2^{r/2}$  for all  $n = 2^{2m}$  and all even  $r \in \mathbb{N}$  satisfying  $r < 2m$ . The proof of this approximation property uses the fact that the diameter of each macrocell  $D_{r,k}$  equals  $\sqrt{2}/2^{r/2}$  [6, Lem. 4.2]. The next theorem states the two-parameter large deviation principle for  $W_{n,r}$  with respect to the product measures  $P_n$ . The approximation property  $d(Y_n, W_{n,r}) \leq \sqrt{2}/2^{r/2}$  implies that with



respect to  $P_n$ ,  $Y_n$  satisfies the Laplace principle, and thus the equivalent large deviation principle, with the same rate function  $I_{\theta \times \rho}$  [6, Lem. 4.3]. Subtleties involved in invoking the Laplace principle are discussed in the proof of that lemma.

**Theorem 11.2.** *With respect to the product measures  $P_n$ , the sequence  $W_{n,r}$  satisfies the following two-parameter large deviation principle on  $\mathcal{P}_\theta(T^2 \times \mathcal{Y})$  with rate function  $I_{\theta \times \rho}$ : for any closed subset  $F$  of  $\mathcal{P}_\theta(T^2 \times \mathcal{Y})$*

$$\limsup_{r \rightarrow \infty} \limsup_{n \rightarrow \infty} \frac{1}{n} \log P_n \{W_{n,r} \in F\} \leq -I_{\theta \times \rho}(F)$$

and for any open subset  $G$  of  $\mathcal{P}_\theta(T^2 \times \mathcal{Y})$

$$\liminf_{r \rightarrow \infty} \liminf_{n \rightarrow \infty} \frac{1}{n} \log P_n \{W_{n,r} \in G\} \geq -I_{\theta \times \rho}(G).$$

Our purpose in introducing the doubly indexed process  $W_{n,r}$  is the following. The local averaging over the sets  $D_{r,k}$  introduces a spatial scale that is intermediate between the macroscopic scale of the torus  $T^2$  and the microscopic scale of the microcells  $M(s)$ . As a result,  $W_{n,r}$  can be written in the form

$$W_{n,r}(dx \times dy) = dx \otimes \sum_{k=1}^{2^r} 1_{D_{r,k}}(x) L_{n,r,k}(dy), \quad (11.6)$$

where

$$L_{n,r,k}(dy) = L_{n,r,k}(\zeta, dy) = \frac{1}{n/2^r} \sum_{s \in D_{r,k}} \delta_{\zeta(s)}(dy).$$

Since each  $D_{r,k}$  contains  $n/2^r$  lattice sites  $s$ , with respect to  $P_n$  the sequence  $\{L_{n,r,k}, k = 1, \dots, 2^r\}$  is a family of i.i.d. empirical measures. For each  $r$  and each  $k \in \{1, \dots, 2^r\}$  Sanov's Theorem 6.7 implies that as  $n \rightarrow \infty$ ,  $L_{n,r,k}$  satisfies the large deviation principle on  $\mathcal{P}(\mathcal{Y})$  with scaling constants  $n/2^r$  and rate function  $I_\rho$ .

We next motivate the large deviation principle for  $W_{n,r}$  stated in Theorem 11.2. Suppose that  $\mu \in \mathcal{P}_\theta(T^2 \times \mathcal{Y})$  has finite relative entropy with respect to

$\theta \times \rho$  and has the special form

$$\mu(dx \times dy) = dx \otimes \tau(x, dy), \quad \text{where } \tau(x, dy) = \sum_{k=1}^{2^r} 1_{D_{r,k}}(x) \tau_k(dy) \quad (11.7)$$

and  $\tau_1, \dots, \tau_{2^r}$  are probability measures on  $\mathcal{Y}$ . The representation (11.6), Sanov's Theorem, and the independence of  $L_{n,r,1}, \dots, L_{n,r,2^r}$  suggest that

$$\begin{aligned} & \lim_{n \rightarrow \infty} \frac{1}{n} \log P_n \{W_{n,r} \sim \mu\} \\ &= \lim_{n \rightarrow \infty} \frac{1}{n} \log P_n \{L_{n,r,k} \sim \tau_k, k = 1, \dots, 2^r\} \\ &= \frac{1}{2^r} \sum_{k=1}^{2^r} \lim_{n \rightarrow \infty} \frac{1}{n/2^r} \log P_n \{L_{n,r,k} \sim \tau_k\} \\ &\approx -\frac{1}{2^r} \sum_{k=1}^{2^r} I_\rho(\tau_k) = -\int_{T^2} I_\rho(\tau(x, \cdot)) dx \\ &= -\int_{T^2} \int_{\mathcal{Y}} \left( \log \frac{d\tau(x, \cdot)}{d\rho(\cdot)}(y) \right) \tau(x, dy) dx \\ &= -\int_{T^2 \times \mathcal{Y}} \left( \log \frac{d\mu}{d(\theta \times \rho)}(x, y) \right) \mu(dx \times dy) \\ &= -I_{\theta \times \rho}(\mu). \end{aligned}$$

Because of this calculation, the two-parameter large deviation principle for  $W_{n,r}$  with rate function  $I_{\theta \times \rho}$  is certainly plausible, in view of the fact that any measure  $\mu \in \mathcal{P}_\theta(T^2 \times \mathcal{Y})$  can be well approximated, as  $r \rightarrow \infty$ , by a sequence of measures of the form (11.7) [7, Lem. 3.2]. The reader is referred to [6] for an outline of the proof of this two-parameter large deviation principle. The large deviation principle for  $W_{n,r}$  is a special case of a large deviation principle proved in [7] for an extensive class of random measures that includes  $W_{n,r}$  as a special case.

This completes our application of the theory of large deviations to models of two-dimensional turbulence. The asymptotic behavior of these models is stated in Theorem 11.1. One of the main components of the proof is the large deviation principle for the macroscopic variables  $Y_n$ , which in turn follows by approximating  $Y_n$  by the doubly indexed sequence  $W_{n,r}$  and proving the large deviation

principle for this sequence. This proof relies on Sanov's Theorem, which generalizes Boltzmann's 1877 calculation of the asymptotic behavior of multinomial probabilities. Earlier in the paper we used the elementary form of Sanov's Theorem stated in Theorem 3.4 to derive the form of the Gibbs state for the discrete ideal gas and to motivate the version of Cramér's Theorem needed to analyze the Curie-Weiss model [Cor. 6.6]. It is hoped that both the importance of Boltzmann's 1877 calculation and the applicability of the theory of large deviations to problems in statistical mechanics have been amply demonstrated in these lectures. It is also hoped that these lectures will inspire the reader to discover new applications.

## References

- [1] R. R. Bahadur and S. Zabell. Large deviations of the sample mean in general vector spaces. *Ann. Prob.* 7:587–621, 1979.
- [2] J. Barré, D. Mukamel, and S. Ruffo. Ensemble inequivalence in mean-field models of magnetism. T. Dauxois, S. Ruffo, E. Arimondo, M. Wilkens (editors). *Dynamics and Thermodynamics of Systems with Long Interactions*, pp. 45–67. Volume 602 of Lecture Notes in Physics. New York: Springer-Verlag, 2002.
- [3] J. Barré, D. Mukamel, and S. Ruffo. Inequivalence of ensembles in a system with long-range interactions. *Phys. Rev. Lett.* 87:030601, 2001.
- [4] L. Boltzmann. Über die Beziehung zwischen dem zweiten Hauptsatze der mechanischen Wärmetheorie und der Wahrscheinlichkeitsrechnung respektive den Sätzen über das Wärmegleichgewicht (On the relationship between the second law of the mechanical theory of heat and the probability calculus). *Wiener Berichte* 2, no. 76, 373–435, 1877.
- [5] E. P. Borges and C. Tsallis. Negative specific heat in a Lennard-Jones-like gas with long-range interactions. *Physica A* 305:148–151, 2002.
- [6] C. Boucher, R. S. Ellis, and B. Turkington. Derivation of maximum entropy principles in two-dimensional turbulence via large deviations. *J. Stat. Phys.* 98:1235–1278, 2000.
- [7] C. Boucher, R. S. Ellis, and B. Turkington. Spatializing random measures: doubly indexed processes and the large deviation principle. *Ann. Prob.* 27:297–324, 1999. Erratum: *Ann. Prob.* 30:2113, 2002.
- [8] C. Cercignani. *Ludwig Boltzmann: The Man Who Trusted Atoms*. Oxford: Oxford Univ. Press, 1998.
- [9] E. Caglioti, P. L. Lions, C. Marchioro, and M. Pulvirenti. A special class of stationary flows for two-dimensional Euler equations: a statistical mechanical description. *Comm. Math. Phys.* 143:501–525, 1992.

- [10] H. Chernoff. A measure of asymptotic efficiency for tests of a hypothesis based on the sum of observations. *Ann. Math. Stat.* 23:493–507, 1952.
- [11] A. J. Chorin, *Vorticity and Turbulence*, New York: Springer, 1994.
- [12] M. Costeniuc, R. S. Ellis and P. T.-H. Otto. 36 limit theorems for sums of dependent random variables occurring in statistical mechanics. Submitted for publication, 2006.
- [13] M. Costeniuc, R. S. Ellis, and H. Touchette. Complete analysis of phase transitions and ensemble equivalence for the Curie-Weiss-Potts model. *J. Math. Phys.* 46:063301, 25 pages, 2005.
- [14] M. Costeniuc, R. S. Ellis, and H. Touchette. Nonconcave entropies from generalized canonical ensembles. *Phys. Rev. E* 74:010105(R) (4 pages), 2006.
- [15] M. Costeniuc, R. S. Ellis, H. Touchette, and B. Turkington. The generalized canonical ensemble and its universal equivalence with the micro-canonical ensemble. *J. Stat. Phys.* 119:1283-1329, 2005.
- [16] M. Costeniuc, R. S. Ellis, H. Touchette, and B. Turkington. Generalized canonical ensembles and ensemble equivalence. *Phys. Rev. E* 73:026105 (8 pages), 2006.
- [17] H. Cramér. Sur un nouveau théorème-limite de la théorie des probabilités. *Actualités Scientifiques et Industrielles* 736:2–23, 1938. Colloque consacré à la théorie des probabilités, Vol. 3, Hermann, Paris.
- [18] T. Dauxois, P. Holdsworth, and S. Ruffo. Violation of ensemble equivalence in the antiferromagnetic mean-field XY model. *Eur. Phys. J. B* 16:659, 2000.
- [19] T. Dauxois, V. Latora, A. Rapisarda, S. Ruffo, and A. Torcini. The Hamiltonian mean field model: from dynamics to statistical mechanics and back. In T. Dauxois, S. Ruffo, E. Arimondo, and M. Wilkens, editors, *Dynamics*

- and Thermodynamics of Systems with Long-Range Interactions*, volume 602 of *Lecture Notes in Physics*, pp. 458–487, New York: Springer, 2002.
- [20] A. Dembo and O. Zeitouni. *Large Deviations Techniques and Applications*. Second edition. New York: Springer, 1998.
- [21] J.-D. Deuschel and D. W. Stroock. *Large Deviations*. Boston: Academic Press, 1989.
- [22] M. DiBattista, A. Majda, and M. Grote. Meta-stability of equilibrium statistical structures for prototype geophysical flows with damping and driving. *Physica D* 151:271–304, 2000.
- [23] M. DiBattista, A. Majda, and B. Turkington. Prototype geophysical vortex structures via large-scale statistical theory. *Geophys. Astrophys. Fluid Dyn.* 89:235–283, 1998.
- [24] I. H. Dinwoodie and S. L. Zabell. Large deviations for exchangeable random vectors. *Ann. Probab.* 20:1147–1166, 1992.
- [25] R. L. Dobrushin and S. B. Shlosman. Large and moderate deviations in the Ising Model. *Adv. Soviet Math.* 20:91–219, 1994.
- [26] M. D. Donsker and S. R. S. Varadhan. Asymptotic evaluation of certain Markov process expectations for large time, I. *Comm. Pure Appl. Math.* 28:1–47, 1975.
- [27] M. D. Donsker and S. R. S. Varadhan. Asymptotic evaluation of certain Markov process expectations for large time, III. *Comm. Pure Appl. Math.* 29:389–461, 1976.
- [28] M. D. Donsker and S. R. S. Varadhan. Asymptotic evaluation of certain Markov process expectations for large time, IV. *Comm. Pure Appl. Math.* 36:183–212, 1983.
- [29] P. Dupuis and R. S. Ellis. Large deviations for Markov processes with discontinuous statistics, II: random walks. *Probab. Th. Relat. Fields* 91:153–194, 1992.

- [30] P. Dupuis and R. S. Ellis. The large deviation principle for a general class of queueing systems, I. *Trans. Amer. Math. Soc.* 347:2689–2751, 1995.
- [31] P. Dupuis and R. S. Ellis. *A Weak Convergence Approach to the Theory of Large Deviations*. New York: Wiley, 1997.
- [32] P. Dupuis, R. S. Ellis, and A. Weiss. Large deviations for Markov processes with discontinuous statistics, I: general upper bounds. *Ann. Prob.* 19:1280–1297, 1991.
- [33] R. S. Ellis. *Entropy, Large Deviations, and Statistical Mechanics*. New York: Springer, 1985. Reprinted in *Classics of Mathematics* series, 2006.
- [34] R. S. Ellis. Large deviations for a general class of random vectors. *Ann. Prob.* 12:1–12, 1984.
- [35] R. S. Ellis. An overview of the theory of large deviations and applications to statistical mechanics. *Scand. Actuarial J.* No. 1, 97–142, 1995.
- [36] R. S. Ellis, K. Haven, and B. Turkington. Large deviation principles and complete equivalence and nonequivalence results for pure and mixed ensembles. *J. Stat. Phys.* 101:999–1064, 2000.
- [37] R. S. Ellis, K. Haven, and B. Turkington. Nonequivalent statistical equilibrium ensembles and refined stability theorems for most probable flows. *Nonlinearity* 15:239–255, 2002.
- [38] R. S. Ellis, R. Jordan, P. Otto, and B. Turkington. A statistical approach to the asymptotic behavior of a generalized class of nonlinear Schrödinger equations. *Comm. Math. Phys.*, 244:187–208, 2004.
- [39] R. S. Ellis and C. M. Newman. Limit theorems for sums of dependent random variables occurring in statistical mechanics. *Z. Wahrsch. verw. Geb.* 44:117–139, 1978.
- [40] R. S. Ellis and C. M. Newman. The statistics of Curie-Weiss models. *J. Stat. Phys.* 19:149–161, 1978.

- [41] R. S. Ellis, C. M. Newman, and J. S. Rosen. Limit theorems for sums of dependent random variables occurring in statistical mechanics, II: conditioning, multiple phases, and metastability. *Z. Wahrsch. verw. Geb.* 51:153–169, 1980.
- [42] R. S. Ellis, P. Otto, and H. Touchette. Analysis of phase transitions in the mean-field Blume-Emery-Griffiths model. *Ann. Appl. Prob.* 15:2203–2254, 2005.
- [43] R. S. Ellis, H. Touchette, and B. Turkington. Thermodynamic versus statistical nonequivalence of ensembles for the mean-field Blume-Emery-Griffith model. *Physica A* 335:518–538, 2004.
- [44] R. S. Ellis and K. Wang. Limit theorems for the empirical vector of the Curie-Weiss-Potts model. *Stoch. Proc. Appl.* 35:59–79, 1990.
- [45] S. N. Ethier and T. G. Kurtz. *Markov Processes: Characterization and Convergence*. New York: Wiley, 1986.
- [46] W. R. Everdell. *The First Moderns*. Chicago: The University of Chicago Press, 1997.
- [47] G. L. Eyink and H. Spohn. Negative-temperature states and large-scale, long-lived vortices in two-dimensional turbulence. *J. Stat. Phys.* 70:833–886, 1993.
- [48] W. Feller. *An Introduction to Probability Theory and Its Applications*. Vol. I, second edition. New York: Wiley, 1957.
- [49] H. Föllmer and S. Orey. Large deviations for the empirical field of a Gibbs measure. *Ann. Probab.* 16:961–977, 1987.
- [50] J. Gärtner. On large deviations from the invariant measure. *Th. Probab. Appl.* 22:24–39, 1977.
- [51] D. H. E. Gross. Microcanonical thermodynamics and statistical fragmentation of dissipative systems: the topological structure of the  $n$ -body phase space. *Phys. Rep.* 279:119–202, 1997.



- [52] P. Hertel and W. Thirring. A soluble model for a system with negative specific heat. *Ann. Phys. (NY)* 63:520, 1971.
- [53] M. K.-H. Kiessling and J. L. Lebowitz. The micro-canonical point vortex ensemble: beyond equivalence. *Lett. Math. Phys.* 42:43–56, 1997.
- [54] M. K.-H. Kiessling and T. Neukirch. Negative specific heat of a magnetically self-confined plasma torus. *Proc. Natl. Acad. Sci. USA* 100:1510–1514, 2003.
- [55] O. E. Lanford. Entropy and equilibrium states in classical statistical mechanics. In: *Statistical Mechanics and Mathematical Problems*, pp. 1–113. Edited by A. Lenard. *Lecture Notes in Physics* 20. Berlin: Springer, 1973.
- [56] V. Latora, A. Rapisarda, and C. Tsallis. Non-Gaussian equilibrium in a long-range Hamiltonian system. *Phys. Rev. E* 64:056134, 2001.
- [57] D. Lindley. *Boltzmann's Atom: The Great Debate That Launched a Revolution in Physics*. New York: Free Press, 2001.
- [58] D. Lynden-Bell and R. Wood. The gravo-thermal catastrophe in isothermal spheres and the onset of red-giant structure for stellar systems. *Mon. Notic. Roy. Astron. Soc.* 138:495, 1968.
- [59] C. Marchioro and M. Pulvirenti. *Mathematical Theory of Incompressible Nonviscous Fluids*. New York: Springer, 1994.
- [60] J. Michel and R. Robert. Large deviations for Young measures and statistical mechanics of infinite dimensional dynamical systems with conservation law. *Comm. Math. Phys.* 159:195–215, 1994.
- [61] J. Miller. Statistical mechanics of Euler equations in two dimensions. *Phys. Rev. Lett.* 65:2137–2140 (1990).
- [62] J. Miller, P. Weichman and M. C. Cross. Statistical mechanics, Euler's equations, and Jupiter's red spot. *Phys. Rev. A* 45:2328–2359, 1992.

- [63] D. Montgomery and G. Joyce. Statistical mechanics of negative temperature states. *Phys. Fluids* 17:1139–1145, 1974.
- [64] P. Ney. Private communication, 1997.
- [65] S. Olla. Large deviations for Gibbs random fields. *Probab. Th. Rel. Fields* 77:343–359, 1988.
- [66] L. Onsager. Statistical hydrodynamics. *Suppl. Nuovo Cim.* 6:279–287, 1949.
- [67] G. Parisi. *Statistical Field Theory*. Redwood City, CA: Addison-Wesley Publishing Co., 1988.
- [68] R. Robert. Concentration et entropie pour les mesures d’Young. *C. R. Acad. Sci. Paris* 309, Série I:757–760, 1989.
- [69] R. Robert. A maximum-entropy principle for two-dimensional perfect fluid dynamics. *J. Stat. Phys.* 65:531–553, 1991.
- [70] R. Robert and J. Sommeria. Statistical equilibrium states for two-dimensional flows. *J. Fluid Mech.* 229:291–310, 1991.
- [71] R. T. Rockafellar. *Convex Analysis*. Princeton: Princeton Univ. Press, 1970.
- [72] E. Seneta. *Non-Negative Matrices and Markov Chains*. Second edition. New York: Springer, 1981.
- [73] R. A. Smith and T. M. O’Neil. Nonaxisymmetric thermal equilibria of a cylindrically bounded guiding center plasma or discrete vortex system. *Phys. Fluids B* 2:2961–2975, 1990.
- [74] D. W. Stroock. *An Introduction to the Theory of Large Deviations*. New York: Springer, 1984.
- [75] W. Thirring. Systems with negative specific heat. *Z. Physik*, 235:339–352, 1970.

- [76] H. Touchette, R. S. Ellis, and B. Turkington. An introduction to the thermodynamic and macrostate levels of nonequivalent ensembles. *Physica A* 340:138–146, 2004.
- [77] B. Turkington. Statistical equilibrium measures and coherent states in two-dimensional turbulence. *Comm. Pure Appl. Math.* 52:781–809, 1999.
- [78] B. Turkington, A. Majda, K. Haven, and M. DiBattista. Statistical equilibrium predictions of jets and spots on Jupiter. *Proc. Natl. Acad. Sci. USA* 98:12346–12350, 2001.
- [79] A. S. Wightman. Convexity and the notion of equilibrium state in thermodynamics and statistical mechanics. Introduction to R. B. Israel. *Convexity in the Theory of Lattice Gases*. Princeton: Princeton Univ. Press, 1979.