1.1 A science of mind

We spend an enormous number of our waking hours thinking and talking about our thoughts, emotions, and experiences. For example, we wonder: Why did the waiter give me that unusual smile? Did my co-worker see me stealing those office supplies? How can I deflect my unwanted admirer's attention – or attract the attention of someone else? In trying to answer such questions, and in interpreting one another's behavior more generally, we make use of a vast body of lore about how people perceive, reason, desire, feel, and so on. So we say such things as: the waiter is smiling obsequiously because he hopes I will give him a larger tip; my co-worker does know, but he won't tell anyone, because he's afraid I'll reveal his gambling problem; and so on. Formulating such explanations is part of what enables us to survive in a shared social environment.

This everyday understanding of our minds, and those of others, is referred to as "folk psychology." The term is usually taken as picking out our ability to attribute psychological states and to use those attributions for a variety of practical ends, including prediction, explanation, manipulation, and deception. It encompasses our ability to verbally produce accounts couched in the everyday psychological vocabulary with which most of us are conversant: the language of beliefs, desires, intentions, fears, hopes, and so on. Such accounts are the stuff of which novels and gossip are made. Although our best evidence for what people think is often what they say, much of our capacity to read the thoughts of others may also be nonverbal, involving the ability to tell moods and intentions immediately by various bodily cues – an ability we may not be conscious that we have.

Although we have an important stake in the success of our folk psychological attributions and explanations, and while social life as we know it would

be impossible without folk psychology, folk psychology also has obvious shortcomings (Churchland, 1981). Our accounts of one another's behavior are often sketchy, unsystematic, or of merely local utility. Moreover, they leave out whole ranges of abnormal mental phenomena such as autism or Capgras syndrome. We have no folk explanation for how we are able to perceive and navigate our way through a three-dimensional space cluttered with objects, how we integrate what we see with what we hear and touch, how we are able to learn language, how we recognize faces and categories, how our memory works, how we reason and make decisions, and so on. The explanations of these varied mental capacities lie far beyond folk psychology's province. If we want to understand the mind, then we need to find better ways to investigate its structure and function. The sciences of the mind have developed in response to this need.

Science aims at systematic understanding of the world, and psychology is the science that takes mental phenomena in general as its domain. This definition has not always been uncontroversially accepted. Behaviorists such as Watson (1913) and Skinner (1965) held that the only proper subject matter for psychology was the domain of observable behavior, in part on the grounds that minds were mysterious and inaccessible to third-person methods of investigation. Few today take this position. Mental states and processes may not be directly observable, but they can be inferred by a variety of converging techniques. Cognitive psychology in particular typically proceeds by positing such inferred states. Many of these states such as occurrent perceptions and thoughts are accessible via introspection with varying degrees of accuracy, but many are entirely unconscious.

"Phenomena" is a cover term for the body of noteworthy natural regularities to be found in the objects, events, processes, activities, and capacities that a science concerns itself with.¹ Objects can include such things as whole organisms (white rats, the sea slug *Aplysia californica*), artificial behaving systems (a trained neural network, an autonomous mobile robot), or their parts (the brain, particular brain structures such as the hippocampus or the supplementary motor area, a particular control structure in a computer). Here the relevant phenomena are reliable patterns of organization or behavior in these objects – for example, the predictable laminar organization and connectivity patterns in the neocortex. Events and processes include any changes

¹ This usage follows Hacking (1983). See also Bogen and Woodward (1988).

1.1 A science of mind

3

undergone by these objects: the myelination of the frontal lobes in normal development, a rat's learning to run a water maze, a child acquiring the lexicon of her first language, an undergraduate carrying out a motor task in response to a visual stimulus, a patient with dementia retrieving a memory of an event from his teenage years. Activities and capacities include any functions that an object can reliably carry out. Normal humans have the capacity to rapidly estimate quantity, to selectively attend to parts of a complex visual array, to judge which of two events is more likely, to generate expectations about the movement of simple physical objects in their environment, to attribute emotional states to others, and so on.

Mental phenomena encompass attention, learning and memory, concept acquisition and categorization, language acquisition, perception (both accurate and illusory), and emotions and moods, among others. We won't try to be exhaustive. Traditional distinctions among types of mental states have been made along the following lines. Some mental states involve concepts in their formation, expression, and function. These are the types of states associated with higher cognition and knowledge (from which "cognitive" derives its name). Such states include beliefs, desires, intentions, hopes, and plain old thoughts in general. Other sorts of states, such as sensory states, do not necessarily involve concepts in their activation. One can smell a rose without knowing that it is a rose one smells. One can hear a C-sharp on the piano without knowing that it is a C-sharp one hears. Emotions such as fear, love, and anger also form a distinctive class of mental states. Finally, there are moods: general overall feelings of excitement, happiness, sadness, mania, and depression.

Is there anything that all mental phenomena have in common? This is controversial, but one proposal is that they are all *representational*.² The higher cognitive states that involve concepts clearly involve representations that can fit into propositional attitudes and generate knowledge of various facts and states of affairs. Sensory states do not necessarily involve the activation of concepts, but they are still a type of representation on at least some views. They represent the presence of a physically perceptible property and the causal interaction of that property with a sensory system of the body. The sweet taste of sugar represents the interaction of the sugar molecules with

² We discuss the issue of how to distinguish mental phenomena in greater depth in Section 5.4.4.

the taste receptors in the mouth, for instance. Even moods have been portrayed as representations of general chemical states or changes in the body.

One goal of the sciences is to describe, clarify, and organize these phenomena. Consider the changes that the past 50 years have wrought in our understanding of the cognitive capacities of infants and young children, for example. At some point, normal children become able to understand and interpret the behavior of others in terms of their beliefs, intentions, and desires. In a pioneering study, Wimmer and Perner (1983) showed that four-year-olds are able to correctly predict how characters with false beliefs will act, whereas younger children are unable to do so. In one of their now-classic tasks, the child watches one puppet place a piece of candy in a certain location and then leave the room. The other puppet, which was present when the candy was hidden, now moves it to a new hidden location. The first puppet then returns, and the child is asked either where she will look for the candy or where she thinks the candy is. Passing this so-called false belief task involves correctly saying that she will look in the original location, rather in the actual location, since she will be guided not by the candy's actual location, but by her erroneous beliefs about it. Here the phenomenon of interest is the alleged shift from failure to success in this particular test (and related variants). This result was widely interpreted as showing that some components of "theory of mind" - those connected with the attribution of beliefs – are not yet in place prior to age four.³

Surprisingly, though, in recent years it has been shown that even 15-montholds can respond in a way that seems to display understanding of false beliefs (Onishi & Baillargeon, 2005). These infants will look longer at a scene depicting a character searching in a place that she could not know an object is located (because she had earlier seen it hidden elsewhere) than at a scene in which she searched for it in the place where she should expect it to be. Looking time in infants is often taken to be an indicator of surprise or violation of expectancy, an interpretation confirmed by studies across many different stimuli and domains. Thus the 15-month-olds in this study don't seem to expect the characters to have information about the true state of the world; this strongly suggests that they naturally attribute something like false beliefs. Moreover, 16-month-olds will even act on this understanding, trying to help out individuals who are attempting to act on false beliefs by

³ For much more on theory of mind, see Chapter 8.

1.1 A science of mind

5

pointing to the correct location of a hidden toy (Buttelmann, Carpenter, & Tomasello, 2009).

This case illustrates two points. First, what the phenomena are in psychology, as in other sciences, is often nonobvious. That is, one cannot, in general, simply look and see that a certain pattern or regularity exists. Experiment and measurement are essential for the production of many interesting psychological phenomena. Second, phenomena are almost always tied closely to experimental tasks or paradigms. The phenomenon of three-year-olds failing the false belief task and four-year-olds passing it depends greatly on *which* false belief task one uses. If we agree to call the nonverbal Onishi and Baillargeon paradigm a false belief task, we need to explain the seeming contradiction between the phenomena, perhaps in terms of the differing requirements of the tasks (Bloom & German, 2000). Individuating phenomena is intimately tied to individuating tasks and experimental methods.

To see this, consider the Stroop effect. In his classic paper, Stroop (1935) performed three experiments, the first two of which are the most well known. In experiment 1, he asked participants to read color names printed in a variety of differently colored inks. The names were given in a 10×10 grid, and no name was ever paired with the color of ink that it named. The control condition required reading the same names printed in black ink. Subtracting the time to read the experimental versus the control cards, Stroop found that on average it took slightly longer to read the color names printed in differently colored ink, but this difference was not significant. In experiment 2, he required participants to name the color of the ink in the experimental condition, rather than reading the color name. In the control condition, words were replaced with colored squares. Here the difference in reading times was striking: participants were 74% slower to name the ink color from a sample. Conflicting lexical information interferes with color naming.

Although this is the canonical "Stroop effect," the term has been broadened over time to include a range of related phenomena. Stroop-like tasks have been carried out using pictures or numbers versus words, using auditory rather than visual materials, using nonverbal response measures, and so on. Further manipulations have involved varying the time at which the conflicting stimulus is presented (e.g., showing the color sample before the word), and the effect persists. Wherever responding to one kind of information interferes asymmetrically with responding to another that is simultaneously

presented, we have a Stroop-like phenomenon. Much of the literature on the effect has focused on delineating the precise sorts of stimuli, tasks, and populations that display the effect (MacLeod, 1991). But the effect itself is elusive outside the context of these experimental manipulations – certainly it is not a straightforwardly observable behavioral regularity on a par with wincing in response to being kicked. More esoteric phenomena may be reliant on even more sophisticated experimental setups for their elicitation.

In these cases, what psychologists are primarily aiming to do is to *characterize* the phenomena. This may require deploying new experimental paradigms, modifying the parameters of old paradigms, or refining techniques of data collection and analysis. The phenomena themselves are dependent on these techniques of investigation for their existence. Producing and measuring these phenomena involve discovering how various parts of the psychological domain behave when placed in relatively artificial circumstances, under the assumption that this will be importantly revealing about their normal structure and function. This is perhaps the biggest advantage scientific psychology has over its folk counterpart, which tends to be resolutely nonexperimental.

But beyond producing and describing phenomena – that is, saying *what* happens in the world – psychology also aims to explain *how* and *why* they are produced. Where we are dealing with genuine, robust phenomena, we assume, as an initial hypothesis at least, that they are not merely accidental. There ought to be some reason why they exist and take the particular form that they do. It is sometimes maintained that what is distinctive about scientific theorizing, as opposed to other ways of reasoning about the world, is that it involves positing and testing explanations. As we have seen, this can't be the whole story, because making and refining ways in which we might better describe the world are themselves major parts of the scientific enterprise. But the psychological phenomena we discover often turn out to be novel or surprising. Hence better descriptions of the phenomena naturally tend to pull us toward generating explanations for their existence.

1.2 Explanations in psychology

We shouldn't assume that all sciences will deploy the same explanatory strategies. What works to explain geological or astronomical phenomena may not work for psychological phenomena. So we begin by considering four sample

1.2 Explanations in psychology

7

cases of psychological explanation. We should note that these explanations are to varying degrees contested, but the present issue is what they can tell us about the structure of explanations in psychology, rather than whether they are strictly accurate.

1.2.1 Case 1: Psychophysics

Some of the earliest systematic psychological research in the nineteenth century concerned psychophysical phenomena, in particular how the properties of sensations depend on and vary with the properties of the physical stimulus that produces them. Light, sound waves, pressure, temperature, and other ambient energy sources interact with sensory receptors and their associated processing systems to give rise to sensations, and this relationship is presumably systematic rather than random. To uncover this hidden order, early psychophysicists had to solve three problems simultaneously: (1) how to devise empirical strategies for measuring sensations, (2) how to quantify the ways in which those sensations covaried with stimulus conditions, and, finally, (3) how to explain those covariations.

Fechner (1860), following the work of Weber (1834), hit on the method of using "just noticeable differences" (jnd's) to measure units of sensation. A stimulus in some sensory modality (e.g., a patch of light, a tone) is increased in intensity until the perceiver judges that there is a detectable change in the quality of her sensations. The measure of a jnd in physical terms is the difference between the initial and final stimulus magnitude. By increasing stimulus intensity until the next jnd was reached, Fechner could plot the intervals at which a detectable change in a sensation occurred against the stimulus that caused the change.

After laboriously mapping stimulus–sensation pairs in various modalities, Fechner proposed a logarithmic law to capture their relationship formally. Fechner's law states:

$S = k \log(I)$

where *S* is the perceived magnitude of the sensation (e.g., the brightness of a light or the loudness of a sound), *I* is the intensity of the physical stimulus, and k is an empirically determined constant. Because this is a logarithmic

law, geometric increases in stimulus intensity will correspond to arithmetic increases in the strength of sensations.

Although Fechner's law delivers predictions that conform with much of the data, it also fails in some notable cases. Stevens (1957) took a different experimental approach. Rather than constructing scales using jnd's, he asked participants to directly estimate magnitudes of various stimuli using arbitrary numerical values. So an initial stimulus would be given a numerical value, and then later stimuli were given values relative to it, where all of the numerical assignments were freely chosen by the participants. He also asked them to directly estimate stimulus ratios, such as when one stimulus seemed to be twice as intense as another. Using these methods, he showed that the perceived intensity of some stimuli departed from Fechner's law. He concluded that Fechner's assumption that all jnd's are of equal size was to blame for the discrepancy and proposed as a replacement for Fechner's law the power law (now known as Stevens' law):

 $S = kI^a$

where *S* and *I* are perceived magnitude and physical intensity, *k* is a constant, and *a* is an exponent that differs for various sensory modalities and perceivable quantities. The power law predicts that across all quantities and modalities, equal stimulus ratios correspond to equal sensory ratios, and, depending on the exponent, perceived magnitudes may increase more quickly or more slowly than the increase in stimulus intensity.

Stevens (1975, pp. 17–19) gave an elegant argument for why we should expect sensory systems in general to obey a power law. He noted that as we move around and sense the environment, the absolute magnitudes we perceive will vary: the visual angle subtended by the wall of a house changes as one approaches it; the intensity of speech sounds varies as one approaches or recedes. What is important in these cases is not the differences in the stimulus, but the constants, which are given by the ratios that the elements of the stimulus bear to one another. A power law is well suited to capture this, because equal ratios of stimulus intensity correspond to equal ratios of sensory magnitude.

Stevens' law provides a generally better fit for participants' judgments about magnitudes and therefore captures the phenomena of stimulussensation relations better than Fechner's law, although it, too, is only

1.2 Explanations in psychology

9

approximate.⁴ However, both laws provide the same sort of explanation for the relationship between the two: in each case, the laws show that these relationships are not arbitrary, but instead conform to a general formula, which can be expressed by a relatively simple equation. The laws explain the phenomena by showing how they can all be systematically related in a simple, unified fashion. Once we have the law in hand, we are in a position to make predictions about the relationship between unmeasured magnitudes, to the effect that they will probably conform to the regularity set out in the law (even if the precise form of the regularity requires empirically determining the values of k and a).

1.2.2 Case 2: Classical conditioning

Any organism that is to survive for long in an environment with potentially changing conditions needs some way of learning about the structure of events in its environment. Few creatures lead such simple lives that they can be born "knowing" all they will need to survive. The investigation of learning in animals (and later humans) started with the work of Pavlov, Skinner, Hull, and other behaviorists. Given their aversion to mentalistic talk, they tended to think of learning as a measurable change in the observable behavior of a creature in response to some physical stimulus or other. The simplest style of learning is classical (Pavlovian) conditioning. In classical conditioning, we begin with an organism that reliably produces a determinate type of response to a determinate type of stimulus - for example, flinching in response to a mild shock, or blinking in response to a puff of air. The stimulus here is called the unconditioned stimulus (US), and the response the unconditioned response (UR). In a typical experiment, the US is paired with a novel, neutral stimulus (e.g., a flash of light or a tone) for a training period; this is referred to as the conditioned stimulus (CS). After time, under the right training conditions, the CS becomes associated with the US, so that the CS is capable of producing the response by itself; when this occurs, it is called the conditioned response (CR).

There were a number of early attempts to formulate descriptions of how conditioning takes place (Bush & Mosteller, 1951; Hull, 1943). These descriptions take the form of learning rules that predict how the strength of

⁴ For useful discussion on the history and logic of various psychophysical scaling procedures, see Shepard (1978) and Gescheider (1988).

associations among CS and US will change over time under different training regimes. One of the most well-known and best empirically validated learning rules was the "delta rule" presented by Rescorla and Wagner (1972). Formally, the rule says:⁵

$$\Delta \mathbf{A}_{ij} = \alpha_i \beta_j (\lambda_j - \Sigma_i \mathbf{A}_{ij})$$

To grasp what this means, suppose we are on training trial n, and we want to know what the associative strengths will be at the next stage n + 1. Let i stand for the CS and j stand for the US. Then A_{ij} is the strength of the association between i and j, and ΔA_{ij} is the change in the strength of that association as a result of training. The terms α_i and β_j are free parameters that determine the rate at which learning can take place involving the CS and US. The term λ_j is the maximum associative strength that the US can support. Finally, $\Sigma_i A_{ij}$ is the sum of the strength of all of the active CSs that are present during trial n. This is needed because some learning paradigms involve presenting multiple CSs at the same time during training.

The essence of the Rescorla–Wagner rule is to reduce the "surprisingness" of a US. If a CS (i) is not associated strongly with a US (j), then (assuming no other CSs are present), the parenthetical term of the rule will be large, and so the strength of the association between i and j will be correspondingly adjusted. Over time, as its association with the CS increases, the surprisingness of the US decreases, and so less change in strength takes place.

The Rescorla–Wagner rule is one of the most extensively studied learning rules in psychology, and it has some significant virtues: it unifies a large range of phenomena by bringing them under a single, relatively simple formal description; it explains previously discovered phenomena; and it generates surprising and often-confirmed predictions about new phenomena. To get the flavor of this, consider some of its successes: (1) The rule explains why acquisition curves show less change over time, for the reason given in the previous paragraph. (2) Extinction is the loss of response to a CS when it is presented without its paired US. The model explains this by positing that

⁵ Gallistel (1990, Chapter 12) gives an excellent critical discussion of the assumptions underlying the R-W rule and its predecessors. He notes that the R-W rule is cast in terms of associative strengths rather than directly observable response probabilities, which represents a significant change of emphasis over earlier behavioristic rules. For a review of some important behavioral findings concerning conditioning, see Rescorla (1988).