Part I   INTRODUCTION

**1**

# Classical conditioning: data and theories

During classical (or Pavlovian) conditioning, human and animal subjects change the magnitude and timing of their conditioned response (CR), as a result of the contingency between the conditioned stimulus (CS) and the unconditioned stimulus (US).

In this chapter we briefly describe results of a number of classical conditioning paradigms that are discussed in detail in different chapters of the book (see Schmajuk, 2008a, 2008b). Then we introduce different types of learning theories. Finally, we present a number of computational models of classical conditioning.

### Classical conditioning data

*A.    Excitatory conditioning*

  1.  Acquisition. After a number of CS–US pairings, the CS elicits a conditioned response (CR) that increases in magnitude and frequency.
  2.  Partial reinforcement. The US follows the CS only on some trials, and might lead to a lower conditioning asymptote.
  3.  Generalization. A $CS_2$ elicits a CR when it shares some characteristics with a $CS_1$ that has been paired with the US.
  4.  US- and CS-specific CR. The nature of the CR is determined not only by the US, but also by the CS.

*B.    Inhibitory conditioning*

  1.  Conditioned inhibition. Stimulus $CS_2$ acquires inhibitory conditioning with $CS_1$ reinforced trials interspersed with, or followed by, $CS_1$–$CS_2$ nonreinforced trials.

4    Introduction

2.  Extinction of conditioned inhibition. Inhibitory conditioning is extinguished by $CS_2$-US presentations, but not by presentations of $CS_2$ alone.

3.  Differential conditioning. Stimulus $CS_2$ acquires inhibitory conditioning with $CS_1$ reinforced trials interspersed with $CS_2$ nonreinforced trials.

4.  Contingency. A CS becomes inhibitory when the probability that the US will occur in the presence of the CS, p(US/CS), is smaller than the probability that the US will occur in the absence of the CS (p[US/noCS]).

*C.    Preexposure effects*

1.  Latent inhibition (LI). Preexposure to a CS followed by CS–US pairings retard the generation of the CR.

2.  Context preexposure. Preexposure to a context facilitates the acquisition of fear conditioning.

3.  US–preexposure effect. Presentation of the US in a training context, prior to CS–US pairings, retards production of the CR.

4.  Learned irrelevance. Random exposure to the CS and the US retards conditioning even more than combined latent inhibition and US preexposure.

*D.    Compound conditioning*

1.  Relative validity. Conditioning to *X* is weaker when training consists of reinforced *XA* trials alternated with *XB* nonreinforced trials, than when training consists of *XA* trials alternated with *XB* trials each type reinforced half of the time.

2.  Blocking. Conditioning to $CS_1$–$CS_2$ following conditioning to $CS_1$ results in a weaker conditioning to $CS_2$ than that attained with $CS_1$–$CS_2$-US pairings.

3.  Unblocking by increasing the US. Increasing the US during $CS_1$–$CS_2$ conditioning increases responding to the blocked $CS_2$.

4.  Unblocking by decreasing the US. Responding to $CS_2$ can be increased by decreasing the US during $CS_1$–$CS_2$ conditioning.

5.  Overshadowing. Conditioning to $CS_1$–$CS_2$ results in a weaker conditioning to $CS_2$ than that attained with $CS_2$-US pairings.

6.  Potentiation. Conditioning to $CS_1$–$CS_2$ results in a stronger conditioning to $CS_2$ than that attained with $CS_2$-US pairings.

7.  Backward blocking. Conditioning to $CS_1$ following conditioning to $CS_1$–$CS_2$ results in a weaker response to $CS_2$ than that attained with $CS_1$–$CS_2$-US pairings.

8.   Overexpectation. Reinforced $CS_1$–$CS_2$ presentations following inde-
     pendent reinforced $CS_1$ and $CS_2$ presentations, result in a decrement in
     their initial associative strength.

9.   Superconditioning. Reinforced $CS_1$–$CS_2$ presentations following inhibi-
     tory conditioning of $CS_1$ increase $CS_2$ excitatory strength compared
     with the case when it is trained in the absence of $CS_1$.

E.   *Recovery from compound conditioning*

1.   Recovery from latent inhibition. Presentation of the US in the context
     of preexposure and conditioning results in renewed responding to the
     preexposed CS.

2.   Recovery from overshadowing. Extinction of the $CS_1$ results in increased
     responding to the overshadowed $CS_2$.

3.   Recovery from forward blocking. Extinction of the blocker $CS_1$ results
     in increased responding to the blocked $CS_2$.

4.   Recovery from backward blocking. Extinction of the blocker $CS_1$ results
     in increased responding to the blocked $CS_2$.

F.   *Extinction*

1.   Extinction. When CS–US pairings are followed by presentations of the
     CS alone, or by unpaired CS and US presentations, the CR decreases.

2.   External disinhibition. Presenting a novel stimulus immediately before
     a previously extinguished CS might produce renewed responding.

3.   Spontaneous recovery. Presentation of the CS after some time after the
     subject stopped responding might yield renewed responding.

4.   Renewal. Presentation of the CS in a novel context might yield renewed
     responding.

5.   Reinstatement. Presentation of the US in the context of extinction and
     testing might yield renewed responding.

6.   Reacquisition. CS–US presentations following extinction might result
     in faster or slower reacquisition.

7.   Partial reinforcement extinction effect (PREE). Extinction is slower fol-
     lowing partial than following continuous reinforcement.

G.   *Nonlinear combinations of multiple stimuli*

1.   Positive patterning. Reinforced $CS_1$–$CS_2$ presentations intermixed with
     nonreinforced $CS_1$ and $CS_2$ presentations result in stronger responding
     to $CS_1$–$CS_2$ than to the sum of the individual responses to $CS_1$ and $CS_2$.

6    Introduction

2.  Negative patterning. Nonreinforced $CS_1$–$CS_2$ presentations intermixed with reinforced $CS_1$ and $CS_2$ presentations result in weaker responding to $CS_1$–$CS_2$ than to the sum of the individual responses to $CS_1$ and $CS_2$.

*H.    Occasion setting*

1.  Simultaneous feature-positive discrimination. Reinforced simultaneous $CS_1$–$CS_2$ presentations, alternated with nonreinforced presentations of $CS_2$, result in stronger responding to $CS_1$–$CS_2$ than to $CS_2$ alone. In this case, $CS_1$ gains a strong excitatory association with the US.

2.  Serial feature-positive discrimination. Reinforced successive $CS_1$–$CS_2$ presentations, alternated with nonreinforced presentations of $CS_2$, result in stronger responding to $CS_1$–$CS_2$ than to $CS_2$ alone. In this case, $CS_1$ acts as an occasion setter.

3.  Simultaneous feature-negative discrimination. Nonreinforced simultaneous $CS_1$–$CS_2$ presentations, alternated with reinforced presentations of $CS_2$, result in weaker responding to $CS_1$–$CS_2$ than to $CS_2$ alone. In this case, $CS_1$ gains a strong inhibitory association with the US.

4.  Serial feature-negative discrimination. Nonreinforced successive $CS_1$–$CS_2$ presentations, alternated with reinforced presentations of $CS_2$, result in weaker responding to $CS_1$–$CS_2$ than to $CS_2$ alone. In this case, $CS_1$ acts as an occasion setter.

*I.    Temporal properties*

1.  Interstimulus interval (ISI) effects. Conditioning is negligible with short ISIs, increases dramatically at an optimal ISI, and gradually decreases with increasing ISIs.

2.  Intertrial interval (ITI) effects. Conditioning to the CS increases with longer ITIs.

3.  Timing of the CR. The CR peak tends to be located around the end of the ISI.

4.  Temporal specificity of blocking. Blocking is observed when the blocked CS is paired in the same temporal relationship with the US as the blocking CS.

5.  Temporal specificity of occasion setting. A serial feature-positive discrimination is best when the feature–target interval during testing matches the training interval.

*J.    Combination of multiple conditioning events*

1.  Sensory preconditioning. When $CS_1$–$CS_2$ pairings are followed by $CS_1$–US pairings, presentation of $CS_2$ generates a CR.

2.  Second-order conditioning. When $CS_1$–US pairings are followed by $CS_1$–$CS_2$ pairings, presentation of $CS_2$ generates a CR.

### Learning theories

Some classical conditioning theories stress the importance of mechanisms that act at the time of the presentation of the CS and the US. These theories assume that the association between events $CS_i$ and $CS_k$, $V_{CSi, CSk}$, represents the *prediction* that the $CS_i$ will be followed by $CS_k$ (Dickinson, 1980). Neural network, or connectionist theories frequently assume that the association between $CS_i$ and $CS_k$ is represented by the efficacy of the synapses, $V_{CSi, CSk}$, that connect a presynaptic neural population excited by $CS_i$ with a postsynaptic neural population that is excited by $CS_k$ (event $k$ might be another CS, or the US). When $CS_k$ is the US, this second population controls the generation of the conditioned response (CR).

Following Hebb's (1949) ideas, changes in synaptic strength, $V_{CSi, CSk}$, might be described by $\Delta V_{CSi, CSk} = f(CS_i)f(CS_k)$, where $f(CS_i)$ represents the presynaptic activity, and $f(CS_k)$ the postsynaptic activity. Different $f(CS_i)$ and $f(CS_k)$ functions have been proposed. Learning rules for $V_{CSi, CSk}$ either assume variations in the effectiveness of $CS_i$, $f(CS_i)$, the US, $f(CS_k)$ (Dickinson & Mackintosh, 1978), or both.

#### *Variations in the effectiveness of the CS during learning*

Attentional theories assume that the effectiveness of $CS_i$ to form $CS_i$–US associations (associability) depends on the magnitude of the "internal representation" of $CS_i$. In neural network terms, attention may be interpreted as the modulation of the CS representation that activates the presynaptic neuronal population involved in associative learning. Attentional theories include Makintosh's (1975), Grossberg's (1975) and Pearce and Hall's (1980) theories.

#### *Variations in the effectiveness of the US during learning*

A popular rule, proposed independently in psychological (Rescorla & Wagner, 1972) and neural network (Widrow & Hoff, 1960) domains, has been termed the "delta" rule. The delta rule describes changes in the synaptic connections between the two neural populations by way of minimizing the squared value of the difference between the output of the population controlling the CR generation, and the US. According to the "simple" delta rule, $CS_i$–US associations are changed until the difference between the US intensity and the "aggregate prediction" of the US computed upon all CSs present at a given moment, $(US - \Sigma_j V_{CSj,US}CS_j)$, is zero. The term $(US - \Sigma_j V_{CSj,US}CS_j)$ can be interpreted as the effectiveness of the US to become associated with the CS.

Schmajuk and DiCarlo (1992) introduced a model (the SD model) that, by employing a "generalized" delta rule (also known as backpropagation, see Rumelhart, Hinton & Williams, 1986) to train a layer of hidden units that *configure* simple CSs, is able to solve exclusive-or problems, and hence, negative patterning.

*Variations in the effectiveness of both the CS and the US during learning*

In order to account for a wider range of classical conditioning paradigms, some theories have combined variations in the effectiveness of both the CS and the US. For example, Frey and Sears (1978) proposed a model of classical conditioning that assumed variations in the effectiveness of both the CS and the US: $f(CS_i)$ is modulated by $V_{i,US}$. Wagner (1978) suggested that $CS_i$–US associations are determined by (a) $f(US) = (US - \Sigma_j V_{j,US} CS_j)$ as in the Rescorla–Wagner model; and (b) $f(CS_i) = (CS_i - V_{i,CX} CX)$, where CX represents the context, and $V_{i,CX}$ the strength of the CX–$CS_i$ association. Other theories that incorporate changes in the effectiveness of both the CS and the US include Wagner's (1981) sometimes opponent process (SOP) theory, Schmajuk, Lam and Gray's (1996) attentional–associative theory, Le Pelley's (2004) hybrid model, and Harris's (2006) elementary model.

## Performance theories

Some classical conditioning theories stress the importance of mechanisms that act during performance to control the generation of the CR. Examples of this approach are Miller's comparator hypothesis (e.g. Miller & Schachtman, 1985), Wagner's (1981) SOP model, Schmajuk, Lam and Gray's (1996) attentional–associative model, and Harris's (2006) elementary model.

*CS–US and CS–CS associations and decision processes during performance*

The comparator hypothesis (Miller & Schachtman, 1985; Miller & Matzel, 1988; Denniston *et al.*, 2001; Stout & Miller, 2007) suggests that the magnitude of the CR is determined by a comparator that uses the CS–CS and CS–US associations of the CSs present at a given time as inputs.

*CS–CS associations and inference generation during performance*

During classical conditioning, animals learn to expect (predict) that a CS is followed by another CS, or by the US. Tolman (1932) proposed that multiple expectancies (predictions) can be integrated into larger units, through a reasoning process called inference. One simple example of inference formation is sensory preconditioning (see Bower & Hilgard, 1981, page 330). As summarized above, sensory preconditioning consists of a first phase in which

two conditioned stimuli, $CS_1$ and $CS_2$, are paired together in the absence of the US. In a second phase, $CS_1$ is paired with the US. Finally, when $CS_2$ is presented alone, it generates a CR: the animal has inferred that $CS_2$ predicts the US. Tolman hypothesized that a large number of expectancies can be combined into a cognitive map (see Chapters 2, 7 and 16).

Dickinson (1980) suggested that knowledge can be represented in declarative or procedural form. Whereas in the declarative form knowledge is represented as a description of the relationships between events (knowing that), in the procedural form, knowledge is represented as the prescription of what should be done in a given situation (knowing how). Examples of declarative knowledge are classical CS–CS associations ($CS_1$ precedes $CS_2$) or CS–US ($CS_2$ precedes the US) associations. An example of procedural knowledge is the operant *S–R* association (if *S* is present, then do *R*). Dickinson indicated that declarative, but not procedural, knowledge can be integrated through inference rules.

By including CS–CS associations, some models of classical conditioning are able to generate inferences and, therefore, to describe sensory preconditioning. For instance, Gelperin, Hopfield and Tank (1985; Gelperin, 1986) proposed an autoassociative recurrent network capable of describing stimulus–stimulus associations during classical conditioning. The network can simulate first- and second-order conditioning, extinction, sensory preconditioning and blocking in the terrestrial slug, *Limax maximus*. Schmajuk (1987) proposed a dual memory architecture that incorporates an autoassociative nonrecurrent network capable of cognitive mapping in classical conditioning. The network separately computes CS–CS and CS–US predictions, and combines them to generate new expectancies. For instance, if CS(*A*) predicts (is associated to) CS(*B*), and CS(*B*) predicts the US, the network *infers* that CS(*A*) also predicts the US. Schmajuk defined first-order predictions as the prediction of the US by CS(*B*), and higher-order predictions as the predictions involving a chain of two or more predictions. The network describes complex classical conditioning paradigms such as sensory preconditioning, second-order conditioning, compound conditioning and serial-compound conditioning. Similarly, the models presented by Schmajuk and Moore (1988, 1989) and Schmajuk, Lam and Gray (1996) are able to describe sensory preconditioning and second-order conditioning.

### Computational models of classical conditioning

As suggested by Hinzman (1991), the inherent unreliability of verbal intuitive reasoning for relating hypotheses and experimental results favors theories that provide precise quantitative descriptions. Furthermore, only formal models can be simultaneously examined at different levels. At the behavioral

level, simulated behavioral results are compared with experimental data describing behavior. At the neuroanatomical level, interconnections among neural elements in the model are compared with neuroanatomical data, and the model performance is compared with animal performance after lesioning. At the computational level, simulated activity of the neural elements of the model is compared with the activity of single neuron or neural population activity. At the neurophysiological level, model performance is compared with animal performance after inducing long-term changes (e.g. lesions) or short-term changes (e.g. drug infusions) in different brain areas.

Below, we review in detail some computational models of classical conditioning that have been applied to a number of the conditioning paradigms described before. Some other models (e.g. Brandon, Vogel & Wagner, 2000; Kruschke, 2001; Pearce, 1994) are described later in the book, when they are relevant to the experimental results being discussed.

### The Rescorla–Wagner (1972) model

In their classic article, Rescorla and Wagner (1972) indicated that the impetus for their new theoretical model was not new data which clearly disconfirmed existing theories, but rather the accumulation of a pattern of data which appeared to invite a more integrated account. The salient pattern of data the authors referred to was a set of observations involving Pavlovian conditioning with compound CSs. The central notion of the theory was that organisms only learn when the actual value of the US differs from its expected value. By proposing the novel principle that this expected value of the US is computed as a linear combination of the associative strength of all active CSs, the effect of reinforcement or nonreinforcement on the associative strength of a CS depends upon the existing associative strength, not only of that CS, but also of other CSs concurrently present.

Rescorla and Wagner (1972) proposed to formalize the basic idea of their theory by modifying Hull's (1943) account of the growth of habit-strength (stimulus–response associations), as described by Bush and Mosteller's (1955) linear operator. In the Rescorla–Wagner (RW) model, variations in the strength of the CS–US association, $V_{i,US}$, are given by $\Delta V_{i,US} = \alpha_i \beta_{US}(\lambda_{US} - B_{US})$, where $\alpha_i$ represents the salience of $CS_i$, $\beta_{US}$ represents the learning rate parameter corresponding to the US, and $B_{US}$ is the linear combination of the prediction of the US by all active CSs. $B_{US}$ is given by $B_{US} = \Sigma_j V_{j,US}$. By this equation, CSs compete to gain association with the US. The conditioned response (CR) was assumed to be proportional to $B_{US}$.

As observed by Sutton and Barto (1981), a rule similar to the RW equation had been described in the neural network field by Widrow and Hoff (1960). This rule, termed the "delta" rule (Rumelhart, Hinton & Williams, 1986), describes

changes in the CS–US associations by way of minimizing the squared value of the difference between the predicted and observed values of the US. Most interestingly, as indicated by Duda and Hart (1973), the rule is able to solve simultaneous systems of equations; a power that provides a different perspective of the processes that take place during classical conditioning.

The RW model correctly described many Pavlovian conditioning phenomena such as acquisition and extinction of conditioned excitation, partial reinforcement, conditioned inhibition, overshadowing, blocking, unblocking by increasing the US strength, overprediction, generalization, US–preexposure effect and contingency effects. The success of the model in making specific correct predictions, inaugurated the modern era of experimental psychology.

In spite of its significant achievements, the RW model was unable to describe several aspects of classical conditioning including (a) the effects of temporal parameters, such as stimulus duration, interstimulus intervals (ISI) or intertrial intervals (ITI); (b) Pavlovian paradigms whose solution require a nonlinear combination of the prediction of the US by all active CSs, such as negative patterning; (c) conditioned inhibition not being extinguished by presentations of the inhibitory CS alone; (d) latent inhibition; (e) backward blocking; and (f) the recovery from blocking and overshadowing.

### The Van Hamme and Wasserman (1994) version of the Rescorla and Wagner (1972) model

Van Hamme and Wasserman (1994) offered a modified version of the Rescorla and Wagner (1972) model that is able to explain some of the results mentioned above. Van Hamme and Wasserman (1994) proposed that the association of a CS with the US decreases when the CS is absent ($\Delta V_{i,US} < 0$), instead of staying constant, as in the original model ($\Delta V_{i,US} = 0$ because $\alpha_i = 0$). The model can explain the effects of extinction of the companion CS overshadowing and blocking.

Dickinson and Burke (1996) observed that the Van Hamme and Wasserman (1994) rule did not specify when an absent CS was allowed to decrease its association with the US. Following previous suggestions (Chapman, 1991; Markman, 1989; Tassoni, 1995), they indicated that the expectation of an absent CS, via its (within-compound) association with a present CS, could serve that purpose.

The Van Hamme and Wasserman (1994) version of the Rescorla–Wagner (1972) rule is able to describe that (a) extinction of the blocking CS results in the recovery of the response to the blocked CS (Blaisdell *et al.*, 1999); (b) extinction of the overshadowing CS results in the recovery of the response to the overshadowed CS (Matzel *et al.*, 1985); but cannot explain (c) extinction of the context following latent inhibition (LI) results in the recovery of the response