# 1 THE WAVE FUNCTION

## 1.1 THE SCHRÖDINGER EQUATION

Imagine a particle of mass $m$, constrained to move along the $x$ axis, subject to some specified force $F(x, t)$ (Figure 1.1). The program of *classical* mechanics is to determine the position of the particle at any given time: $x(t)$. Once we know that, we can figure out the velocity ($v = dx/dt$), the momentum ($p = mv$), the kinetic energy $\left(T = (1/2)mv^2\right)$, or any other dynamical variable of interest. And how do we go about determining $x(t)$? We apply Newton's second law: $F = ma$. (For *conservative* systems—the only kind we shall consider, and, fortunately, the only kind that *occur* at the microscopic level—the force can be expressed as the derivative of a potential energy function,[1] $F = -\partial V/\partial x$, and Newton's law reads $m\, d^2x/dt^2 = -\partial V/\partial x$.) This, together with appropriate initial conditions (typically the position and velocity at $t = 0$), determines $x(t)$.

Quantum mechanics approaches this same problem quite differently. In this case what we're looking for is the particle's **wave function**, $\Psi(x, t)$, and we get it by solving the **Schrödinger equation**:

$$i\hbar\frac{\partial \Psi}{\partial t} = -\frac{\hbar^2}{2m}\frac{\partial^2 \Psi}{\partial x^2} + V\Psi. \tag{1.1}$$

Here $i$ is the square root of $-1$, and $\hbar$ is Planck's constant—or rather, his *original* constant ($h$) divided by $2\pi$:

$$\hbar = \frac{h}{2\pi} = 1.054572 \times 10^{-34} \,\text{J s}. \tag{1.2}$$

The Schrödinger equation plays a role logically analogous to Newton's second law: Given suitable initial conditions (typically, $\Psi(x, 0)$), the Schrödinger equation determines $\Psi(x, t)$ for all future time, just as, in classical mechanics, Newton's law determines $x(t)$ for all future time.[2]

## 1.2 THE STATISTICAL INTERPRETATION

But what exactly *is* this "wave function," and what does it do for you once you've *got* it? After all, a particle, by its nature, is localized at a point, whereas the wave function (as its name

---

[1] Magnetic forces are an exception, but let's not worry about them just yet. By the way, we shall assume throughout this book that the motion is nonrelativistic ($v \ll c$).

[2] For a delightful first-hand account of the origins of the Schrödinger equation see the article by Felix Bloch in *Physics Today*, December 1976.
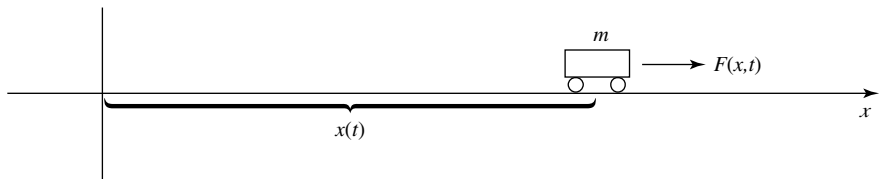
**Figure 1.1:** A "particle" constrained to move in one dimension under the influence of a specified force.
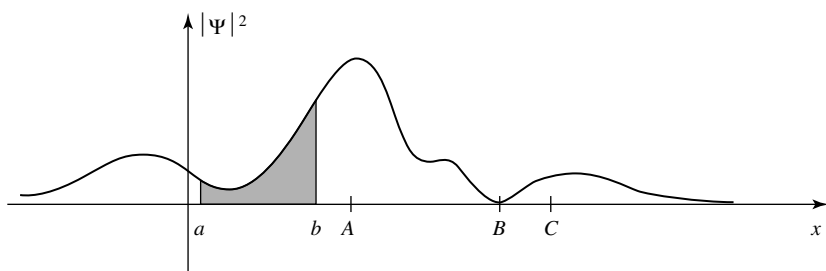


**Figure 1.2:** A typical wave function. The shaded area represents the probability of finding the particle between $a$ and $b$. The particle would be relatively likely to be found near $A$, and unlikely to be found near $B$.

suggests) is spread out in space (it's a function of $x$, for any given $t$). How can such an object represent the state of a *particle*? The answer is provided by Born's **statistical interpretation**, which says that $|\Psi(x, t)|^2$ gives the *probability* of finding the particle at point $x$, at time $t$—or, more precisely,[3]

$$\int_a^b |\Psi(x, t)|^2 \, dx = \left\{ \begin{array}{l} \text{probability of finding the particle} \\ \text{between } a \text{ and } b, \text{ at time } t. \end{array} \right\} \tag{1.3}$$

Probability is the *area* under the graph of $|\Psi|^2$. For the wave function in Figure 1.2, you would be quite likely to find the particle in the vicinity of point $A$, where $|\Psi|^2$ is large, and relatively *un*likely to find it near point $B$.

The statistical interpretation introduces a kind of **indeterminacy** into quantum mechanics, for even if you know everything the theory has to tell you about the particle (to wit: its wave function), still you cannot predict with certainty the outcome of a simple experiment to measure its position—all quantum mechanics has to offer is *statistical* information about the *possible* results. This indeterminacy has been profoundly disturbing to physicists and philosophers alike, and it is natural to wonder whether it is a fact of nature, or a defect in the theory.

Suppose I *do* measure the position of the particle, and I find it to be at point $C$.[4] *Question:* Where was the particle just *before* I made the measurement? There are three plausible answers

---

[3] The wave function itself is complex, but $|\Psi|^2 = \Psi^* \Psi$ (where $\Psi^*$ is the complex conjugate of $\Psi$) is real and non-negative—as a probability, of course, *must* be.

[4] Of course, no measuring instrument is perfectly precise; what I *mean* is that the particle was found *in the vicinity of* $C$, as defined by the precision of the equipment.

to this question, and they serve to characterize the main schools of thought regarding quantum indeterminacy:

1. The **realist** position: *The particle was at C.* This certainly seems reasonable, and it is the response Einstein advocated. Note, however, that if this is true then quantum mechanics is an *incomplete* theory, since the particle *really was* at $C$, and yet quantum mechanics was unable to tell us so. To the realist, indeterminacy is not a fact of nature, but a reflection of our ignorance. As d'Espagnat put it, "the position of the particle was never indeterminate, but was merely unknown to the experimenter."[5] Evidently $\Psi$ is not the whole story—some additional information (known as a **hidden variable**) is needed to provide a complete description of the particle.

2. The **orthodox** position: *The particle wasn't really anywhere.* It was the act of measurement that forced it to "take a stand" (though how and why it decided on the point $C$ we dare not ask). Jordan said it most starkly: "Observations not only *disturb* what [is] to be measured, they *produce* it . . . We *compel* [the particle] to assume a definite position."[6] This view (the so-called **Copenhagen interpretation**), is associated with Bohr and his followers. Among physicists it has always been the most widely accepted position. Note, however, that if it is correct there is something very peculiar about the act of measurement—something that almost a century of debate has done precious little to illuminate.

3. The **agnostic** position: *Refuse to answer.* This is not quite as silly as it sounds—after all, what sense can there be in making assertions about the status of a particle *before* a measurement, when the only way of knowing whether you were right is precisely to *make* a measurement, in which case what you get is no longer "before the measurement"? It is metaphysics (in the pejorative sense of the word) to worry about something that cannot, by its nature, be tested. Pauli said: "One should no more rack one's brain about the problem of whether something one cannot know anything about exists all the same, than about the ancient question of how many angels are able to sit on the point of a needle."[7] For decades this was the "fall-back" position of most physicists: they'd try to sell you the orthodox answer, but if you were persistent they'd retreat to the agnostic response, and terminate the conversation.

Until fairly recently, all three positions (realist, orthodox, and agnostic) had their partisans. But in 1964 John Bell astonished the physics community by showing that it makes an *observable* difference whether the particle had a precise (though unknown) position prior to

---

[5]  Bernard d'Espagnat, "The Quantum Theory and Reality" (*Scientific American*, November 1979, p. 165).
[6]  Quoted in a lovely article by N. David Mermin, "Is the moon there when nobody looks?" (*Physics Today*, April 1985, p. 38).
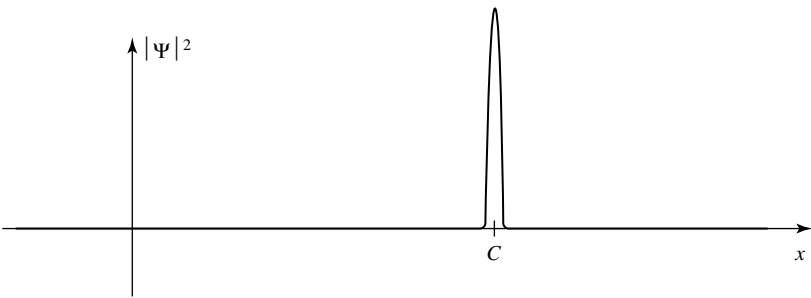[7]  Ibid., p. 40.

**Figure 1.3:** Collapse of the wave function: graph of $|\Psi|^2$ immediately *after* a measurement has found the particle at point $C$.

the measurement, or not. Bell's discovery effectively eliminated agnosticism as a viable option, and made it an *experimental* question whether 1 or 2 is the correct choice. I'll return to this story at the end of the book, when you will be in a better position to appreciate Bell's argument; for now, suffice it to say that the experiments have decisively confirmed the orthodox interpretation:[8] a particle simply does not have a precise position prior to measurement, any more than the ripples on a pond do; it is the measurement process that insists on one particular number, and thereby in a sense *creates* the specific result, limited only by the statistical weighting imposed by the wave function.

What if I made a *second* measurement, immediately after the first? Would I get $C$ again, or does the act of measurement cough up some completely new number each time? On this question everyone is in agreement: A repeated measurement (on the same particle) must return the same value. Indeed, it would be tough to prove that the particle was really found at $C$ in the first instance, if this could not be confirmed by immediate repetition of the measurement. How does the orthodox interpretation account for the fact that the second measurement is bound to yield the value $C$? It must be that the first measurement radically alters the wave function, so that it is now sharply peaked about $C$ (Figure 1.3). We say that the wave function **collapses**, upon measurement, to a spike at the point $C$ (it soon spreads out again, in accordance with the Schrödinger equation, so the second measurement must be made quickly). There are, then, two entirely distinct kinds of physical processes: "ordinary" ones, in which the wave function evolves in a leisurely fashion under the Schrödinger equation, and "measurements," in which $\Psi$ suddenly and discontinuously collapses.[9]

---

[8] This statement is a little too strong: there exist viable nonlocal hidden variable theories (notably David Bohm's), and other formulations (such as the **many worlds** interpretation) that do not fit cleanly into any of my three categories. But I think it is wise, at least from a pedagogical point of view, to adopt a clear and coherent platform at this stage, and worry about the alternatives later.

[9] The role of measurement in quantum mechanics is so critical and so bizarre that you may well be wondering what precisely *constitutes* a measurement. I'll return to this thorny issue in the Afterword; for the moment let's take the naive view: a measurement is the kind of thing that a scientist in a white coat does in the laboratory, with rulers, stopwatches, Geiger counters, and so on.

**Example 1.1**

**Electron Interference.** I have asserted that particles (electrons, for example) have a wave nature, encoded in $\Psi$. How might we check this, in the laboratory?

The classic signature of a wave phenomenon is *interference*: two waves *in phase* interfere constructively, and out of phase they interfere destructively. The wave nature of light was confirmed in 1801 by Young's famous double-slit experiment, showing interference "fringes" on a distant screen when a monochromatic beam passes through two slits. If essentially the same experiment is done with *electrons*, the same pattern develops,[10] confirming the wave nature of electrons.

Now suppose we decrease the intensity of the electron beam, until only one electron is present in the apparatus at any particular time. According to the statistical interpretation each electron will produce a spot on the screen. Quantum mechanics cannot predict the precise *location* of that spot—all it can tell us is the *probability* of a given electron landing at a particular place. But if we are patient, and wait for a hundred thousand electrons—one at a time—to make the trip, the accumulating spots reveal the classic two-slit interference pattern (Figure 1.4).[11]
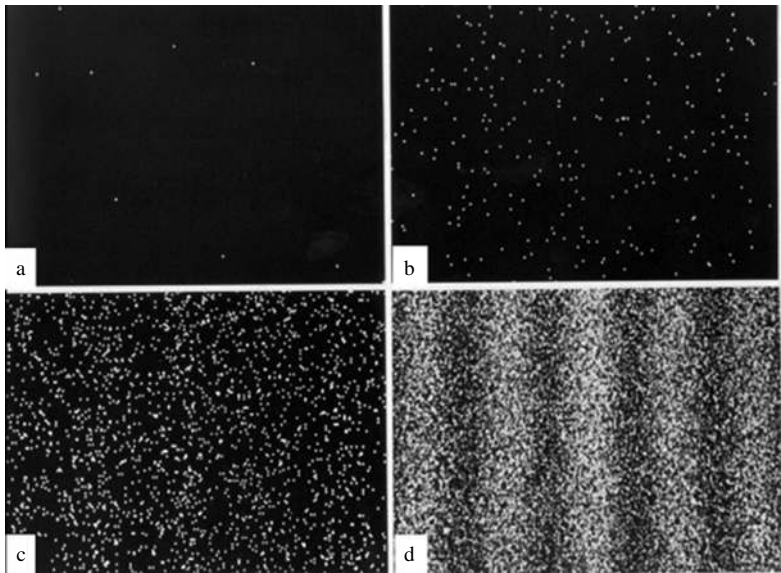


**Figure 1.4:** Build-up of the electron interference pattern. (a) Eight electrons, (b) 270 electrons, (c) 2000 electrons, (d) 160,000 electrons. Reprinted courtesy of the Central Research Laboratory, Hitachi, Ltd., Japan.

---

[10] Because the wavelength of electrons is typically very small, the slits have to be extremely close together. Historically, this was first achieved by Davisson and Germer, in 1925, using the atomic layers in a crystal as "slits." For an interesting account, see R. K. Gehrenbeck, *Physics Today*, January 1978, page 34.

[11] See Tonomura et al., *American Journal of Physics*, Volume 57, Issue 2, pp. 117–120 (1989), and the amazing associated video at www.hitachi.com/rd/portal/highlight/quantum/doubleslit/. This experiment can now be done with much more massive particles, including "Bucky-balls"; see M. Arndt, et al., *Nature* **401**, 680 (1999). Incidentally, the same thing can be done with light: turn the intensity so low that only one "photon" is present at a time and you get an identical point-by-point assembly of the interference pattern. See R. S. Aspden, M. J. Padgett, and G. C. Spalding, *Am. J. Phys.* **84**, 671 (2016).

Of course, if you close off one slit, or somehow contrive to detect which slit each electron passes through, the interference pattern disappears; the wave function of the emerging particle is now entirely different (in the first case because the boundary conditions for the Schrödinger equation have been changed, and in the second because of the collapse of the wave function upon measurement). But with both slits open, and no interruption of the electron in flight, each electron interferes with itself; it didn't pass through one slit or the other, but through both at once, just as a water wave, impinging on a jetty with two openings, interferes with itself. There is nothing mysterious about this, once you have accepted the notion that particles obey a wave equation. The truly *astonishing* thing is the blip-by-blip assembly of the pattern. In any classical wave theory the pattern would develop smoothly and continuously, simply getting more intense as time goes on. The quantum process is more like the pointillist painting of Seurat: The picture emerges from the cumulative contributions of all the individual dots.[12]

## 1.3  PROBABILITY

### 1.3.1  Discrete Variables

Because of the statistical interpretation, probability plays a central role in quantum mechanics, so I digress now for a brief discussion of probability theory. It is mainly a question of introducing some notation and terminology, and I shall do it in the context of a simple example.

Imagine a room containing fourteen people, whose ages are as follows:

one person aged 14,
one person aged 15,
three people aged 16,
two people aged 22,
two people aged 24,
five people aged 25.

If we let $N(j)$ represent the number of people of age $j$, then

$N(14) = 1,$
$N(15) = 1,$
$N(16) = 3,$
$N(22) = 2,$
$N(24) = 2,$
$N(25) = 5,$

---

[12] I think it is important to distinguish things like interference and diffraction that would hold for any wave theory from the uniquely quantum mechanical features of the measurement process, which derive from the statistical interpretation.
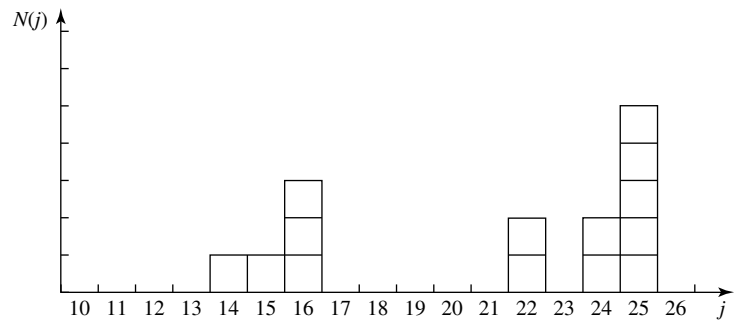
**Figure 1.5:** Histogram showing the number of people, $N(j)$, with age $j$, for the example in Section 1.3.1.

while $N(17)$, for instance, is zero. The *total* number of people in the room is

$$N = \sum_{j=0}^{\infty} N(j). \tag{1.4}$$

(In the example, of course, $N = 14$.) Figure 1.5 is a histogram of the data. The following are some questions one might ask about this distribution.

**Question 1**   If you selected one individual at random from this group, what is the **probability** that this person's age would be 15?

**Answer**   One chance in 14, since there are 14 possible choices, all equally likely, of whom only one has that particular age. If $P(j)$ is the probability of getting age $j$, then $P(14) = 1/14$, $P(15) = 1/14$, $P(16) = 3/14$, and so on. In general,

$$P(j) = \frac{N(j)}{N}. \tag{1.5}$$

Notice that the probability of getting *either* 14 *or* 15 is the *sum* of the individual probabilities (in this case, 1/7). In particular, the sum of *all* the probabilities is 1—the person you select must have *some* age:

$$\sum_{j=0}^{\infty} P(j) = 1. \tag{1.6}$$

**Question 2**   What is the **most probable** age?

**Answer**   25, obviously; five people share this age, whereas at most three have any other age. The most probable $j$ is the $j$ for which $P(j)$ is a maximum.

**Question 3**   What is the **median** age?

**Answer**   23, for 7 people are younger than 23, and 7 are older. (The median is that value of $j$ such that the probability of getting a larger result is the same as the probability of getting a smaller result.)

**Question 4**   What is the **average** (or **mean**) age?

**Answer**

$$\frac{(14) + (15) + 3(16) + 2(22) + 2(24) + 5(25)}{14} = \frac{294}{14} = 21.$$

In general, the average value of $j$ (which we shall write thus: $\langle j \rangle$) is

$$\langle j \rangle = \frac{\sum j N(j)}{N} = \sum_{j=0}^{\infty} j P(j).$$  (1.7)

Notice that there need not be anyone with the average age or the median age—in this example nobody happens to be 21 or 23. In quantum mechanics the average is usually the quantity of interest; in that context it has come to be called the **expectation value**. It's a misleading term, since it suggests that this is the outcome you would be most likely to get if you made a single measurement (*that* would be the *most probable value*, not the average value)—but I'm afraid we're stuck with it.

**Question 5**    What is the average of the *squares* of the ages?
**Answer**    You could get $14^2 = 196$, with probability 1/14, or $15^2 = 225$, with probability 1/14, or $16^2 = 256$, with probability 3/14, and so on. The average, then, is

$$\langle j^2 \rangle = \sum_{j=0}^{\infty} j^2 P(j).$$  (1.8)

In general, the average value of some *function* of $j$ is given by

$$\boxed{\langle f(j) \rangle = \sum_{j=0}^{\infty} f(j) P(j).}$$  (1.9)

(Equations 1.6, 1.7, and 1.8 are, if you like, special cases of this formula.) *Beware:* The average of the squares, $\langle j^2 \rangle$, is *not* equal, in general, to the square of the average, $\langle j \rangle^2$. For instance, if the room contains just two babies, aged 1 and 3, then $\langle j^2 \rangle = 5$, but $\langle j \rangle^2 = 4$.

Now, there is a conspicuous difference between the two histograms in Figure 1.6, even though they have the same median, the same average, the same most probable value, and the same number of elements: The first is sharply peaked about the average value, whereas the second is broad and flat. (The first might represent the age profile for students in a big-city classroom, the second, perhaps, a rural one-room schoolhouse.) We need a numerical measure
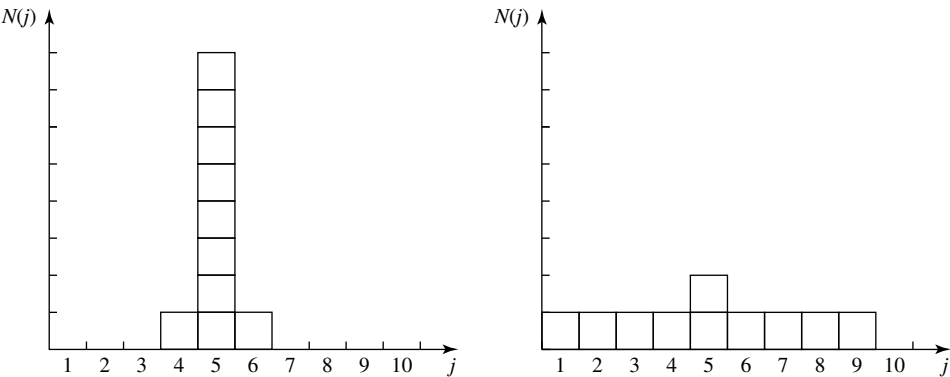


**Figure 1.6:** Two histograms with the same median, same average, and same most probable value, but different standard deviations.

of the amount of "spread" in a distribution, with respect to the average. The most obvious way to do this would be to find out how far each individual is from the average,

$$\Delta j = j - \langle j \rangle, \tag{1.10}$$

and compute the average of $\Delta j$. Trouble is, of course, that you get *zero*:

$$\langle \Delta j \rangle = \sum (j - \langle j \rangle) P(j) = \sum j P(j) - \langle j \rangle \sum P(j)$$
$$= \langle j \rangle - \langle j \rangle = 0.$$

(Note that $\langle j \rangle$ is constant—it does not change as you go from one member of the sample to another—so it can be taken outside the summation.) To avoid this irritating problem you might decide to average the *absolute value* of $\Delta j$. But absolute values are nasty to work with; instead, we get around the sign problem by *squaring* before averaging:

$$\sigma^2 \equiv \langle (\Delta j)^2 \rangle. \tag{1.11}$$

This quantity is known as the **variance** of the distribution; $\sigma$ itself (the square root of the average of the square of the deviation from the average—gulp!) is called the **standard deviation**. The latter is the customary measure of the spread about $\langle j \rangle$.

There is a useful little theorem on variances:

$$\sigma^2 = \langle (\Delta j)^2 \rangle = \sum (\Delta j)^2 P(j) = \sum (j - \langle j \rangle)^2 P(j)$$
$$= \sum \left( j^2 - 2j \langle j \rangle + \langle j \rangle^2 \right) P(j)$$
$$= \sum j^2 P(j) - 2 \langle j \rangle \sum j P(j) + \langle j \rangle^2 \sum P(j)$$
$$= \langle j^2 \rangle - 2 \langle j \rangle \langle j \rangle + \langle j \rangle^2 = \langle j^2 \rangle - \langle j \rangle^2.$$

Taking the square root, the standard deviation itself can be written as

$$\sigma = \sqrt{\langle j^2 \rangle - \langle j \rangle^2}. \tag{1.12}$$

In practice, this is a much faster way to get $\sigma$ than by direct application of Equation 1.11: simply calculate $\langle j^2 \rangle$ and $\langle j \rangle^2$, subtract, and take the square root. Incidentally, I warned you a moment ago that $\langle j^2 \rangle$ is not, in general, equal to $\langle j \rangle^2$. Since $\sigma^2$ is plainly non-negative (from its definition 1.11), Equation 1.12 implies that

$$\langle j^2 \rangle \geq \langle j \rangle^2, \tag{1.13}$$

and the two are equal only when $\sigma = 0$, which is to say, for distributions with no spread at all (every member having the same value).

### 1.3.2 Continuous Variables

So far, I have assumed that we are dealing with a *discrete* variable—that is, one that can take on only certain isolated values (in the example, $j$ had to be an integer, since I gave ages only in years). But it is simple enough to generalize to *continuous* distributions. If I select a random person off the street, the probability that her age is *precisely* 16 years, 4 hours, 27 minutes, and 3.333... seconds is *zero*. The only sensible thing to speak about is the probability that her age lies in some *interval*—say, between 16 and 17. If the interval is sufficiently short, this probability is *proportional to the length of the interval*. For example, the chance that her age is between 16 and 16 plus *two* days is presumably twice the probability that it is between 16 and

16 plus *one* day. (Unless, I suppose, there was some extraordinary baby boom 16 years ago, on exactly that day—in which case we have simply chosen an interval too long for the rule to apply. If the baby boom lasted six hours, we'll take intervals of a second or less, to be on the safe side. Technically, we're talking about *infinitesimal* intervals.) Thus

$$\left\{ \begin{array}{l} \text{probability that an individual (chosen} \\ \text{at random) lies between } x \text{ and } (x + dx) \end{array} \right\} = \rho(x)\, dx. \tag{1.14}$$

The proportionality factor, $\rho(x)$, is often loosely called "the probability of getting $x$," but this is sloppy language; a better term is **probability density**. The probability that $x$ lies between $a$ and $b$ (a *finite* interval) is given by the integral of $\rho(x)$:

$$P_{ab} = \int_a^b \rho(x)\, dx, \tag{1.15}$$

and the rules we deduced for discrete distributions translate in the obvious way:

$$\int_{-\infty}^{+\infty} \rho(x)\, dx = 1, \tag{1.16}$$

$$\langle x \rangle = \int_{-\infty}^{+\infty} x\rho(x)\, dx, \tag{1.17}$$

$$\langle f(x) \rangle = \int_{-\infty}^{+\infty} f(x)\, \rho(x)\, dx, \tag{1.18}$$

$$\sigma^2 \equiv \left\langle (\Delta x)^2 \right\rangle = \left\langle x^2 \right\rangle - \langle x \rangle^2. \tag{1.19}$$

---

**Example 1.2**

Suppose someone drops a rock off a cliff of height $h$. As it falls, I snap a million photographs, at random intervals. On each picture I measure the distance the rock has fallen. *Question:* What is the *average* of all these distances? That is to say, what is the *time average* of the distance traveled?[13]

**Solution:** The rock starts out at rest, and picks up speed as it falls; it spends more time near the top, so the average distance will surely be less than $h/2$. Ignoring air resistance, the distance $x$ at time $t$ is

$$x(t) = \frac{1}{2}gt^2.$$

The velocity is $dx/dt = gt$, and the total flight time is $T = \sqrt{2h/g}$. The probability that a particular photograph was taken between $t$ and $t + dt$ is $dt/T$, so the probability that it shows a distance in the corresponding range $x$ to $x + dx$ is

$$\frac{dt}{T} = \frac{dx}{gt}\sqrt{\frac{g}{2h}} = \frac{1}{2\sqrt{hx}}\, dx.$$

---

[13] A statistician will complain that I am confusing the average of a *finite sample* (a million, in this case) with the "true" average (over the whole continuum). This can be an awkward problem for the experimentalist, especially when the sample size is small, but here I am only concerned with the *true* average, to which the sample average is presumably a good approximation.