

## SPECIAL INVITED PAPER

### A NEW LOOK AT INDEPENDENCE<sup>1</sup>

BY MICHEL TALAGRAND

*C.N.R.S. and Ohio State University*

The concentration of measure phenomenon in product spaces is a far-reaching abstract generalization of the classical exponential inequalities for sums of independent random variables. We attempt to explain in the simplest possible terms the basic concepts underlying this phenomenon, the basic method to prove concentration inequalities and the meaning of several of the most useful inequalities.

**1. Introduction.** What is the most important theorem of probability? The following statement could be a reasonable candidate:

(1.1) In a long sequence of tossing a fair coin, it is likely that heads will come up nearly half of the time.

This rather imprecise statement could serve as an introduction to the study of laws of large numbers. These are limit theorems. A commonly heard piece of conventional wisdom (which certainly should not be hastily dismissed) asserts, however, that the “age of computing” coming upon us will shift much of the focus of mathematics from the infinite to the discrete. A precise discrete statement of (1.1) is as follows:

Consider an independent sequence of Bernoulli random variables  $(\varepsilon_i)_{i \leq N}$  [i.e.,  $P(\varepsilon_i = 1) = P(\varepsilon_i = -1) = \frac{1}{2}$ ]. Then for all  $t \geq 0$  we have the following inequality [which will be proved in (4.7) below]:

$$(1.2) \quad P\left(\left|\sum_{i \leq N} \varepsilon_i\right| \geq t\right) \leq 2 \exp\left(-\frac{t^2}{2N}\right).$$

To relate (1.2) to (1.1), we simply observe that if  $B_N$  is the number of 1's in the sequence  $(\varepsilon_i)_{i \leq N}$ , then  $\sum_{i \leq N} \varepsilon_i = 2B_N - N$ , so that (1.2) is equivalent to

$$(1.3) \quad P\left(\left|B_N - \frac{N}{2}\right| \geq t\right) \leq 2 \exp\left(\frac{-2t^2}{N}\right).$$

Inequality (1.2) is possibly the simplest occurrence of the concentration of measure phenomenon that will be explored in the present paper. Upon evoking generalizations of (1.2), the words “exponential inequalities” and the names of Chernoff, Bennett, Prokhorov and Hoeffding (and more) come to mind. The generalizations of (1.2) we have in mind, however, require a change

---

Received October 1994.

<sup>1</sup>Based on the Rietz Lecture delivered by the author at the 58th Annual Meeting of the IMS in Montreal, Canada, in July 1995.

AMS 1991 subject classifications. Primary 60E15, 28A35.

Key words and phrases. Exponential inequalities, concentration of measure, product spaces.

of perspective: simply to think of the random variable  $X = \sum_{i \leq N} \varepsilon_i$  as a *function* of the individual variables  $\varepsilon_i$  and to state (1.2) [or rather (1.1)] as

$$(1.4) \quad X \text{ is essentially constant } (= 0).$$

This statement seems pretty offensive, since the fluctuations of  $X$  are of order  $\sqrt{N}$ , which is hardly zero. This impression is misleading and is simply created by the fact that we do not look at  $X$  on the proper scale. As  $X$  can take values as large as  $N$ , this should be the scale at which one should measure  $X$ , in which case (1.4) is indeed true (i.e.,  $X/N$  is essentially zero).

In words, the form of the concentration of measure phenomenon we will study could be stated as follows:

$$(1.5) \quad \begin{array}{l} \text{A random variable that depends (in a "smooth" way)} \\ \text{on the influence of many independent variables (but} \\ \text{not too much on any of them) is essentially constant.} \end{array}$$

This statement will of course be quantified by inequalities such as (1.2). Most of these inequalities will be of exponential type, so another (shameless ...) way to advertise the results of the present paper is by the following:

$$(1.6) \quad \begin{array}{l} \text{A random variable that smoothly depends on the} \\ \text{influence of many independent random variables} \\ \text{satisfies Chernoff-type bounds.} \end{array}$$

It should be self-evident why a statement such as (1.6) is important. Of special interest is the case where the random variable is defined in an indirect or a complicated way and where explicit computations are all but impossible. A typical situation is the case where the random variable is the solution of a (stochastic) optimization problem, in which case it is sometimes rather hard to say anything at all about it. The body of inequalities underlying the imprecise statement (1.6) has by now been applied to a variety of such optimization problems and has in each occurrence improved and streamlined previous results. These problems include, in particular, stochastic versions of famous questions such as bin packing, the traveling salesman problem and, not surprisingly, models for randomness in physics, such as percolation theory and models for disordered matter in statistical mechanics. (Many applications have also been given to more classical areas of probability such as probability in Banach spaces [9] and empirical processes theory [27].) While going through a large number of applications would have been a fair attempt at impressing upon the reader the importance of the present material, I have resisted the temptation. The main reason is that the abstract inequalities that form the core of the paper (and, in particular, the one presented in Section 6) are sufficiently powerful so that once the basic mechanism of their application is understood, this application becomes mostly a routine matter. The two examples presented in Section 6 should be a sufficient illustration. Numerous other applications are presented in [29], and I hope that the reader will be interested enough by the present essay to ask for more and will be immediately at ease while plunging into this considerably more detailed work.

While the topic of giving a meaning to (1.6) has now become almost a theory in itself, it is a rather pleasant fact that the proof of the main results is very simple. However, how can such simply obtained results have such drastic consequences? The answer lies, of course, in using a good point of view. This requires several layers of abstraction. While the key ideas are again very simple once understood, this is not necessarily the case beforehand. Therefore, these ideas will be explained in considerable detail, and I must apologize should I insist too much on trivialities; triviality is apparently in the eye of the beholder [32]. The true motivation for insisting upon the abstract ideas is that it is while pursuing abstract principles that the main discoveries have been made, and thereby this appears to be the best way of fostering further advances.

The idea of concentration of measure (which was discovered by V. Milman) is arguably one of the great ideas of analysis in our times. While its impact on Probability is only a small part of the whole picture, this impact should not be ignored. The present paper represents my best attempt to explain in the simplest way I can achieve what this is all about, without *ever* doing anything technical. Due to this exacting requirement of simplicity (and even more to space limitation), the present work is very far from being a complete account of what is known. (We refer for this to [28–30]). I hope, however, that it will be informative for the casual reader and will even possibly induce him/her to learn more about this ever-fascinating topic.

**2. The Gromov–Milman formulation.** The Gromov–Milman [4, 17] formulation is rather simple and very effective. It is also our first step toward increased abstraction and the opportunity to stress a number of key features.

First of all, to examine (1.5) it will be convenient, in contrast with a long-standing tradition, to specify the underlying probability space. The probabilistic notion of independence is intimately related to the notion of product measure, and product measures will be the focus of our interest.

Consider a probability space  $(\Omega, \Sigma, \mu)$  and a power  $(\Omega^N, P)$ , where  $P = \mu^{\otimes N}$ . One could consider different factors, but it would not truly increase the generality. The coordinate functions are probabilistically independent, and any sequence of probabilistically independent functions can be realized as above. Thus to study (1.5), we will study functions defined on a product of probability spaces provided with a product measure.

How should we define the fact that a function depends smoothly on the argument? A reasonable answer seems to be that a small variation of the argument produces a small change in the value of the function. The most natural way to define a small variation of the argument is to assume that the underlying space is provided with a distance. Fortunately, a product space  $\Omega^N$  is provided with a natural distance: the Hamming distance given by

$$(2.1) \quad d(x, y) = \text{card}\{i \leq N; x_i \neq y_i\},$$

where  $x = (x_i)_{i \leq N}$  and  $y = (y_i)_{i \leq N}$ . This initial success should not hide a basic limitation: unless the factor  $\Omega$  is provided with some kind of structure,

it seems difficult to define a genuinely different distance than (2.1). Much of Sections 6–8 will be devoted to showing how to bypass this limitation.

The basic object in the Gromov–Milman formulation of the concentration of measure phenomenon is a (Polish) metric space  $(X, d)$  provided with a (Borel) probability  $P$ . It is not required here that  $X$  be a product space, so that this formulation is considerably more general than the special case of product spaces. Quite naturally, in view of the preceding discussion, the class of well-behaved functions will be the class of 1-Lipschitz functions, that is, functions  $f$  from  $X$  to  $\mathbb{R}$  that satisfy

$$(2.2) \quad \forall x, y \in X, \quad |f(x) - f(y)| \leq d(x, y),$$

and the object is to find situations where the Lipschitz functions are essentially constant. How do we identify the value of the constant? It turns out that the most convenient choice is through a median  $M_f$  of  $f$ , that is, a number such that

$$P(f \leq M_f) \geq \frac{1}{2}, \quad P(f \geq M_f) \geq \frac{1}{2}.$$

The statement that  $f$  is essentially constant is then quantified by a bound for

$$P(|f - M_f| \geq t),$$

for  $t > 0$ .

Consider the set  $A = \{f \leq M_f\}$ . Thus  $P(A) \geq \frac{1}{2}$ . Consider the set

$$(2.3) \quad A_t = \{x \in X; \inf\{d(x, y); y \in A\} \leq t\} = \{x; d(x, A) \leq t\}.$$

It follows from (2.2) that

$$x \in A_t \Rightarrow f(x) \leq t + M_f,$$

so that

$$P(f > M_f + t) \leq 1 - P(A_t).$$

This simple observation has accomplished a key fact: it has reduced the study of functions  $f$  to the study of sets, which are genuinely simpler objects, and the central concern now is to show that when  $P(A) \geq \frac{1}{2}$ , the “enlargement” of  $A$  defined by (2.3) has probability close to 1. This question (in the setting of product spaces) will be the central objective of the paper. To quantify this phenomenon in their setting, Gromov and Milman introduce the *concentration function*  $\alpha(P, t)$  (depending also on  $X, d$ ) as the smallest number such that

$$P(A) \geq \frac{1}{2} \Rightarrow 1 - P(A_t) \leq \alpha(P, t).$$

The above discussion should then make clear that, for any 1-Lipschitz function  $f$ ,

$$(2.4) \quad P(|f - M_f| > t) \leq 2\alpha(P, t).$$

If we define the Lipschitz constant  $\|f\|_{\text{Lip}}$  of any function  $f$  on  $X$  as the smallest number such that

$$\forall x, y \in X, \quad |f(x) - f(y)| \leq \|f\|_{\text{Lip}} d(x, y),$$

homogeneity and (2.4) imply

$$(2.5) \quad P(|f - M_f| > t) \leq 2\alpha\left(P, \frac{t}{\|f\|_{\text{Lip}}}\right).$$

The point of these definitions is that, in a variety of situations, the function  $\alpha(P, t)$  decreases very fast as  $t$  increases. We will summarize this in a somewhat imprecise manner by the statement that concentration of measure holds in that case. The origin of this terminology is that, whenever one considers a set  $A$  with  $P(A) \geq \frac{1}{2}$ , most of the points of  $X$  are close to  $A$ ; thus  $P$  “concentrates” around each such set  $A$ .

In Section 5, we will prove somewhat more than the following proposition.

**PROPOSITION 2.1.** *If  $X$  is the product of  $N$  probability spaces,  $P$  is a product measure and  $X$  is provided with the Hamming distance  $d$ , the concentration function satisfies*

$$(2.6) \quad \alpha(P, t) \leq 2 \exp\left(-\frac{t^2}{N}\right).$$

In order to compare this with (1.2), on the space  $X = \{-1, 1\}^N$ , we consider the function  $f$  that is the sum of the coordinates and we observe that [when  $X$  is provided with the Hamming distance given by (2.1)]  $\|f\|_{\text{Lip}} = 2$ . Since  $M_f = 0$  by symmetry, combining (2.4) and (2.6) yields

$$(2.7) \quad P\left(\left|\sum_{i \leq N} \varepsilon_i\right| \geq t\right) \leq 4 \exp\left(-\frac{t^2}{4N}\right).$$

This is not quite as good as (1.2), but still captures its main features.

A prime example of a space where concentration of measure holds is the sphere  $S_N$  of  $\mathbb{R}^{N+1}$  equipped with its geodesic distance  $d$  and normalized Haar measure  $Q_N$ . In that case, P. Lévy proved in 1919 that for any (regular) set  $A$  of  $S_N$  we have

$$(2.8) \quad Q_N(A_t) \geq Q_N(C_t),$$

where  $C$  is a cap of the same measure as  $A$ . (This is a true isoperimetric inequality.) It follows, in particular, through a simple computation that

$$\alpha(P_N, t) \leq \left(\frac{\pi}{8}\right)^{1/2} \exp\left(-\frac{(N-1)}{2}t^2\right).$$

Keeping in mind that in (2.6) the diameter of  $X$  is  $N$ , while in (2.8) it is 1, one sees a great similarity between these inequalities. Around 1970, Milman understood that (2.8) was the key to the famous theorem of Dvoretzky on almost Euclidean sections of convex bodies [15, 16]. Subsequently, Milman most vigorously promoted the concept of concentration of measure (see, e.g., [1] for early thoughts about product spaces), and his ideas had a considerable influence. This concept now plays an important role in the local theory of Banach spaces and the dominant role in probability in Banach space. (This

author is, in particular, pleased to acknowledge that his contributions in this direction have their ultimate source in Milman's philosophy.)

More in line with the topic of the present paper is the case where  $X = \mathbb{R}^N$  is provided with the Euclidean distance and where  $P = \gamma_N$  is the canonical Gaussian measure. Thus  $\gamma_N$  is a product measure when each factor is provided with the canonical Gaussian measure  $\gamma_1$  on  $\mathbb{R}$ , of density  $(2\pi)^{-1/2} \exp(-t^2/2)$ . The importance of this situation stems from the fact that all Gaussian measures (such as Wiener measure) can be suitably approximated by  $\gamma_N$  and that inequalities proved for  $\gamma_N$  can rather trivially be transferred to them.

The Gaussian measure on  $\mathbb{R}^N$  can be seen as the limit of the projection of the dilatation of  $Q_M$  by a factor  $\sqrt{M}$  on  $\mathbb{R}^N$  as  $M \rightarrow \infty$ , a fact known as Poincaré's lemma. It can then be deduced from (2.8) that

$$(2.9) \quad \alpha(\gamma_N, t) \leq \int_t^\infty \frac{1}{\sqrt{2\pi}} \exp\left(\frac{-u^2}{2}\right) du \leq \frac{1}{2} \exp\left(\frac{-t^2}{2}\right),$$

a fact of considerable importance [8, 9].

**3. Classical isoperimetry and rearrangements.** Inequalities such as (2.6) will be called concentration inequalities, and it is instructive to discuss the relationship of such an inequality with classical isoperimetry. The most recognized isoperimetric inequality is likely to be the following statement.

(3.1) Of the bodies of a given volume in  $\mathbb{R}^N$ , the Euclidean ball is the one with the smallest surface area.

This formulation needs the notion of surface area, which in the present case can be defined (when  $\partial A$  is smooth enough) as

$$(3.2) \quad \text{Vol}_{N-1}(\partial A) = \lim_{t \rightarrow 0} \frac{\text{Vol}_N(A_t \setminus A)}{t},$$

where  $A_t$  is the set of points within Euclidean distance  $t$  of  $A$ .

As it turns out, (3.1) is equivalent to a lesser-known formulation that does not require the notion of surface area.

(3.3) Among the bodies  $A$  of a given volume in  $\mathbb{R}^N$ , the ones for which  $A_t$  has minimum volume are the Euclidean balls.

It should be clear through (3.2) that (3.3) implies (3.1) as  $t \rightarrow 0$ . Conversely, bounding below  $d \text{Vol}_N(A_t)/dt$  through (3.1) and (3.2) and integrating yield (3.3). The topic of Section 2 connects with (3.3) for the *large* values of  $t$ . This is uninteresting when  $N = 3$ , but it would be disastrous to stop there because our intuition does *not* function beyond the case  $N \leq 3$ .

In the Gaussian case, the statement corresponding to (3.3) is as follows:

(3.4) Among the sets  $A$  of given measure (for  $\gamma_N$ ), the ones for which  $\gamma_N(A_t)$  are minimal are the half-spaces

(cf. [8, 9] and the references therein).

Using this when the half-space is orthogonal to a basis vector yields

$$(3.5) \quad \gamma_N(A) = \gamma_1((-\infty, a]) \Rightarrow \gamma_N(A_t) \geq \gamma_1((-\infty, a + t]),$$

from which (2.9) follows in the case  $a = 0$ .

An inequality such as (3.5) is extremely satisfactory. It is optimal and points to the so-called extremal sets on which equality occurs (here the half-spaces). It apparently is impossible to obtain a very simple proof of (3.5), or indeed of any inequality with the same quality of precision. The only known approach is based on rearrangements. Starting with  $A$ , one constructs a set  $T(A)$  which is somewhat more regular than  $A$ , such that  $\gamma_N(A) = \gamma_N(T(A))$ , while  $\gamma_N(T(A)_t) \leq \gamma_N(A_t)$ . One then iterates the construction in such a way that the iterates “converge” to an extremal set. [See [3] for a proof of (3.5) in this spirit.] This is somewhat delicate. More important, it seems that the method is bound to failure unless the extremal sets have a reasonably simple structure. This does not appear to be the case in a number of situations of crucial interest. Therefore, it is of primary importance to find other methods.

To finish this section, we will describe a result which, while not in the main line of the paper, is connected by several key features. This result is of the same nature as (3.1), but in a setting where it was not obvious how to define “surface area.” The space is  $\Omega = \{-1, 1\}^N$  provided with the uniform measure  $P_N$ .

Given  $x \in \Omega$  and  $i \leq N$ , we define the point  $T_i x$  obtained by changing the sign of the  $i$ th component of  $x$ . Given a subset  $A$  of  $\Omega$  and  $x \in A$ , we define

$$h_A(x) = \text{card}\{i \leq N; T_i(x) \notin A\}.$$

Thus  $h_A(x)$  counts “the number of directions along which one can leave  $A$  from  $x$ .” The following theorem was motivated by a result of Margulis [10].

**THEOREM 3.1** [25]. *For some universal constant  $K$  and all subsets  $A$  of  $\Omega$ , we have*

$$(3.6) \quad \int_A \sqrt{h_A(x)} dP_N(x) \geq \frac{1}{K} P_N(A)(1 - P_N(A)) \sqrt{\log \frac{1}{P_N(A)(1 - P_N(A))}}.$$

The philosophy of the result is that the left-hand side is a measure of the “surface area” of  $A$ .

Thus, (3.6) provides a lower bound for the surface area of  $A$ , given the “volume”  $P_N(A)$  of  $A$ . To understand better the nature of this lower bound, we first state the Gaussian version of (3.1), which follows from (3.5) the way (3.1) follows from (3.3). We have

$$(3.7) \quad \gamma_N(A) = \gamma_1((-\infty, a]) \Rightarrow s_{N-1}(A) \geq \frac{1}{\sqrt{2\pi}} \exp\left(\frac{-a^2}{2}\right),$$

where the “Gaussian surface area”  $s_{N-1}(A)$  is defined as

$$s_{N-1}(A) = \lim_{t \rightarrow 0} t^{-1} \gamma_N(A_t \setminus A).$$

If we remember that, for  $a \leq -1$ ,  $\gamma_1([-\infty, a])$  is of order  $(1/|a|) \exp(-a^2/2)$ , we see that (3.7) implies

$$(3.8) \quad s_{N-1}(A) \geq \frac{1}{K} \gamma_N(A)(1 - \gamma_N(A)) \sqrt{\log \frac{1}{\gamma_N(A)(1 - \gamma_N(A))}}.$$

The similarity between (3.6) and (3.8) is no accident. It arises from the fact that  $(\mathbb{R}^N, \gamma_N)$  is essentially a “quotient” of  $(\{-1, 1\}^{N'}, P_{N'})$  when  $N' \gg N$  (so that isoperimetry in the latter cannot be better than in the former). To see this, we simply observe that when  $M$  is large,  $\gamma_1$  is close to the image of  $P_M$  under the map  $(x_i)_{i \leq M} \rightarrow M^{-1/2} \sum_{i \leq M} x_i$ , by the central limit theorem, so that  $\gamma_N$  is close to an image of  $P_{NM}$ . Thus, (3.6) can be seen as extending some aspects of (3.7).

One important feature of (3.6) (proved by induction over  $N$ ) is that while it provides a bound of the correct order, it avoids the considerably more difficult “extremal” problem of determining the infimum of the left-hand side given  $P_N(A)$ . As already mentioned, this feature is shared by many of the inequalities we will present.

As should be expected from the discussion relating (3.6) and (3.8) and as is easy to see, both sides of (3.6) are of the same order when  $A$  is a set of the type

$$A_{n,k} = \left\{ (x_i)_{i \leq N}; \sum_{i \leq n} x_i \leq k \right\}.$$

An important feature of (3.6) is that it is “dimension independent,” that is, does not depend on  $N$  (a feature already present in the original result of Margulis). Combinatorialists have considered the problem of finding which subsets of  $\{-1, 1\}^N$  have the smallest boundary  $\partial A$  [defined, e.g., as the set of points of  $A$  for which  $h_A(x) > 0$ ], but their measure of the size of  $\partial A$  is simply  $P_N(\partial A)$ . This formulation, however, is *not* dimension independent. In particular, the sets  $A_{n,n/2}$ , for  $n \leq N$ , play essentially the same role with respect to (3.6), and for each of them both sides of (3.6) are of the same order. However, the size of their boundaries, when measured with the “dimension dependent” quantity  $P_N(\partial A)$  is very different, and only  $A_{N,N/2}$  has a boundary of the smallest possible order among all sets of measure about  $\frac{1}{2}$ . This matter of independence of dimension will be a crucial feature of the result of Section 5, where it will be discussed again.

**4. Martingales.** The martingale method has been important in exploring concentration of measure in cases that are not accessible to the rearrangement methods described in the previous section. It is elegant, robust and simple. Even when rearrangement could be used, the martingale method sometimes gives comparable results in a much simpler fashion.

In contrast with the approach of the previous sections, which concentrate on “enlargements” of large sets, the martingale method deals directly with functions. The basic idea is that if  $(\Sigma_i)_{i \leq n}$  is an increasing sequence of



$\sigma$ -algebras, such that  $\Sigma_0$  is trivial, and if  $f$  is  $\Sigma_n$  measurable, then

$$(4.1) \quad f - Ef = \sum_{1 \leq i \leq n} d_i,$$

where  $d_i = E(f|\Sigma_i) - E(f|\Sigma_{i-1})$ . Thus  $(d_i)$  is a martingale difference sequence, that is,  $d_i$  is  $\Sigma_i$ -measurable, and  $E(d_i|\Sigma_{i-1}) = 0$ . The next step is to get bounds on

$$(4.2) \quad P\left(\left|\sum_{i \leq n} d_i\right| \geq t\right).$$

This could be the time to observe that the martingale will give bounds for

$$P(|f - Ef| \geq t),$$

in contrast with (2.4), which involves  $M_f$ . In practice, it is easier to deal with  $Ef$  rather than  $M_f$ . However, it must be pointed out that, under (2.4),

$$|M_f - Ef| \leq E|f - M_f| \leq 2 \int_0^\infty \alpha(P, u) du,$$

so that (2.4) implies

$$P\left(|f - Ef| \geq t + 2 \int_0^\infty \alpha(P, u) du\right) \leq 2\alpha(P, t)$$

and the occurrence of  $M_f$  in (2.4) is only a secondary nuisance when  $\alpha(P, t)$  is very small.

While there is an extensive and deep theory of martingale inequalities, the inequalities required to bound (4.2) are simple martingale adaptations of classical exponential inequalities for sums of independent random variables. Namely, one bounds  $E \exp(\lambda \sum_{i \leq n} d_i)$  in order to use Chebyshev's exponential inequality

$$(4.3) \quad P(Z \geq t) \leq \inf_{\lambda} (\exp(-\lambda t) E \exp \lambda Z).$$

To do this, we observe that

$$(4.4) \quad \begin{aligned} E\left(\exp \lambda \sum_{i \leq n} d_i\right) &= E\left(\left(\exp \lambda \sum_{i \leq n-1} d_i\right) E(\exp \lambda d_n | \Sigma_{n-1})\right) \\ &\leq E\left(\exp \lambda \sum_{i \leq n-1} d_i\right) \|E(\exp \lambda d_n | \Sigma_{n-1})\|_\infty, \end{aligned}$$

so that, by iteration,

$$(4.5) \quad E \exp \lambda \sum_{i \leq n} d_i \leq \prod_{1 \leq i \leq n} \|E(\exp \lambda d_i | \Sigma_{i-1})\|_\infty.$$

The key of the success of this method lies in an efficient control of  $d_i$ . Probably the most important case is when one controls  $\|d_i\|_\infty$ . In that case, it is a simple matter to show that

$$\|E(\exp \lambda d_i | \Sigma_{i-1})\|_\infty \leq \exp \frac{\lambda}{2} \|d_i\|_\infty^2,$$

which, when combined with (4.3) and (4.5), yields

$$(4.6) \quad P\left(\left|\sum_{i \leq n} d_i\right| \geq t\right) \leq 2 \exp\left(-\frac{2t^2}{\sum_{i \leq n} \|d_i\|_\infty^2}\right),$$

a result usually referred to as Azuma's inequality. In the case where  $d_i = a_i \varepsilon_i$  [ $(\varepsilon_i)_{i \leq n}$  are independent Bernoulli random variables], (4.6) specializes to the so-called sub-Gaussian inequality

$$(4.7) \quad P\left(\left|\sum_{i \leq n} a_i \varepsilon_i\right| \geq t\right) \leq 2 \exp\left(-\frac{t^2}{2 \sum_{i \leq n} a_i^2}\right),$$

a very important fact that contains (1.2) as a special case.

The use of martingales in the spirit above was apparently first made by Yurinskii [33] in the case  $f = \|\sum_{i \leq n} Y_i\|$ , where  $Y_i$  are independent Banach space-valued r.v.'s. In this case, taking for  $\Sigma_i$  the  $\sigma$ -algebra generated by  $Y_1, \dots, Y_i$ , the key observation is that  $d_i$  is estimated by

$$d_i \leq \|Y_i\| + E(\|Y_i\| \mid \Sigma_{i-1}).$$

An important step was performed by Maurey [13], who discovered how to use (4.5) in a situation where neither the choice of  $\Sigma_i$  nor the control of  $d_i$  is obvious. The generality of the method was understood by Schechtman [21]. It yields concentration of measure in several important situations (cf. Chapter 1 of the beautiful book [17]).

In more applied fields, (4.6) was used independently by Shamir and Spencer [22] in studying the chromatic number of random graphs, by Rhee and Talagrand [20] in studying stochastic bin packing and the stochastic traveling salesman problem and later by Pastur and Shcherbina [19] in statistical mechanics. Since then, it has literally swept the world (see, e.g., [14]).

For all its qualities, the martingale method has a great drawback: it does not seem to yield results of optimal order in several key situations. In particular, it seems unable to obtain even a weak version of concentration of measure phenomenon in Gaussian space, as described in Section 3, and does not allow the main inequalities of the present paper to be obtained. For this reason, a new method needed to be invented. It will be explained and demonstrated in the rest of the paper.

**5. Approximation by one point.** In this section we will prove (2.6). The reason for the title of the section is that (2.6) means that when  $P(A) \geq \frac{1}{2}$  most points of  $\Omega^N$  belong to  $A_t$  for  $t$  not too large, which in turn means they can be well approximated by at least one point of  $A$ .

Inequality (2.6) can (essentially) be obtained using (4.6). A special case, with identical proof, is obtained in [17]. In fact, given a set  $A$  with  $P(A) \geq \frac{1}{2}$ , it suffices to apply (4.6) to the function  $f(x) = d(x, A)$ , where  $d$  denotes the Hamming distance and where  $\Sigma_i$  is generated by the first  $i$  coordinates. Then

one can show that  $|d_i| \leq 2$ , so that by (4.6),

$$P(|f - Ef| \geq t) \leq 2 \exp\left(-\frac{t^2}{2N}\right).$$

Now, when  $t = Ef$ , the left-hand side is at least  $\frac{1}{2}$ , since  $P(f \leq 0) = P(A) \geq \frac{1}{2}$  and we get  $t = Ef \leq (2N \log 4)^{1/2}$  so that

$$P(f \geq t + (2N \log 4)^{1/2}) \leq 2 \exp\left(-\frac{t^2}{2N}\right),$$

a weak form of (2.6) that is of comparable strength.

Somewhat weaker statements than (2.6) were also discovered independently through a completely different approach in information theory; see, for example, [11]. The reason for which we choose (2.6) to explain our basic approach is simply that because the meaning of what we try to prove is easy to understand, the reader should be better able to concentrate on the mechanism of the proof.

The most natural way to prove (2.6) seems to be by induction over  $N$ . Thus, starting with  $A \subset \Omega^N$ , one should try to define sets in  $\Omega^{N-1}$  to which the induction hypothesis can be applied. These sets will not necessarily be of measure greater than or equal to  $\frac{1}{2}$ , so that it is necessary to use as induction hypothesis a statement valid whatever the value of  $P(A)$ . In view of what is done for martingales, it is natural to try to bound  $E \exp td(x, A)$ , where  $d(x, A)$  is the Hamming distance of  $x$  and  $A$ . [One might object that  $d(x, A)$  need not be measurable, but measurability questions are irrelevant here and will be ignored.] It is remarkable that about the simplest bound one can expect for  $E \exp td(x, A)$  turns out to be suitable.

PROPOSITION 5.1.

$$(5.1) \quad E \exp td(\cdot, A) \leq \frac{1}{P(A)} \exp \frac{t^2 N}{4}.$$

*In particular,*

$$(5.2) \quad P(d(\cdot, A) \geq k) \leq \frac{1}{P(A)} \exp\left(-\frac{k^2}{N}\right).$$

We observe that (5.2) follows from (5.1) by the Chebyshev exponential inequality, and that (5.2) implies (2.6).

The first key feature of the method of proof (which we will simply call the induction method) is that it will reduce the proof of a statement such as (5.1) concerning  $\Omega^N$  to the proof of a statement concerning only functions on  $\Omega$ . Most of the time the proof of this statement is easy; sometimes it is a bit harder, but its very elementary nature ensures success with sufficient effort.

The second key feature is that (as of today) the method of proof has turned out to be almost miraculously sharp in *every* situation. The reasons for this success are not entirely clear at present.

In the present case, the induction method reduces the proof of Proposition 5.1 to the following lemma.

LEMMA 5.2. *Consider a measurable function  $g$  on  $\Omega$ . Then we have*

$$(5.3) \quad \int_{\Omega} \min\left(e^t, \frac{1}{g(\omega)}\right) d\mu(\omega) \int_{\Omega} g(\omega) d\mu(\omega) \leq \exp \frac{t^2}{4}.$$

PROOF. We observe that

$$\min\left(e^t, \frac{1}{g(\omega)}\right) \leq 1 + e^t(1 - g(\omega)),$$

so that the left-hand side of (5.3) is at most

$$a(1 + e^t(1 - a)),$$

where  $a = \int g d\mu$ . The maximum over  $a$  is

$$\left(\frac{e^{t/2} + e^{-t/2}}{2}\right)^2.$$

Now  $(e^u + e^{-u})/2 \leq e^{u^2/2}$ , as is clear from a power series expansion.  $\square$

The proof of Proposition 5.1 goes by induction over  $N$ . The case  $N = 1$  follows from the application of (5.3) to  $g = 1_A$ .

Suppose now that the result has been proved for  $N$ , and let us prove it for  $N + 1$ . Consider  $A \subset \Omega^{N+1} = \Omega^N \times \Omega$ . For  $\omega \in \Omega$ , we set

$$(5.4) \quad A(\omega) = \{x \in \Omega^N; (x, \omega) \in A\}$$

and

$$B = \{x \in \Omega^N; \exists \omega \in \Omega, (x, \omega) \in A\}.$$

With obvious notation, we have

$$d((x, \omega), A) \leq d(x, A(\omega)).$$

Indeed, if  $y \in A(\omega)$ , then  $(y, \omega) \in A$ , and the number of coordinates where  $(y, \omega)$  and  $(x, \omega)$  differ is the number of coordinates where  $x$  and  $y$  differ. Thus, by the induction hypothesis, we have

$$(5.5) \quad \int_{\Omega^N} \exp(td((x, \omega), A)) dP(x) \leq \frac{\exp(t^2 N/4)}{P(A(\omega))}.$$

We also observe that

$$d((x, \omega), A) \leq d(x, B) + 1.$$

Indeed, if  $y \in B$ , then for some  $\omega' \in \Omega$  we have  $(y, \omega') \in A$ , and the numbers of coordinates at which  $(x, \omega)$  and  $(y, \omega')$  differ is at most one more than

the number of coordinates at which  $x$  and  $y$  differ. Thus, by the induction hypothesis, we have

$$\int_{\Omega^N} \exp(td((x, \omega), A)) dP(x) \leq \frac{\exp(t^2 N/4)}{P(B)},$$

and combining this with (5.5) we get

$$\int_{\Omega^N} \exp(td((x, \omega), A)) dP(x) \leq \exp\left(\frac{t^2 N}{4}\right) \min\left(\frac{\exp(t)}{P(B)}, \frac{1}{P(A(\omega))}\right).$$

Integrating with respect to  $\omega$ , we have

$$\begin{aligned} & \int_{\Omega^{N+1}} \exp(td((x, \omega), A)) dP(x) d\mu(\omega) \\ & \leq \exp\left(\frac{t^2 N}{4}\right) \int_{\Omega} \min\left(\frac{e^t}{P(B)}, \frac{1}{P(A(\omega))}\right) d\mu(\omega). \end{aligned}$$

To complete the induction, it suffices to show, by Fubini's theorem, that

$$\int_{\Omega} \min\left(\frac{e^t}{P(B)}, \frac{1}{P(A(\omega))}\right) d\mu(\omega) \leq \frac{\exp(t^4/4)}{P \otimes \mu(A)} = \frac{\exp(t^4/4)}{\int_{\Omega} P(A(\omega)) d\mu(\omega)}.$$

This follows from Lemma 5.2 applied to the function  $g(\omega) = P(A(\omega))/P(B)$ .

One way to express the fundamental difference between the induction method of Proposition 5.1 and the martingale method is that the martingale method looks "forward" while the induction method conditions with respect to the last coordinate and looks "backward," taking full advantage of the fact that the measures obtained after conditioning are identical product measures.

An interesting extension of (5.2) is obtained by allowing a term  $P(A)^\alpha$ ,  $\alpha \geq 1$ , rather than  $P(A)$  in (5.2); that is,

$$(5.6) \quad P(d(\cdot, A) \geq k) \leq \frac{1}{P(A)^\alpha} \exp\left(-\frac{2k^2}{N} \frac{\alpha}{1+\alpha}\right).$$

The proof is similar to that of (5.2), but requires more calculus. The point of (5.6) is that, as  $\alpha \rightarrow \infty$ , we obtain almost the best possible exponent  $-2k^2/N$ . [The claim that this is the best possible follows from the analysis of the situation of (1.2) that will be done in the next section.]

**6. Approximation by many points.** In order to evaluate the result of Section 5, let us analyze a situation equivalent to that of (1.2). We take  $\Omega^N = \{0, 1\}^N$ , provided with the uniform measure, and

$$A = \left\{ x = (x_i)_{i \leq N}; \sum_{i \leq N} x_i \leq \frac{N}{2} \right\}$$

and we assume for simplicity that  $N$  is even.

Consider  $x \in \{0, 1\}^N$ ,  $m = m(x) = \sum_{i \leq N} x_i$  and assume that  $m > N/2$ . We claim that  $d(x, A) = m - N/2$ . To prove that  $d(x, A) \geq m - N/2$ , we observe that the function  $y \rightarrow \sum_{i \leq N} y_i$  is 1-Lipschitz. On the other hand, if  $y \in A$  is

such that for all  $i$ ,  $y_i \leq x_i$  (which we summarize by the statement  $y \leq x$ ), we have  $d(x, y) = m - N/2$ . Thus, if  $k > 0$ ,

$$\{d(\cdot, A) \geq k\} = \left\{x; \sum_{i \leq N} x_i \geq k + \frac{N}{2}\right\}.$$

The central limit theorem shows that, for  $k = t/2\sqrt{N}$ ,

$$P(d(\cdot, A) \geq k) \sim \gamma_1((t, \infty)) \sim \exp - \frac{t^2}{2} \sim \exp - \frac{2k^2}{N}$$

(neglecting polynomial terms in front of the exponential), so that (5.2) is sharp (except for the factor 2 in the exponent). The previous discussion could seem redundant since the derivation of (2.7) from (2.6) already shows that (2.6) (except for the numerical factor in the exponent) is sharp. There is, however, a detail of crucial importance. The definition of  $A_k$  only means that if  $x \in A_k$ , there is *one*  $y$  in  $A$  for which  $d(x, y) \leq k$ . However, in the preceding example *every*  $y$  in  $A$  with  $y \leq x$  satisfies  $d(x, y) = k$ .

Given  $x \in \{0, 1\}^N$ ,  $y \in A$ , it is rather natural to measure the “failure” in approximating  $x$  by  $y$  by the set  $\{i \leq N, x_i \neq y_i\}$  or, equivalently, by the element  $h(x, y)$  of  $\mathbb{R}^N$  such that

$$(6.1) \quad h(x, y)_i = \begin{cases} 0, & \text{if } x_i = y_i, \\ 1, & \text{if } x_i \neq y_i. \end{cases}$$

To take into account the fact that the elements  $y$  of  $A$  that approximate  $x$  well do not “miss” the same coordinates of  $x$ , it is natural to investigate how small an average of points  $h(x, y)$  can be. In the present case it is natural to average over all  $y \leq x$ , with equal weight. This average is clearly equal to

$$h(x) = \frac{m - N/2}{m} x.$$

We now observe that the Euclidean norm  $\|h(x)\|_2$  of  $h(x)$  satisfies

$$\|h(x)\|_2 = \left(m - \frac{N}{2}\right) \frac{1}{\sqrt{m}} \approx \left(m - \frac{N}{2}\right) \sqrt{\frac{2}{N}}$$

since  $m(x) \sim N/2$  (with overwhelming probability). Now (1.2) implies that

$$P\left(\left|m - \frac{N}{2}\right| \geq t\right) \leq 2 \exp(-2t^2),$$

so we get that (essentially)

$$P(\|h(x)\|_2 \geq t) \leq \exp(-t^2).$$

Quite remarkably, the dimension  $N$  has disappeared from this formula. Well, maybe there is some kind of coincidence there, so let us now investigate a more general example, where  $\Omega^N$  is provided with the probability  $P$  such that the law of the coordinates is independent and has expectation  $p$ ,  $0 < p < 1$ . In jargon,

$$P = ((1 - p)\delta_0 + p\delta_1)^{\otimes N}.$$

Assume again for simplicity that  $pN$  is an integer and that  $p \leq \frac{1}{2}$ , and define

$$A = \left\{ x = (x_i)_{i \leq N}; \sum_{i \leq N} x_i \leq pN \right\}.$$

For  $x \in \{0, 1\}^N$ ,  $m = \sum_{i \leq N} x_i$ , we again have  $d(y, A) = m - pN$ . We should observe that (1.2) is now very inaccurate. Indeed, by the classical bounds on the tails of the binomial law [5] we have something like

$$(6.2) \quad P(m - pN \geq t) \leq \exp\left(-\frac{t^2}{2p(1-p)N + \text{smaller term}}\right)$$

(for  $t \leq 2pN$ ), which is much better than (5.2) as  $p \rightarrow 0$ .

On the other hand, proceeding as in the case  $p = \frac{1}{2}$ , we get

$$h(x) = \frac{m - Np}{m} x,$$

so that

$$\|h(x)\|_2 = (m - Np)\sqrt{m} \simeq (m - Np)\sqrt{\frac{p}{N}},$$

and combining with (6.2) yields

$$(6.3) \quad P(\|h(x)\|_2 \geq t) \leq \exp\left(-\frac{t^2}{2(1-p)}\right) \leq \exp\left(-\frac{t^2}{2}\right).$$

Quite remarkably, not only  $N$ , but also  $p$  has vanished from this inequality: it can no longer be an accident, but only a special case of a general fact.

Consider now a probability space  $\Omega$ . For  $x, y \in \Omega^N$ , we define  $h(x, y) \in \mathbb{R}^N$  by

$$h(x, y)_i = \begin{cases} 1, & \text{if } x_i \neq y_i, \\ 0, & \text{if } x_i = y_i. \end{cases}$$

For a subset  $A$  of  $\Omega^N$ , we define

$$U'_A(x) = \{h(x, y); y \in A\} \subset \mathbb{R}^N.$$

Define  $V'_A(x)$  as the convex hull of  $U'_A(x)$ , and define  $f(A, x)$  as the Euclidean distance from zero to  $V'_A(x)$ . Thus  $f(A, x)$  measures "how far  $x$  is from  $A$ ."

Consider a product measure  $P$  on  $\Omega^N$ .

**THEOREM 6.1.** *We have*

$$(6.4) \quad \int \exp \frac{1}{4} f^2(A, x) dP(x) \leq \frac{1}{P(A)}.$$

*In particular,*

$$(6.5) \quad P(f(A, \cdot) \geq t) \leq \frac{1}{P(A)} \exp\left(\frac{-t^2}{4}\right).$$

Compared with (6.3), we observe a loss of a factor of 2 in the exponent. This loss can, however, be almost recovered when one replaces in (6.4) the term  $P(A)^{-1}$  by  $P(A)^{-\alpha}$  [as in (5.6)].

Theorem 6.1 shares several important features with Theorem 3.1. [Somehow I feel that when  $\Omega^N = \{0, 1\}^N$ , provided with the uniform probability, Theorem 6.1 is to Theorem 3.1 what (3.3) is to (3.1), although I do not know how to make this idea precise.] The most important feature is that it is dimension independent so that (in contrast to Proposition 5.1) it is useful to study (e.g.) infinite series.

The key to the proof of Theorem 6.1 is the following lemma. The proof is elementary calculus and need not be reproduced here.

LEMMA 6.2. *Consider  $0 \leq r \leq 1$ . Then*

$$(6.6) \quad \inf_{0 \leq \lambda \leq 1} r^{-\lambda} \exp \frac{(1-\lambda)^2}{4} \leq 2 - r.$$

Before we start the proof of Theorem 6.1, we need an equivalent way to define  $f(A, x)$ . This way is less transparent, but technically more convenient. We set

$$\begin{aligned} U_A(x) &= \{(s_i)_{i \leq N} \in \{0, 1\}^N; \exists y \in A, s_i = 0 \Rightarrow x_i = y_i\} \\ &= \{(s_i)_{i \leq N} \in \{0, 1\}^N; \exists y \in A, \forall i \leq N, s_i \geq h(x, y)_i\}. \end{aligned}$$

For convenience, if  $s_i \geq h(x, y)_i$ , for each  $i \leq N$ , we say that  $y$  witnesses that  $s \in U_A(x)$ . Thus, viewing  $U_A(x)$  as a subset of  $\mathbb{R}^N$ , we have  $U_A(x) \supset U'_A(x)$ . We denote by  $V_A(x)$  the convex hull of  $U_A(x)$ . It should be clear that

$$\forall z \in V_A(x), \quad \exists z' \in U'_A(x), \quad \forall i \leq N, \quad z_i \geq z'_i,$$

so that  $f(A, x)$  is also the distance from 0 to  $V_A(x)$ .

We now prove Theorem 6.1 by induction upon  $N$ . We leave to the reader the easy case  $N = 1$ . For the induction step from  $N$  to  $N+1$ , consider a subset  $A$  of  $\Omega^{N+1}$  and its projection  $B$  on  $\Omega^N$ . For  $\omega \in \Omega$ , we set as usual

$$A(\omega) = \{x \in \Omega^N; (x, \omega) \in A\}.$$

Consider  $x \in \Omega^N$ ,  $\omega \in \Omega$  and  $z = (x, \omega)$ . The basic observation is that

$$\begin{aligned} s \in U_{A(\omega)}(x) &\Rightarrow (s, 0) \in U_A(z), \\ t \in U_B(x) &\Rightarrow (t, 1) \in U_A(x). \end{aligned}$$

For the first claim, if  $y \in A(\omega)$  witnesses that  $s \in U_{A(\omega)}(x)$ , then  $(y, \omega) \in A$  and witnesses that  $(s, 0) \in U_A(z)$ . For the second claim, if  $y \in B$  witnesses that  $t \in U_B(x)$ , then for some  $\omega'$  we have  $(y, \omega') \in A$ , and this point witnesses that  $(t, 1) \in U_A(x)$ . Thus, for  $s \in U_{A(\omega)}(x)$  and  $t \in U_B(x)$ ,  $0 \leq \lambda \leq 1$ , we have  $(\lambda s + (1-\lambda)t, 1-\lambda) \in U_A(z)$ . The convexity of the function  $u \rightarrow u^2$  shows that

$$(6.7) \quad f^2(A, z) \leq (1-\lambda)^2 + \lambda f^2(A(\omega), x) + (1-\lambda) f^2(B, x).$$



The main trick of the proof is to resist the temptation to optimize now over  $\lambda$ . By Hölder's inequality and the induction hypothesis, we have

$$\begin{aligned} & \int \exp \frac{1}{4} f^2(A, (x, \omega)) dP(x) \\ & \leq \exp \frac{1}{4} (1 - \lambda)^2 \left( \int_{\Omega^N} \exp \frac{1}{4} f^2(A(\omega), x) dP(x) \right)^\lambda \\ & \quad \times \left( \int_{\Omega^N} \exp \frac{1}{4} f^2(B, x) dP(x) \right)^{1-\lambda} \\ & \leq \exp \frac{1}{4} (1 - \lambda)^2 \left( \frac{1}{P(A(\omega))} \right)^\lambda \left( \frac{1}{P(B)} \right)^{1-\lambda} \\ & = \frac{1}{P(B)} \exp \frac{1}{4} (1 - \lambda)^2 \left( \frac{P(A(\omega))}{P(B)} \right)^{-\lambda}. \end{aligned}$$

This inequality holds for all  $0 \leq \lambda \leq 1$ . Using (6.6) with  $r = P(A(\omega))/P(B) \leq 1$ , we get

$$\int_{\Omega^N} \exp \frac{1}{4} f^2(A, (x, \omega)) dP(x) \leq \frac{1}{P(B)} \left( 2 - \frac{P(A(\omega))}{P(B)} \right).$$

Integrating with respect to  $\omega$  and using Fubini's theorem yields

$$\begin{aligned} \int \exp \frac{1}{4} f^2(A, \cdot) d(P \otimes \mu) & \leq \frac{1}{P(B)} \left( 2 - \frac{P \otimes \mu(A)}{P(B)} \right) \\ & \leq \frac{1}{P \otimes \mu(A)}, \end{aligned}$$

since  $x(2 - x) \leq 1$  for all  $x$  real.

While Theorem 6.1 turns out to be a principle of considerable power, it takes some effort to realize this. One efficient way to use Theorem 6.1 is through the following observation.

**LEMMA 6.3.** *Consider  $x \in \Omega^N$ . Then given any sequence  $(\alpha_i)_{i \leq N}$ , we can find  $y$  in  $A$  such that*

$$(6.8) \quad \sum_{i \leq N} \{\alpha_i; x_i \neq y_i\} \leq f(A, x) \sqrt{\sum_{i \leq N} \alpha_i^2}.$$

**PROOF.** The linear functional  $\bar{\alpha}: s \rightarrow \sum_{i \leq N} \alpha_i s_i$  on  $\mathbb{R}^N$ , provided with the Euclidean norm, has a norm  $\sqrt{\sum_{i \leq N} \alpha_i^2}$ . Since  $V'_A(x)$  contains a point at distance  $f(A, x)$  to the origin, the infimum of  $\bar{\alpha}$  on  $V'_A(x)$  is at most  $f(A, x) \sqrt{\sum_{i \leq N} \alpha_i^2}$ , but since  $V'_A(x)$  is the convex hull of  $U'_A(x)$ , the infimum of  $\bar{\alpha}$  on  $U'_A(x)$  is also at most  $f(A, x) \sqrt{\sum_{i \leq N} \alpha_i^2}$ , which is the statement of the lemma.  $\square$

Theorem 6.1 has proved efficient in stochastic combinatorial optimization, so we describe a typical application. Consider a sequence  $X_1, \dots, X_N$  of independent r.v.'s, uniformly distributed over  $[0, 1]$ . We are interested in  $L_N = L_N(X_1, \dots, X_N)$ , where  $L_N(X_1, \dots, X_N)$  denotes the longest increasing subsequence of the sequence  $X_1, \dots, X_N$  of numbers. To reduce to the setting of sets in product spaces, we consider  $\Omega = [0, 1]$  and, for  $x = (x_i)_{i \leq N} \in \Omega^N$ , we set  $L(x) = L_N(x_1, \dots, x_N)$ .

For  $a > 0$ , we set

$$A(a) = \{x \in \Omega^N; L_N(x) \leq a\}.$$

The basic observation follows:

LEMMA 6.4. *For all  $x \in \Omega^N$ , we have*

$$(6.9) \quad a \geq L_N(x) - f(A(a), x)\sqrt{L_N(x)}.$$

*In particular,*

$$(6.10) \quad L_N(x) \geq a + v \Rightarrow f(A(a), x) \geq \frac{v}{\sqrt{a+v}}.$$

PROOF. For simplicity, we write  $b = L_N(x)$ . By definition, we can find a subset  $I$  of  $\{1, \dots, N\}$  of cardinality  $b$  such that if  $i, j \in I$ ,  $i < j$ , then  $x_i < x_j$ . By Lemma 6.3 (taking  $\alpha_i = 1$ , if  $i \in I$ , and  $\alpha_i = 0$ , otherwise), there exists  $y \in A(a)$  such that  $\text{card } J \leq f(A(a), x)\sqrt{b}$ , where  $J = \{i \in I; y_i \neq x_i\}$ . Thus  $(x_i)_{i \in I \setminus J}$  is an increasing subsequence of  $y$ . Since  $y \in A(a)$ , we have  $\text{card}(I \setminus J) \leq a$ , which proves (6.9).

To prove (6.10), we observe that by (6.9) we have

$$f(A(a), x) \geq \frac{L_N(x) - a}{\sqrt{L_N(x)}}$$

and that the function  $u \rightarrow (u - a)/\sqrt{u}$  increases for  $u \geq a$ .  $\square$

We denote by  $M$  ( $= M_N$ ) a median of  $L_N$ .

THEOREM 6.5. *For all  $u > 0$ , we have*

$$(6.11) \quad P(L_N \geq M + u) \leq 2 \exp\left(-\frac{u^2}{4(M+u)}\right),$$

$$(6.12) \quad P(L_N \leq M - u) \leq 2 \exp\left(-\frac{u^2}{4M}\right).$$

PROOF. To prove (6.11), we combine (6.12) with  $M = a$  and (6.5). To prove (6.12), we use (6.10) with  $a = M - u$ ,  $v = u$ , to see that

$$L_N(x) \geq M \Rightarrow f(A(M - u), x) \geq \frac{u}{\sqrt{M}},$$

so that

$$(6.13) \quad P\left(f(A(M-u), x) \geq \frac{u}{\sqrt{M}}\right) \geq \frac{1}{2}.$$

On the other hand, by (6.5),

$$(6.14) \quad P\left(f(A(M-u), x) \geq \frac{u}{\sqrt{M}}\right) \leq \frac{1}{P(A(M-u))} \exp\left(-\frac{u^2}{4M}\right).$$

Comparing (6.13) and (6.14) gives the required bound on  $P(A(M-u))$ .  $\square$

It is known (and very easy to see) that  $M_N$  is of order  $\sqrt{N}$ , so that Theorem 6.4 proves that the fluctuations of  $L_N$  are not larger than  $N^{1/4}$ . Simulation [18] suggests, however, that the correct order of magnitude is smaller. Such a phenomenon cannot occur from a deficiency of Theorem 6.1, but rather from the specifics of the situation. We would like to suggest a plausible explanation of what happens.

We conjecture that (in most situations) a random sequence  $(X_1, \dots, X_N)$  has many subsequences of (nearly) maximal length. To see the relevance of this, let us go back to the proof of (6.9). Consider  $b \leq L_N(x)$ . Consider the family  $\mathcal{J}$  of subsets  $I$  of  $\{1, \dots, N\}$  of cardinality  $b$  such that  $i, j \in I$ ,  $i < j$  implies  $x_i < x_j$ . Consider the family  $\mathcal{H}$  of functions on  $\{1, \dots, N\}$  that consists of the indicators of sets of  $\mathcal{J}$ . Consider an element  $(\alpha_i)_{i \leq N}$  in the convex hull of  $\mathcal{H}$  and let  $\sigma = (\sum_{i \leq N} \alpha_i^2)^{1/2}$ . When the family  $\mathcal{J}$  is "rich," we can expect that there is an averaging-out effect and that the sequence  $(\alpha_i)_{i \leq N}$  can be chosen such that  $\sigma^2 \ll b$ . Using Lemma 6.3 we can find  $y$  in  $A$  with

$$\sum_{i \leq N} \{\alpha_i; x_i \neq y_i\} \leq \sigma f(A(a), x).$$

Thus, we can find  $I$  in  $\mathcal{J}$  such that

$$\text{card}\{i \in I; x_i \neq y_i\} \leq \sigma f(A(a), x).$$

As in the proof of (6.9), this shows that  $b - \sigma f(A(a), x) \geq a$ . Thus, if  $b$  is close to  $L(x)$  and  $\sigma^2 \ll b$ , this allows us to improve upon (6.9). Deciding whether the phenomenon described above occurs or not is unrelated to the methods of the present paper and would certainly require a better understanding of the specifics of random sequences.

The reader must have observed that in Lemma 6.4 we do not use the full power of Lemma 6.3; rather, instead of using (6.8) for all sequences of numbers  $(\alpha_i)$ , we used it only for sequences of 0's and 1's. It seems reasonable to assert that Theorem 6.4 uses Theorem 6.1 at the very limit of its area of competence. This can also be seen by the fact that martingale methods can prove an inequality almost as good as (6.11) and (6.12) [2]. By contrast, martingale methods seem powerless to approach the applications where Theorem 6.1 is used at full power, such as in the following theorem.

**THEOREM 6.6.** Consider a real-valued function  $f$  defined on  $[-1, 1]^N$ . We assume that, for each real number  $a$ ,

$$(6.15) \quad \text{the set } \{f \leq a\} \text{ is convex.}$$

Consider a convex set  $B \subset [-1, 1]^N$ , consider  $\sigma > 0$  and assume that the restriction of  $f$  to  $B$  has a Lipschitz constant at most  $\sigma$ ; that is,

$$(6.16) \quad \forall x, y \in B, \quad |f(x) - f(y)| \leq \sigma \|x - y\|,$$

where  $\|x\|$  denotes the Euclidean norm of  $x$ .

Consider independent random variables  $(X_i)_{i \leq N}$  valued in  $[-1, 1]$ , and consider the random variable

$$h = f(X_1, \dots, X_N).$$

Then, if  $M$  is a median of  $h$ , we have, for all  $t > 0$ , that

$$(6.17) \quad P(|h - M| \geq t) \leq 4b + \frac{4}{1 - 2b} \exp\left(-\frac{t^2}{16\sigma^2}\right),$$

where we assume

$$b = P((X_1, \dots, X_N) \notin B) < \frac{1}{2}.$$

Certainly the reader should first consider the special case  $B = [-1, 1]^N$ , where  $b = 0$  and where (6.17) reads

$$(6.18) \quad P(|h - M| \geq t) \leq 4 \exp\left(-\frac{t^2}{16\sigma^2}\right).$$

To understand this inequality better, we will compare it with the Gaussian case. Let us now assume that  $f$  is defined on all  $\mathbb{R}^N$  and has Lipschitz constant  $\sigma$ . Set  $h' = f(Y_1, \dots, Y_N)$ , where the sequence  $Y_1, \dots, Y_N$  is independent standard normal. Combining (2.5) and (2.9), we have

$$(6.19) \quad P(|h' - M'| \geq t) \leq \exp\left(-\frac{t^2}{2\sigma^2}\right),$$

where  $M'$  is of course a median of  $h'$ . Thus, what (6.18) does is to prove an inequality similar to (6.19) for random variables that need no longer be Gaussian (but rather are bounded) and this under only the pretty mild restriction (6.15).

**PROOF OF THEOREM 6.6.** Let us fix  $a \in \mathbb{R}$ , and consider the set  $A(a) = \{f \leq a\} \cap B$ . The key observation is that, for any  $x$  in  $[-1, 1]^N$ , we have

$$(6.20) \quad d(x, A(a)) \leq 2f(A(a), x).$$

Indeed, if  $y \in A(a)$ , we have

$$\forall i \leq N, \quad |x_i - y_i| \leq 2h(x, y)_i$$

[where  $h(x, y)_i$  is defined in (6.1)] because the left-hand side is at most 2 and is zero when the right-hand side is not 2. Thus, for any points  $(y^k)_{k \leq M}$  of  $A(a)$  and convex coefficients  $(\alpha_k)_{k \leq M}$ , we have, for each  $i \leq N$ , that

$$\left| x_i - \sum_k \alpha_k y_i^k \right| \leq 2 \sum_k \alpha_k h(x, y^k)_i,$$

so that, since  $A(a)$  is convex,

$$d(x, A(a)) \leq \left\| x - \sum_k \alpha_k y^k \right\| \leq 2 \left( \sum_{i \leq N} \left( \sum_k \alpha_k h(x, y^k)_i \right)^2 \right)^{1/2},$$

from which (6.20) follows by definition of  $V'_A$ .

Now, if  $x \in B$ , it follows from (6.16) that

$$f(x) \leq a + \sigma d(x, A(a)) \leq a + 2\sigma f(A(a), x).$$

Thus, if we denote by  $P$  the law  $(X_1, \dots, X_N)$  on  $\Omega^N = [-1, 1]^N$ , (6.5) implies

$$P(f \geq a + t) \leq b + \frac{1}{P(A(a))} \exp\left(-\frac{t^2}{16\sigma^2}\right).$$

Taking  $a = M$ , we get  $P(A(a)) \geq \frac{1}{2} - b$ , so that

$$P(f \geq M + t) \leq b + \frac{1}{\frac{1}{2} - b} \exp\left(-\frac{t^2}{16\sigma^2}\right).$$

Taking  $a + t = M$ , we get

$$\frac{1}{2} \leq b + \frac{1}{P(A(M - t))} \exp\left(-\frac{t^2}{16\sigma^2}\right),$$

so that

$$P(A(M - t)) = P(f \leq M - t) \leq 2b + 2 \exp\left(-\frac{t^2}{16\sigma^2}\right). \quad \square$$

COMMENTS. (i) Certainly the reader has observed the similarity of this proof with the proof of Theorem 6.5.

(ii) We have not been very cautious with the coefficients of  $b$ . This is not needed because in applications  $b$  is extremely small.

Here is an important corollary.

**THEOREM 6.7.** *Consider independent (real) random variables  $(X_i)_{i \leq N}$  valued in  $[-1, 1]$  and vectors  $(v_i)_{i \leq N}$  in a Banach space  $Y$ . Define*

$$\sigma^2 = \sup \left\{ \sum_{i \leq N} y^*(v_i)^2; y^* \in Y^*, \|y^*\| \leq 1 \right\},$$

where  $Y^*$  is the dual of  $Y$ . Then, if  $M$  denotes a median of  $\|\sum_{i \leq N} X_i v_i\|$ , we have

$$(6.21) \quad P\left(\left|\left\|\sum_{i \leq N} X_i v_i\right\| - M\right| \geq t\right) \leq 4 \exp\left(-\frac{t^2}{16\sigma^2}\right).$$

REMARK. The most important case is where the r.v.'s  $X_i$  are Bernoulli, that is,  $P(X_i = 1) = P(X_i = -1) = \frac{1}{2}$ .

PROOF. We observe that  $\|\sum_{i \leq N} X_i v_i\| = f(X_1, \dots, X_N)$ , where, for  $x = (x_i)_{i \leq N}$  in  $\mathbb{R}^N$ , we set

$$f(x) = \left\|\sum_{i \leq N} x_i v_i\right\|.$$

By the Hahn–Banach theorem,

$$\left\|\sum_{i \leq N} x_i v_i\right\| = \sup\left\{y^*\left(\sum_{i \leq N} x_i v_i\right); y^* \in Y^*; \|y^*\| \leq 1\right\}.$$

Now, by Cauchy–Schwarz,

$$y^*\left(\sum_{i \leq N} x_i v_i\right) = \sum_{i \leq N} x_i y^*(v_i) \leq \left(\sum_{i \leq N} x_i^2\right)^{1/2} \left(\sum_{i \leq N} y^*(v_i)^2\right)^{1/2} \leq \sigma \|x\|.$$

Thus, by the triangle inequality,

$$|f(x) - f(y)| \leq f(x - y) \leq \sigma \|x - y\|$$

and thus (6.21) is a specialization of (6.18).  $\square$

We now give another application of Theorem 6.6 to the Hopfield model of associative memory [6]. Consider two integers  $M, N$ . For  $x = (x_{i,k})_{i \leq N, k \leq M} \in \mathbb{R}^{MN}$  and for  $\varepsilon = (\varepsilon_i)_{i \leq N} \in \{-1, 1\}^N$ , we set

$$H(x, \varepsilon) = \frac{1}{2N} \sum_{k \leq M} \left(\sum_{i \leq N} x_{i,k} \varepsilon_i\right)^2$$

(the factor  $1/2N$  is customary but unimportant).

Given a subset  $A$  of  $\{-1, 1\}^N$ , we set

$$f(x) = f_N(x) = \frac{1}{\beta} \log\left(\sum_{\varepsilon \in A} \exp \beta H(x, \varepsilon)\right).$$

The quantity of interest is the random variable  $h_N = f_N(\eta)$ , when  $\eta = (\eta_{i,k})_{i \leq N, k \leq M}$  and when  $(\eta_{i,k})_{i \leq N, k \leq M}$  are independent Bernoulli r.v.'s [ $P(\eta_{i,k} = 1) = \frac{1}{2} = P(\eta_{i,k} = -1)$ ]. In the case  $A = \{-1, 1\}^N$ ,  $h_N$  is the free energy of the Hopfield model (at temperature  $T = 1/\beta$ ) and its study is extremely difficult, yet one has the following general result.

**THEOREM 6.8.** Denoting by  $m_N$  a median of  $h_N$ , for some universal constant  $K$  and all  $0 \leq t \leq (N + M)$ , we have

$$P(|h_N - m_N| \geq t) \leq 12 \exp\left(-\frac{t^2}{K(N + M)}\right).$$

**PROOF.** The proof relies on Theorem 6.6, applied to the function  $f$  on  $[-1, 1]^{NM}$ . It is not clear whether  $f$  is convex, but certainly  $\exp \beta f$  is convex, and this implies (6.15). Consider a parameter  $L$  and set

$$B = \{x \in [-1, 1]^{NM}; \forall \varepsilon \in A, H(x, \varepsilon) \leq L\}$$

so that  $B$  is convex. Consider now  $x$  and  $y$  in  $B$ . We try to prove (6.16). We observe that, given  $\varepsilon \in A$ ,

$$2N(H(x, \varepsilon) - H(y, \varepsilon)) = \sum_{k \leq M} \left( \sum_{i \leq N} (x_{i,k} - y_{i,k}) \varepsilon_i \right) \left( \sum_{i \leq N} x_{i,k} \varepsilon_i + \sum_{i \leq N} y_{i,k} \varepsilon_i \right).$$

Thus, by Cauchy–Schwarz,

$$|H(x, \varepsilon) - H(y, \varepsilon)| \leq \frac{1}{2N} UV,$$

where

$$U^2 = \sum_{k \leq M} \left( \sum_{i \leq N} (x_{i,k} - y_{i,k}) \varepsilon_i \right)^2,$$

$$V^2 = \sum_{k \leq M} \left( \sum_{i \leq N} x_{i,k} \varepsilon_i + \sum_{i \leq N} y_{i,k} \varepsilon_i \right)^2.$$

Using the inequality  $(a + b)^2 \leq 2(a^2 + b^2)$ , we see that

$$V^2 \leq 4N(H(x, \varepsilon) + H(y, \varepsilon)) \leq 8NL.$$

Using Cauchy–Schwarz, we see that

$$U^2 \leq N \sum_{\substack{k \leq M \\ i \leq N}} (x_{i,k} - y_{i,k})^2 = N \|x - y\|^2,$$

so that, finally, for each  $\varepsilon$  in  $A$  we have

$$|H(x, \varepsilon) - H(y, \varepsilon)| \leq \|x - y\| \sqrt{2L}.$$

It is then very simple to see that this implies

$$|f(x) - f(y)| \leq \sqrt{2L} \|x - y\|.$$

Thus, by (6.17) we have

$$(6.22) \quad P(|h_N - m_N| \geq t) \leq 4b + \frac{4}{1 - 2b} \exp\left(-\frac{t^2}{32L}\right),$$

where  $b = P(\eta \notin B)$ . To choose  $L$ , we note that by (4.7) we have, for each  $k$ ,

$$E \exp \frac{1}{4N} \left( \sum_{i \leq N} \eta_{i,k} \varepsilon_i \right)^2 \leq K_1,$$

where  $K_1$  is a universal constant, so that, by independence,

$$E \exp \frac{1}{2} H(\eta, \varepsilon) \leq K_1^M$$

and, by Chebyshev's inequality,

$$P(H(\eta, \varepsilon) > L) \leq K_1^M e^{-L/2}.$$

Thus, if  $L = 4N + 2M(1 + \log K_1)$ , we have

$$P(H(\eta, \varepsilon) \geq L) \leq e^{-2N-M},$$

so that

$$P(\exists \varepsilon \in \{-1, 1\}^N; H(\eta, \varepsilon) \geq L) \leq e^{-(N+M)}.$$

Thus  $b = P(\eta \notin B) \leq e^{-(N+M)}$ . Since  $b \leq \frac{1}{4}$ , we get from (6.22) that

$$P(|h_N - m_N| \geq t) \leq 4e^{-(M+N)} + 8 \exp\left(-\frac{t^2}{K(N+M)}\right),$$

where  $K$  is a universal constant. The result follows.  $\square$

**7. Approximation by very many points.** Let us go back to the discussion at the beginning of Section 6. Given  $x \in \{0, 1\}^N$ , what we have used is that the functions  $h(x, y)$  ( $y \in A$ ) have a small average  $h(x) = (1 - N/2m)x$ , where  $m = \sum_{i \leq N} x_i$ . The existence of this average, however, does not fully reflect the multitude of points of  $A$  that approximate  $x$ . Indeed, to obtain such an average it would essentially suffice to have about  $m/(m - N/2)$  elements  $y \leq x$  in  $A$  such that the sets  $\{i; x_i \neq y_i\}$  are disjoint.

The result we will present in this section is a considerable strengthening of Theorem 6.1; it, however, requires a further leap into abstraction.

The basic idea is identical to that of Theorem 6.1. Given  $x, y$  in  $\Omega^N$ , we associate an object  $\nu_{x,y}$  and we express that  $x$  is close to  $A$  if there is a convex combination of the objects  $\nu_{x,y}$ ,  $y \in A$ , that is "small." In Section 6, the object  $\nu_{x,y}$  was the indicator of the set  $\{i; x_i \neq y_i\}$ . In the present section, we use a higher dimensional object so that it is harder to make averages of objects  $\nu_{x,y}$  "small" and so that, in turn, the existence of such small averages yields stronger information.

Consider a number  $0 < \theta < 1$  and the probability measure

$$\nu = ((1 - \theta)\delta_0 + \theta\delta_1)^{\otimes N}$$

on  $\{0, 1\}^N$ . Thus  $\nu$  is the law of an independent sequence  $(\eta_i)_{i \leq N}$  with  $E\eta_i = \theta$ ,  $\eta_i \in \{0, 1\}$ . Given  $x, y$  in  $\Omega^N$ , we consider the measure  $\nu_{x,y}$  on  $\{0, 1\}^N$  such



that  $\nu_{x,y}$  is the image of  $\nu$  by the map  $T$  of  $\{0, 1\}^n$  that “flips the coordinates”  $i$  for which  $x_i \neq y_i$ , that is,

$$T(u)_i = \begin{cases} u_i, & \text{if } x_i = y_i, \\ 1 - u_i, & \text{if } x_i \neq y_i. \end{cases}$$

In other words,  $\nu_{x,y}$  is the law of an independent sequence  $\eta_i \in \{0, 1\}$  such that  $E\eta_i = \theta$ , if  $x_i = y_i$ , and  $E\eta_i = 1 - \theta$ , if  $x_i \neq y_i$ . Thus, if  $\theta \neq \frac{1}{2}$ , the more coordinates of  $x$  and  $y$  are different, the more different  $\nu_{x,y}$  is from  $\nu$ .

To measure how far a probability  $\mu$  on  $\{0, 1\}^n$  is from  $\nu$ , we will use the quantity

$$\int \left(\frac{d\mu}{d\nu}\right)^2 d\nu,$$

where the integral is over  $\{0, 1\}^N$ . We observe that since  $\int (d\mu/d\nu) d\nu = 1$ , by Cauchy–Schwarz, we have  $\int (d\mu/d\nu)^2 d\nu \geq 1$ .

For a subset  $A$  of  $\Omega^N$ , we set

$$m(A, x) = \inf \left\{ \int \left(\frac{d\mu}{d\nu}\right)^2 d\nu; \mu \in \text{conv}\{\nu_{x,y}; y \in A\} \right\}.$$

**THEOREM 7.1 [29].** *Assume that  $\beta = |\theta - \frac{1}{2}| < \frac{1}{6}$  and define*

$$\alpha = \frac{32\beta^2}{1 - 36\beta^2}.$$

*Then for all subsets  $A$  of  $\Omega^N$  we have*

$$(7.1) \quad \int_{\Omega^N} m(A, x) dP(x) \leq \frac{1}{P(A)^\alpha}.$$

Certainly the condition  $|\theta - \frac{1}{2}| < \frac{1}{6}$  looks strange. The previous result is, however, a good illustration of the wonders of the induction method. Analysis of the canonical example presented at the beginning of Section 6 allows one to show that the left-hand side of (7.1) can stay bounded when  $P(A) \geq \frac{1}{2}$  independently of  $\Omega, A, N$  only when  $|\theta - \frac{1}{2}| \leq \frac{1}{6}$ . We have no intuitive explanation to offer as the reason for this “phase transition.”

As will be demonstrated later, Theorem 7.1 is in many respects a considerable strengthening of Theorem 6.1. However, it would have been hard to discover Theorem 7.1 this way, and the motivation came from the convolution problem of [24], that we recall now. Consider, on the group  $G_N = \{-1, 1\}^N$ , the Haar measure  $\lambda$  and the measure

$$\nu = \{(1 - \theta)\delta_{-1} + \theta\delta_1\}^{\otimes N}.$$

Consider the convolution operator  $T: f \rightarrow f * \nu$  from  $L^1(\lambda)$  to  $L^1(\lambda)$ . The conjecture means that  $T$  displays some regularization properties, as follows.

CONJECTURE 7.2. Consider  $f \in L^1(G)$ ,  $f \geq 0$  and  $\int f d\lambda = 1$ . Then, for all  $t \geq 0$ , we have

$$(7.2) \quad \lambda(\{Tf \geq t\}) \leq \frac{K}{t\sqrt{\log(e+t)}},$$

where  $K$  is a universal constant.

The logarithmic factor in (7.2) is apparently related to the logarithmic factor in (3.6).

The idea of Theorem 7.1 was simply that  $Tf(x) = \nu_x(f)$ , where the probability  $\nu_x$  is the translation of  $\nu$  by  $x$ . Thus, if  $A = \{x; \nu_x(f) \geq t\}$ , it should help to know that for many  $y$  we have  $\nu_y$  close to the set  $\{\nu_x, x \in A\}$  (a fact whose formulation led to Theorem 7.1). We have, however, been unable to carry out the idea.

The progress that Theorem 7.1 represents over Theorem 6.1 is exemplified by the following result, where we set

$$I_k = \{\mathbf{i} = (i_1, \dots, i_k); 1 \leq i_1 < i_2 < \dots < i_k \leq N\}.$$

PROPOSITION 7.3. Let us fix  $\theta$  with  $|\theta - \frac{1}{2}| < \frac{1}{6}$ . Assume that  $m(A, x) \leq e^t$ . Then for each  $k \geq 1$  and each family  $(\alpha_i)_{i \in I_k}$  there exists  $y$  in  $A$  such that if

$$J_k = J_k(x, y) = \{\mathbf{i} \in I_k; \forall \ell \leq k, x_{i_\ell} \neq y_{i_\ell}\},$$

then

$$\sum_{\mathbf{i} \in J_k} \alpha_{\mathbf{i}} \leq C^k t^{k/2} \left( \sum_{\mathbf{i} \in I_k} \alpha_{\mathbf{i}}^2 \right)^{1/2},$$

where  $C$  depends on  $\theta$  only.

To understand this result better, let us specialize to the case  $k = 1$ . Thus, given numbers  $(\alpha_i)_{i \leq N}$ , we can find  $y \in A$  such that

$$(7.3) \quad \sum_{i \leq N} \{\alpha_i; x_i \neq y_i\} \leq C \sqrt{\log m(A, x)} \left( \sum_{i \leq N} \alpha_i^2 \right)^{1/2}.$$

To compare (7.3) with (6.8), we have to compare the set where  $C\sqrt{\log m(A, \cdot)}$  is large with the set where  $f(A, \cdot)$  is large. We note that from (7.1) we have

$$(7.4) \quad P(C\sqrt{\log m(A, \cdot)} \geq u) \leq \frac{1}{P(A)^\alpha} \exp\left(-\frac{u^2}{C^2}\right),$$

whereas from (6.4) we have

$$(7.5) \quad P(f(A, \cdot) \geq u) \leq \frac{1}{P(A)} \exp\left(-\frac{u^2}{4}\right).$$

Thus (6.8) and (7.3) are comparable [with the exception of worse constants in (7.4) versus (7.5)], but the conclusion of Proposition 7.3 holds for any  $k \geq 1$ .

**8. Control by several points.** Let us start by providing motivation. Suppose that we are given a sequence  $(Y_i)_{i \leq N}$  of nonnegative r.v.'s and that we know that

$$P\left(\sum_{i \leq N} Y_i \leq M\right) \geq \frac{1}{2}.$$

We attempt to find bounds for the tail probabilities  $P(\sum_{i \leq N} Y_i \geq t)$ . The basic observation is as follows. Set  $A = \{\sum_{i \leq N} Y_i \leq M\}$ . Consider  $\omega'$  and  $\omega$  in  $A$ . Set

$$I = \{i \leq N; Y_i(\omega) = Y_i(\omega')\}$$

so that

$$(8.1) \quad \sum_{i \in I} Y_i(\omega') = \sum_{i \in I} Y_i(\omega) \leq M$$

by *positivity* of  $Y_i$ . Consider now  $\omega^1, \dots, \omega^q$  in  $A$  and set

$$(8.2) \quad J = \{i \leq N; \exists \ell \leq q, Y_i(\omega') = Y_i(\omega^\ell)\}.$$

Thus, by (8.1) and the positivity of  $Y_i$ , we have

$$(8.3) \quad \sum_{i \leq N} Y_i(\omega') \leq qM + \sum_{i \notin J} Y_i.$$

One then hopes that if  $\text{card}\{i \notin J\}$  is small, we will be able to control the last term. This discussion should provide motivation for the following discussion. Consider an integer  $q \geq 2$ . For  $x \in \Omega^N$  and subsets  $A_1, \dots, A_q$  of  $\Omega^N$  we define

$$(8.4) \quad \begin{aligned} f(A_1, \dots, A_q, x) \\ = \inf\{\text{card}\{i \leq N; x_i \in \{y_i^1, \dots, y_i^q\}\}: y^1 \in A_1, \dots, y^q \in A_q\}. \end{aligned}$$

What we are really interested in is the case  $A_1 = A_2 = \dots = A_q$ , but the proof by induction requires considering different sets. Later we will prove the following basic fact about  $f(A_1, \dots, A_q, x)$ .

**THEOREM 8.1.** *If  $P$  is a product measure, we have*

$$(8.5) \quad \int q^{f(A_1, \dots, A_q, x)} dP(x) \leq \frac{1}{\prod_{i \leq q} P(A_i)}$$

and, in particular,

$$(8.6) \quad P(f(A, \dots, A, \cdot) \geq k) \leq \frac{1}{q^k P(A)^q}.$$

Combining with (8.3) we see that if  $S_k$  denotes the sum of the largest  $k$  terms of the sequence  $(Y_i)_{i \leq N}$ , we have

$$(8.7) \quad P\left(\sum_{i \leq N} Y_i \geq qM + t\right) \leq \frac{2^q}{q^k} + P(S_k \geq t).$$

Hopefully the last term can be controlled by classical methods, and it remains only to optimize (8.7) over the various parameters.

Certainly the previous method seems an overkill to study the tails of  $\sum_{i \leq N} Y_i$ . Suppose, however, that we now have a function  $f$  on  $\Omega^N$  and functions  $(Y_i)_{i \leq N}$  such that if  $A = \{f \leq M\}$ , where  $M$  is the median of  $f$ , the following holds. Given  $x \in \Omega^N$ ,  $y^1, \dots, y^q \in A$  and

$$J = \{i \leq N; \exists \ell \leq q, x_i = y_i^\ell\},$$

then

$$(8.8) \quad f(x) \leq qM + S_k,$$

where  $k = \text{card}\{i \leq N; i \notin J\}$  and  $S_k$  is the sum of the  $k$  largest terms of the sequence  $(Y_i(x_i))_{i \leq N}$ . Then

$$(8.9) \quad P(f \geq qM + t) \leq \frac{2^q}{q^k} + P(S_k \geq t).$$

To give the most important example of this situation, let us consider the case where

$$f(x) = E_\varepsilon \left\| \sum_{i \leq N} \varepsilon_i Z_i(x_i) \right\|,$$

for functions  $Z_i$  from  $\Omega$  to a Banach space, and where  $(\varepsilon_i)_{i \leq N}$  are independent Bernoulli r.v.'s. The key observation is that the function

$$E_\varepsilon \left\| \sum_{i \in I} \varepsilon_i Z_i(x_i) \right\|$$

is an increasing function of  $I$ , as is seen by taking the expectation with respect to certain  $\varepsilon_i$ 's inside rather than outside the norm. Thus, when  $J = \bigcup_{\ell \leq q} I_\ell$ , by the triangle inequality we have

$$E_\varepsilon \left\| \sum_{i \in J} \varepsilon_i Z_i(x_i) \right\| \leq \sum_{\ell \leq q} E_\varepsilon \left\| \sum_{i \in J} \varepsilon_i Z_i(x_i) \right\|$$

and thus

$$f(x) \leq \sum_{\ell \leq q} E_\varepsilon \left\| \sum_{i \in J} \varepsilon_i Z_i(x_i) \right\| + \sum_{i \notin J} \|Z_i(x_i)\|.$$

This implies that (8.8) holds for  $Y_i = \|Z_i\|$ . An important contribution of Ledoux [7] made clear that controlling  $f$  is the main step in controlling  $\|\sum_{i \leq n} \varepsilon_i Z_i(x_i)\|$ . This approach using inequality (8.9) has been very successful, as demonstrated in [9], and it is remarkable that its proof is now so simple.

The key fact of the proof of Theorem 8.1 is the following simple statement about functions.

LEMMA 8.2. Consider a function  $g$  on  $\Omega$ , such that  $1/q \leq g \leq 1$ . Then

$$(8.10) \quad \int_{\Omega} \frac{1}{g} d\mu \left( \int_{\Omega} g d\mu \right)^q \leq 1.$$

PROOF. Observing that  $\log x \leq x - 1$ , to prove that  $ab^q \leq 1$  it suffices to show that  $a + qb \leq q + 1$ . Thus, it suffices to show that

$$\int_{\Omega} \frac{1}{g} d\mu + q \int_{\Omega} g d\mu \leq q + 1,$$

but this is obvious since  $x^{-1} + qx \leq q + 1$  for  $q^{-1} \leq x \leq 1$ .  $\square$

COROLLARY 8.3. Consider functions  $g_i$  on  $\Omega$ ,  $0 \leq g_i \leq 1$ . Then

$$(8.11) \quad \int_{\Omega} \min_{i \leq q} \left( q, \frac{1}{g_i} \right) d\mu \prod_{i \leq q} \int_{\Omega} g_i d\mu \leq 1.$$

PROOF. Set  $g = (\min_{i \leq q} (q, g_i^{-1}))^{-1}$ , observe that  $g_i \leq g$  and use (8.10).  $\square$

We now prove Theorem 8.1 by induction over  $N$ . For  $N = 1$ , the result follows from (8.11) taking  $g_i = 1_{A_i}$ .

We assume now that Theorem 8.1 has been proved for  $N$  and we prove it for  $N + 1$ . Consider sets  $A_1, \dots, A_q$  of  $\Omega^{N+1}$ . For  $\omega \in \Omega$ , we define the sets  $A_i(\omega)$  as in (5.4) and we consider the projection  $B_i$  of  $A_i$  on  $\Omega^N$ . The basic observation is that

$$(8.12) \quad f(A_1, \dots, A_q, (x, \omega)) \leq 1 + f(B_1, \dots, B_q, x)$$

and that, whenever  $j \leq q$ ,

$$(8.13) \quad f(A_1, \dots, A_q, (x, \omega)) \leq f(C_1, \dots, C_q, x),$$

where  $C_i = B_i$  for  $i \neq j$ ,  $C_j = A_j(\omega)$ .

To prove Theorem 8.1, we observe that, using (8.12) and induction hypothesis, we have

$$\int q^{f(A_1, \dots, A_q, (x, \omega))} dP(x) \leq q \frac{1}{\prod_{i \leq q} P(B_i)}$$

while using (8.13) we get

$$\int q^{f(A_1, \dots, A_q, (x, \omega))} dP(x) \leq q \frac{1}{\prod_{i \leq q} P(C_i)}.$$

Thus, setting  $g_i(\omega) = P(A_i(\omega))/P(B_i)$ , we have

$$\int q^{f(A_1, \dots, A_q, (x, \omega))} dP(x) \leq \frac{1}{\prod_{i \leq q} P(B_i)} \int \min \left( q, \min_{i \leq q} \frac{1}{g_i(\omega)} \right) d\mu(\omega).$$

Using now the Fubini theorem and (8.11), we have

$$\int q^{f(A_1, \dots, A_q(x, \omega))} dP(x) d\mu(\omega) \leq \frac{1}{\prod_{i \leq q} P(B_i) \int g_i d\mu},$$

which finishes the proof since  $\int g_i d\mu = P(A_i)/P(B_i)$ .

**9. Penalties.** Roughly speaking, the Hamming distance measures how far  $x$  is from  $A$  by counting the smallest number of coordinates of  $x$  that cannot be captured by a point of  $A$ . Thus we get one penalty for each coordinate we miss. A natural extension of this idea is to consider a nonnegative function  $h$  on  $\Omega \times \Omega$  and, for  $x \in \Omega^N$ ,  $A \subset \Omega^N$ , to consider

$$(9.1) \quad f_h(A, x) = \inf \left\{ \sum_{i \leq N} h(x_i, y_i); y \in A \right\}$$

as a way to measure the “distance” from  $x$  to  $A$ .

It is reasonable to require

$$(9.2) \quad \forall \omega \in \Omega, \quad h(\omega, \omega) = 0.$$

Thus, the case of the Hamming distance is simply  $h(x, y) = 1_{\{x \neq y\}}$ .

We observe that, since  $x, y$  do not play the same role, we will not require  $h$  to be symmetric. In contrast with the work of Sections 5–8 that requires *no* structure on  $\Omega$ , Definition 9.1 does require a minimum of structure, namely, the existence of the function  $h$ . On the other hand, this opens the door to a theory whose complexity certainly would not have been suspected beforehand.

Certainly one needs some control on the size of  $h$ . The most obvious way to achieve this is through moment conditions on  $h$ . A typical result is as follows.

**THEOREM 9.1.** *Set*

$$(9.3) \quad \begin{aligned} \|h\|_\infty &= \sup\{h(x, y); x, y \in \Omega\}, \\ \|h\|_2^2 &= \iint_{\Omega^2} h^2(\omega, \omega') d\mu(\omega) d\mu(\omega'). \end{aligned}$$

*Then, for each subset  $A$  of  $\Omega^N$ , we have*

$$(9.4) \quad P(f_h(A, \cdot) \geq u) \leq \frac{1}{P(A)} \exp\left(-\min\left(\frac{u^2}{8N\|h\|_2^2}, \frac{u}{2\|h\|_\infty}\right)\right).$$

We do not know how to obtain sharp numerical constants in (9.4). Inequality (9.4) generalizes Bernstein’s inequality the way Theorem 9.1 generalizes (1.3). If  $g$  is a function on  $\Omega$ , setting  $h(x, y) = |g(x) - g(y)|$ , it is an interesting exercise to recover from (9.4) a qualitatively correct version of Bernstein’s inequality (i.e., only the numerical constants are different).

It is arguable that Theorem 9.1 does not represent a truly new phenomenon. It turns out, however, that in Theorem 9.1 what matters is not really  $h$ , but rather the following functional, defined for all subsets  $B$  of  $\Omega$ :

$$(9.5) \quad h(\omega, B) = \inf\{h(\omega, \omega'); \omega' \in B\}.$$

**THEOREM 9.2.** *Assume that for each subset  $B$  of  $\Omega$  we have*

$$(9.6) \quad \int_{\Omega} \exp 2h(x, B) d\mu(x) \leq \frac{e}{\mu(B)}.$$

*Then for  $t \leq 1$  and each subset  $A$  of  $\Omega^N$  we have*

$$(9.7) \quad \int_{\Omega^N} \exp tf_h(A, x) dP(x) \leq \frac{\exp(t^2 N)}{P(A)}.$$

*In particular, if  $u \leq 2N$  we have*

$$(9.8) \quad P(f_h(A, \cdot) \geq u) \leq \frac{1}{P(A)} \exp\left(-\frac{u^2}{4N}\right).$$

The point of (9.6) is that taking the infimum in (9.5) has a dramatic effect and that condition (9.6) is less stringent than the control of

$$\int_{\Omega^2} \exp 2h(x, y) d\mu(x) d\mu(y)$$

one would expect would be required in order to obtain something like (9.8).

We illustrate this in the case where  $\Omega$  is itself a product of  $m$  spaces and where  $h(x, y) = ad(x, y)$ , where  $d$  is the Hamming distance on  $\Omega$  and  $a$  is a parameter. It follows from (5.1) that (9.6) holds for  $a = 2m^{-1/2}$ . On the other hand, if  $\|h\|_2$  is given by (9.3), then for this value of  $a$ ,  $\|h\|_2$  is of order  $\sqrt{m}$ , so that there is a loss of a factor  $\sqrt{m}$  in the exponent in (9.4) compared to (9.8).

To give a vivid illustration of what Theorem 9.2 can prove, consider a product space  $\Omega^N$ . Consider a subset  $A$  of  $\Omega^N$ ,  $P(A) \geq \frac{1}{2}$ . Then for most elements  $x$  of  $\Omega^N$ , we can find an element  $y$  of  $A$  such that the set  $I = \{i \leq N; x_i \neq y_i\}$  has a cardinal of order  $\sqrt{N}$ . This is the content of Proposition 5.1; but now if we view  $N$  as built from  $N_1$  blocks of length  $N_2$  ( $N = N_1 N_2$ ), we can moreover require that  $I$  meets only about  $\sqrt{N_1}$  blocks.

One of the most interesting phenomena related to the theory of penalties occurs under a condition somewhat stronger than (9.6). However, rather than stating the most general theorem (it requires some effort to understand the hypothesis), we will only state the most important case. In that case,  $\Omega = \mathbb{R}$ ,  $\mu$  has a density  $\frac{1}{2}e^{-|x|}$  with respect to Lebesgue measure, and the function  $h$  is given by

$$h(x, y) = \min(|x - y|, (x - y)^2).$$

**THEOREM 9.3.** *For some universal constant  $K$  and each subset  $A$  of  $\mathbb{R}^N$ , we have*

$$(9.9) \quad \int_{\Omega^N} \exp\left(\frac{1}{K} f_h(A, x)\right) dP(x) \leq \frac{1}{P(A)}.$$

The most obviously remarkable feature of this theorem is that (9.9) does not depend on  $N$ . The depth of Theorem 9.3 can, however, better be measured

by the fact that it does constitute a kind of improvement upon what was previously known about Gaussian measure. To see this, consider the nondecreasing map  $\varphi$  from  $\mathbb{R}$  to  $\mathbb{R}$  that transforms  $\mu$  into the one-dimensional Gaussian measure  $\gamma_1$ . It is a simple fact to see that, for some universal constant  $K$ , we have

$$(9.10) \quad (\varphi(x) - \varphi(y))^2 \leq Kh(x, y).$$

On the other hand, the map  $\psi$  from  $\mathbb{R}^N$  to  $\mathbb{R}^N$  given by  $\psi((x_i)_{i \leq N}) = (\varphi(x_i))_{i \leq N}$  transforms  $P$  into  $\gamma_N$ .

Consider now  $B \subset \mathbb{R}^N$ . Thus

$$\gamma_N(B) = P(\psi^{-1}(B))$$

Now, by (9.10), we have

$$d(\psi(A), \psi(x))^2 \leq Kf_h(A, x),$$

where  $d(B, y)$  is the Euclidean distance from  $B$  to  $y$  and thus from (9.9),

$$\int_{\Omega^N} \exp\left(\frac{1}{K} d(\psi(\psi^{-1}(B)), \psi(x))^2\right) dP(x) \leq \frac{1}{\gamma_N(B)},$$

so that (for a new constant  $K$ )

$$\int_{\mathbb{R}^N} \exp\left(\frac{1}{K} d(B, y)^2\right) d\gamma_N(y) \leq \frac{1}{\gamma_N(B)}.$$

Therefore, for  $t > 0$ ,

$$(9.11) \quad \gamma_N(d(B, \cdot) \geq t) \leq \frac{1}{\gamma_N(B)} \exp\left(-\frac{t^2}{K}\right).$$

In the case  $\gamma_N(B) = \frac{1}{2}$ , this is a weak form of (2.9).

It turns out that, for many applications, (9.11) rather than (2.9) suffices. In particular, it is now clearly understood that (9.11) is one of the central facts that allows us to characterize continuity and boundedness of Gaussian processes [23]. The importance of Theorem 9.3 is that it allows extending these characterizations to more general processes [26].

One of the most intriguing further aspects of the theory of penalties is that the roles of  $x$  and  $y$  in the penalty function  $h(x, y)$  are highly asymmetric. This is particularly apparent when the idea of penalty function is combined with the method of Section 6, a topic for which we must refer to [29].

In conclusion, we have tried to make the reader aware that there are unexpectedly subtle phenomena related to concentration of measure in product spaces. That such a rich theory should exist at all with such minimal structure is certainly remarkable, as is remarkable the width of its applications. It is not clear to me at present where the potential for future advances, if any, lies. A worthy project would be a systematic development of the "transportation method" that very recently arose from the work of Marton [12]. This method is a potentially serious competitor to the induction method presented here. It allows, in some cases, an easier computation of the best constants and an easier approach to Theorem 9.1 [30], but whether it can lead to genuinely new



results in the independent case is unclear at present. In a different direction, an obvious research question is whether there exists at all a usable theory beyond the case of product measures; see, for example, [29] for the case of the symmetric group (that resembles a product) and of [12] for certain Markov chains.

**Acknowledgment.** The author is indebted to Michel Ledoux and Gilles Godefroy for many useful comments.

*Note added in proof.* After this paper was written, further progress was made, in particular on the material of Section 7 [31].

## REFERENCES

- [1] AMIR, D. and MILMAN, V. D. (1980). Unconditional and symmetric sets in  $n$ -dimensional normed spaces. *Israel J. Math.* **37** 3–20.
- [2] BOLLABÁS, B. and BRIGHTWELL, G. (1992). The height of a random partial order: concentration of measure. *Ann. Appl. Probab.* **2** 1009–1018.
- [3] EHRHARD, A. (1983). Symétrisation dans l'espace de Gauss. *Math. Scand.* **53** 281–301.
- [4] GROMOV, M. and MILMAN, V. D. (1983). A topological application of the isoperimetric inequality. *Amer. J. Math.* **105** 843–854.
- [5] Hoeffding, W. (1963). Probability inequalities for sums of bounded random variables. *J. Amer. Statist. Assoc.* **58** 13–30.
- [6] HOPFIELD, J. J. (1982). Neural networks and physical systems with emergent collective computational abilities. *Proc. Nat. Acad. Sci. USA* **79** 2554–2558.
- [7] LEDOUX, M. (1985). Gaussian randomization and the law of the iterated logarithm in Banach spaces. Unpublished manuscript.
- [8] LEDOUX, M. (1993). Les inégalités isopérimétriques en analyse et probabilités. Séminaire Bourbaki, juin 1993. *Astérisque* **216** 343–375.
- [9] LEDOUX, M. and TALAGRAND, M. (1991). *Probability in Banach Spaces*. Springer, Berlin.
- [10] MARGULIS, G. A. (1977). Probabilistic characteristics of graphs with large connectivity. *Problems Inform. Transmission* **10** 174–179.
- [11] MARTON, K. (1986). A simple proof of the blowing-up lemma. *IEEE Trans. Inform. Theory* **IT-32** 445–446.
- [12] MARTON, K. (1994). A concentration of measure inequality for contracting Markov chains. Unpublished manuscript.
- [13] MAUREY, B. (1979). Construction de suites symétriques. *C. R. Acad. Sci. Paris* **288** 679–681.
- [14] MCDIARMID, C. (1989). On the method of bounded differences. In *Surveys in Combinatorics, 1989* (J. Simons, ed.) 148–188. Cambridge Univ. Press.
- [15] MILMAN, V. D. (1971). A new proof of the theorem of A. Dvoretzky on sections of convex bodies. *Funct. Anal. Appl.* **5** 28–37.
- [16] MILMAN, V. D. (1988). The heritage of P. Levy in geometrical functional analysis. *Astérisque* **157/158** 273–301.
- [17] MILMAN, V. and SCHECHTMAN, G. (1986). *Asymptotic Theory of Finite Dimensional Normed Spaces. Lecture Notes in Math.* **1200**. Springer, Berlin.
- [18] ODLYZKO, A. (1993). Private communication.
- [19] PASTUR, L. and SHCHERBINA, M. (1991). Absence of self-averaging of the order parameter in the Sherrington–Kirkpatrick model. *J. Statist. Phys.* **62** 1–19.
- [20] RHEE, W. and TALAGRAND, M. (1987). Martingales inequality and NP-complete problems. *Math. Oper. Res.* **12** 177–181.
- [21] SCHECHTMAN, G. (1981). Levy type inequality for a class of metric spaces. In *Martingale Theory in Harmonic Analysis and Banach Spaces Lecture Notes in Math.* **939** 211–215. Springer, Berlin.

- [22] SHAMIR, E. and SPENCER, J. (1987). Sharp concentration of the chromatic number of random graphs  $G_{n,p}$ . *Combinatorica* **7** 121–129.
- [23] TALAGRAND, M. (1987). Regularity of Gaussian processes. *Acta Math.* **159** 99–149.
- [24] TALAGRAND, M. (1989). A conjecture on convolution operators, and operators from  $L^1$  to a Banach lattice. *Israel J. Math.* **68** 82–88.
- [25] TALAGRAND, M. (1993). Isoperimetry, logarithmic Sobolev inequalities on the discrete cube and Margulis' graph connectivity theorem. *Geom. Funct. Anal.* **3** 295–314.
- [26] TALAGRAND, M. (1994). Supremum of some canonical processes. *Amer. J. Math.* **116** 283–325.
- [27] TALAGRAND, M. (1994). Sharper bounds for Gaussian and empirical processes. *Ann. Probab.* **22** 28–76.
- [28] TALAGRAND, M. (1994). Isoperimetry in product spaces: higher level, large sets. Unpublished manuscript.
- [29] TALAGRAND, M. (1995). Concentration of measure and isoperimetric inequalities in product spaces. *Publications Math. IHES* **81** 73–205.
- [30] TALAGRAND, M. (1995). Transportation cost for Gaussian and other product measures. *Geom. Funct. Anal.* To appear.
- [31] TALAGRAND, M. (1995). New concentration inequalities in product spaces. Unpublished manuscript.
- [32] THURSTON, W. (1994). On proof and progress in mathematics. *Bull. Amer. Math. Soc.* **30** 161–177.
- [33] YURINSKII, V. V. (1974). Exponential bounds for large deviations. *Theory Probab. Appl.* **19** 154–155.

EQUIPE D'ANALYSE—TOUR 56  
E.R.A. AU C.N.R.S. NO. 754  
UNIVERSITÉ PARIS VI  
4 PLACE JUSSIEU  
75230 PARIS CEDEX 05  
FRANCE

DEPARTMENT OF MATHEMATICS  
THE OHIO STATE UNIVERSITY  
231 W. 18TH AVENUE  
COLUMBUS, OHIO 43210-1174  
E-mail: talagran@math.ohio-state.edu