

in volume 2. All results will be derived anew, independently, by different methods. This chapter will therefore serve primarily readers who are not in a hurry to proceed with the systematic theory, or readers interested in the spirit of probability theory without wanting to specialize in it. For other readers a comparison of methods should prove instructive and interesting. Accordingly, *the present chapter should be read at the reader's discretion independently of, or parallel to, the remainder of the book.*

1. GENERAL ORIENTATION. THE REFLECTION PRINCIPLE

From a formal point of view we shall be concerned with arrangements of finitely many plus ones and minus ones. Consider $n = p + q$ symbols $\epsilon_1, \dots, \epsilon_n$, each standing either for $+1$ or for -1 ; suppose that there are p plus ones and q minus ones. The partial sum $s_k = \epsilon_1 + \dots + \epsilon_k$ represents the difference between the number of pluses and minuses occurring at the first k places. Then

$$(1.1) \quad s_k - s_{k-1} = \epsilon_k = \pm 1, \quad s_0 = 0, \quad s_n = p - q,$$

where $k = 1, 2, \dots, n$.

We shall use a geometric terminology and refer to rectangular coordinates t, x ; for definiteness we imagine the t -axis is horizontal, the x -axis vertical. The arrangement $(\epsilon_1, \dots, \epsilon_n)$ will be represented by a polygonal line whose k th side has slope ϵ_k and whose k th vertex has ordinate s_k . Such lines will be called paths.

Definition. Let $n > 0$ and x be integers. A path (s_1, s_2, \dots, s_n) from the origin to the point (n, x) is a polygonal line whose vertices have abscissas $0, 1, \dots, n$ and ordinates s_0, s_1, \dots, s_n satisfying (1.1) with $s_n = x$.

We shall refer to n as the *length* of the path. There are 2^n paths of length n . If p among the ϵ_k are positive and q are negative, then

$$(1.2) \quad n = p + q, \quad x = p - q.$$

A path from the origin to an arbitrary point (n, x) exists only if n and x are of the form (1.2). In this case the p places for the positive ϵ_k can be chosen from the $n = p + q$ available places in

$$(1.3) \quad N_{n,x} = \binom{p+q}{p} = \binom{p+q}{q}$$

different ways. For convenience we define $N_{n,x} = 0$ whenever n and x

are not of the form (1.2). With this convention there exist exactly $N_{n,x}$ different paths from the origin to an arbitrary point (n, x) .

Before turning to the principal topic of this chapter, namely the theory of random walks, we illustrate possible applications of our scheme.

Examples. (a) *The ballot theorem.* The following amusing proposition was proved in 1878 by W. A. Whitworth, and again in 1887 by J. Bertrand.

Suppose that, in a ballot, candidate P scores p votes and candidate Q scores q votes, where $p > q$. The probability that throughout the counting there are always more votes for P than for Q equals $(p-q)/(p+q)$.

Similar problems of arrangements have attracted the attention of students of combinatorial analysis under the name of ballot problems. The recent renaissance of combinatorial methods has increased their popularity, and it is now realized that a great many important problems may be reformulated as variants of some generalized ballot problem.³

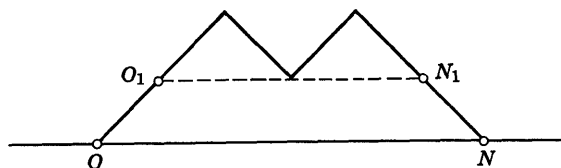


Figure 1. Illustrating positive paths. The figure shows also that there are exactly as many strictly positive paths from the origin to the point $(2n, 0)$ as there are non-negative paths from the origin to $(2n-2, 0)$.

The whole voting record may be represented by a path of length $p + q$ in which $\epsilon_k = +1$ if the k th vote is for P ; conversely, every path from the origin to the point $(p + q, p - q)$ can be interpreted as a record of a voting with the given totals p and q . Clearly s_k is the number of votes by which P leads, or trails, just after the k th vote is cast. The candidate P leads throughout the voting if, and only if, $s_1 > 0, \dots, s_n > 0$, that is, if all vertices lie strictly above the t -axis. (The path from O to N_1 in figure 1 is of this type.) The ballot theorem assumes tacitly that all admissible paths are equally probable. The assertion then reduces to the theorem proved at the end of this section as an immediate consequence of the reflection lemma.

(b) *Galton's rank order test.*⁴ Suppose that a quantity (such as the height

³ A survey of the history and the literature may be found in *Some aspects of the random sequence*, by D. E. Barton and C. L. Mallows [Ann. Math. Statist., vol. 36 (1965), pp. 236-260]. These authors discuss also various applications. The most recent generalization with many applications in queuing theory is due to L. Takacs.

⁴ J. L. Hodges, *Biometrika*, vol. 42 (1955), pp. 261-262.

of plants) is measured on each of r treated subjects, and also on each of r control subjects. Denote the measurements by a_1, \dots, a_r and b_1, \dots, b_r , respectively. To fix ideas, suppose that each group is arranged in decreasing order: $a_1 > a_2 > \dots$ and $b_1 > b_2 > \dots$. (To avoid trivialities we assume that no two observations are equal.) Let us now combine the two sequences into one sequence of $n = 2r$ numbers arranged in decreasing order. For an extremely successful treatment all the a 's should precede the b 's, whereas a completely ineffectual treatment should result in a random placement of a 's and b 's. Thus the efficiency of the treatment can be judged by the number of different a 's that precede the b of the same rank, that is, by the number of subscripts k for which $a_k > b_k$. This idea was first used in 1876 by F. Galton for data referred to him by Charles Darwin. In this case r equaled 15 and the a 's were ahead 13 times. Without knowledge of the actual probabilities Galton concluded that the treatment *was* effective. But, assuming perfect randomness, the probability that the a 's lead 13 times or more equals $\frac{3}{16}$. This means that in three out of sixteen cases a perfectly ineffectual treatment would appear as good or better than the treatment classified as effective by Galton. This shows that a quantitative analysis may be a valuable supplement to our rather shaky intuition.

For an interpretation in terms of paths write $\epsilon_k = +1$ or -1 according as the k th term of the combined sequence is an a or a b . The resulting path of length $2r$ joins the origin to the point $(2r, 0)$ of the t -axis. The event $a_k > b_k$ occurs if, and only if, s_{2k-1} contains at least k plus ones, that is, if $s_{2k-1} > 0$. This entails $s_{2k} \geq 0$, and so the $(2k-1)$ st and the $2k$ th sides are above the t -axis. It follows that the inequality $a_k > b_k$ holds ν times if, and only if, 2ν sides lie above the t -axis. In section 9 we shall prove the unexpected result that the probability for this is $1/(r+1)$, irrespective of ν . (For related tests based on the theory of runs see II, 5.b.)

(c) *Tests of the Kolmogorov-Smirnov type.* Suppose that we observe two populations of the same biological species (animals or plants) living at different places, or that we wish to compare the outputs of two similar machines. For definiteness let us consider just one measurable characteristic such as height, weight, or thickness, and suppose that for each of the two populations we are given a sample of r observations, say a_1, \dots, a_r and b_1, \dots, b_r . The question is roughly whether these data are consistent with the hypothesis that the two populations are statistically identical. In this form the problem is vague, but for our purposes it is not necessary to discuss its more precise formulation in modern statistical theory. It suffices to say that the tests are based on a comparison of the two empirical distributions. For every t denote by $A(t)$ the fraction k/n of subscripts i for which $a_i \leq t$. The function so defined over the

real axis is the *empirical distribution* of the a 's. The empirical distribution B is defined in like manner. A refined mathematical theory originated by N. V. Smirnov (1939) derives the probability distribution of the maximum of the discrepancies $|A(t) - B(t)|$ and of other quantities which can be used for testing the stated hypothesis. The theory is rather intricate, but was greatly simplified and made more intuitive by B. V. Gnedenko who had the lucky idea to connect it with the geometric theory of paths. As in the preceding example we associate with the two samples a path of length $2r$ leading from the origin to the point $(2r, 0)$. To say that the two populations are statistically indistinguishable amounts to saying that ideally the sampling experiment makes all possible paths equally probable. Now it is easily seen that $|A(t) - B(t)| > \xi$ for some t if, and only if, $|s_k| > \xi r$ for some k . The probability of this event is simply the probability that a path of length $2r$ leading from the origin to the point $(0, 2r)$ is not constrained to the interval between $\pm \xi r$. This probability has been known for a long time because it is connected with the ruin problem in random walks and with the physical problem of diffusion with absorbing barriers. (See problem 3.)

This example is beyond the scope of the present volume, but it illustrates how random walks can be applied to problems of an entirely different nature.

(d) *The ideal coin-tossing game and its relation to stochastic processes.* A path of length n can be interpreted as the record of an ideal experiment consisting of n successive tosses of a coin. If $+1$ stands for heads, then s_k equals the (positive or negative) excess of the accumulated number of heads over tails at the conclusion of the k th trial. The classical description introduces the fictitious gambler Peter who at each trial wins or loses a unit amount. The sequence s_1, s_2, \dots, s_n then represents Peter's successive cumulative gains. It will be seen presently that they are subject to chance fluctuations of a totally unexpected character.

The picturesque language of gambling should not detract from the general importance of the coin-tossing model. In fact, the model may serve as a first approximation to many more complicated chance-dependent processes in physics, economics, and learning theory. Quantities such as the energy of a physical particle, the wealth of an individual, or the accumulated learning of a rat are supposed to vary in consequence of successive collisions or random impulses of some sort. For purposes of a first orientation one assumes that the individual changes are of the same magnitude, and that their sign is regulated by a coin-tossing game. Refined models take into account that the changes and their probabilities vary from trial to trial, but even the simple coin-tossing model leads to surprising, indeed to shocking, results. They are of practical importance because they

show that, contrary to generally accepted views, the laws governing a prolonged series of individual observations will show patterns and averages far removed from those derived for a whole population. In other words, currently popular psychological tests would lead one to say that in a population of "normal" coins most individual coins are "maladjusted."

It turns out that the chance fluctuations in coin tossing are typical for more general chance processes with cumulative effects. Anyhow, it stands to reason that if even the simple coin-tossing game leads to paradoxical results that contradict our intuition, the latter cannot serve as a reliable guide in more complicated situations. ◀

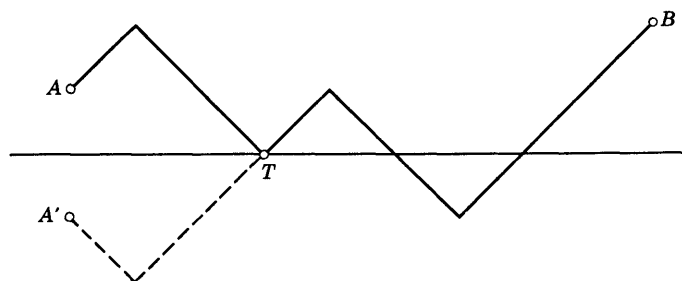


Figure 2. Illustrating the reflection principle.

It is as surprising as it is pleasing that most important conclusions can be drawn from the following simple lemma.

Let $A = (a, \alpha)$ and $B = (b, \beta)$ be integral points in the positive quadrant: $b > a \geq 0$, $\alpha > 0$, $\beta > 0$. By reflection of A on the t -axis is meant the point $A' = (a, -\alpha)$. (See figure 2.) A path from A to B is defined in the obvious manner.

Lemma.⁵ (*Reflection principle.*) *The number of paths from A to B which touch or cross the x -axis equals the number of all paths from A' to B .*

Proof. Consider a path $(s_a = \alpha, s_{a+1}, \dots, s_b = \beta)$ from A to B having one or more vertices on the t -axis. Let t be the abscissa of the first such vertex (see figure 2); that is, choose t so that $s_a > 0, \dots, s_{t-1} > 0$, $s_t = 0$. Then $(-s_a, -s_{a+1}, \dots, -s_{t-1}, s_t = 0, s_{t+1}, s_{t+2}, \dots, s_b)$ is a

⁵ The reflection principle is used frequently in various disguises, but without the geometrical interpretation it appears as an ingenious but incomprehensible trick. The probabilistic literature attributes it to D. André (1887). It appears in connection with the difference equations for random walks in XIV, 9. These are related to some partial differential equations where the reflection principle is a familiar tool called *method of images*. It is generally attributed to Maxwell and Lord Kelvin. For the use of repeated reflections see problems 2 and 3.

path leading from A' to B and having $T = (t, 0)$ as its first vertex on the t -axis. The sections AT and $A'T$ being reflections of each other, there exists a one-to-one correspondence between all paths from A' to B and such paths from A to B that have a vertex on the x -axis. This proves the lemma. ►

As an immediate consequence we prove the result discussed in example (a). It will serve as starting point for the whole theory of this chapter.

The ballot theorem. *Let n and x be positive integers. There are exactly $\frac{x}{n} N_{n,x}$ paths $(s_1, \dots, s_n = x)$ from the origin to the point (n, x) such that $s_1 > 0, \dots, s_n > 0$.*

Proof. Clearly there exist exactly as many admissible paths as there are paths from the point $(1, 1)$ to (n, x) which neither touch or cross the t -axis. By the last lemma the number of such paths equals

$$N_{n-1, x-1} - N_{n-1, x+1} = \binom{p+q-1}{p-1} - \binom{p+q-1}{p}$$

with p and q defined in (1.2). A trite calculation shows that the right side equals $N_{n,x}(p-q)/(p+q)$, as asserted. ►

2. RANDOM WALKS: BASIC NOTIONS AND NOTATIONS

The ideal coin-tossing game will now be described in the terminology of random walks which has greater intuitive appeal and is better suited for generalizations. As explained in the preceding example, when a path (s_1, \dots, s_ρ) is taken as record of ρ successive coin tossings the partial sums s_1, \dots, s_ρ represent the successive cumulative gains. For the geometric description it is convenient to pretend that the tossings are performed at a uniform rate so that the n th trial occurs at epoch⁶ n . The successive partial sums s_1, \dots, s_n will be marked as points on the vertical x -axis; they will be called the positions of a "particle" performing a random walk. Note that the particle moves in unit steps, up or down, on a

⁶ Following J. Riordan, the word *epoch* is used to denote *points* on the time axis because some contexts use the alternative terms (such as moment, time, point) in different meanings. Whenever used mathematically, the word time will refer to an interval or duration. A physical experiment may take some time, but our ideal trials are timeless and occur at epochs.