# Rational models of conditioning

Nick Chater

*Division of Psychology and Language Sciences, University College London, London WC1E 6BT, United Kingdom.*
**n.chater@ucl.ac.uk**
**http://www.psychol.ucl.ac.uk/people/profiles/chater_nick.htm**

**Abstract:** Mitchell et al. argue that conditioning phenomena may be better explained by high-level, rational processes, rather than by non-cognitive associative mechanisms. This commentary argues that this viewpoint is compatible with neuroscientific data, may extend to nonhuman animals, and casts computational models of reinforcement learning in a new light.

Mitchell et al. provide an important critical challenge to the presuppositions underlying current theories of human conditioning, both in psychology and the neurosciences. They suggest that human contingency learning results from reasoning processes over propositional knowledge, rather than from an elementary process of forming associations. This commentary focuses on three questions raised by this analysis, and concludes with a perspective on the origin of contradictory forces in the control of behavior which does not invoke a clash between a cognitive and associative system.

**Multiple neural systems for decision making?** Mitchell et al. argue that behavioral evidence makes a case against a distinct associative learning system. Yet the idea that there are multiple, competing, neural systems underpinning decision making is very widespread within neuroscience. One line of evidence for multiple systems comes from double dissociations in human neuropsychology, and, perhaps most strikingly, from animal lesion studies (see, e.g., Coutureau & Killcross 2003; Killcross & Coutureau 2003). Yet such studies provide only tentative evidence for functionally distinct systems, rather than differential engagement of a single system (Chater 2003; Shallice 1988). Consider an analogy with allergies: Some people cannot eat prawns, but can eat pine nuts; other people can eat pine nuts, but not prawns. But we cannot, of course, conclude that there are two distinct digestive systems that process these different foods. Instead, a *single* digestive system deals almost uniformly with all foods, but exhibits two biochemical "quirks" leading to the selective allergies. Thus, a *single* processing system can in principle yield striking double dissociations of function (Chater, in press). Hence, double dissociations in humans, and animal lesion studies yielding double dissociations, are weak evidence for distinct processing systems. The same caveats apply to studies in which reinforcement learning is selectively impaired not by a lesion, but by a pharmacological intervention (e.g., a dopamine agonist, Pizzagalli et al. 2008). Similar issues arise, too, with neuroimaging studies. Such studies reveal differential neural activity under different task conditions. But such differential activity may nonetheless be entirely compatible with the existence of a single, unitary, decision-making system.

**Is animal conditioning associative?** Mitchell et al.'s account may be correct with regard to people. But perhaps rats really do use dedicated associative learning mechanisms. Indeed, this latter assumption is widespread in the comparative literature (e.g., Mackintosh 1983). Nonetheless, there are at least three reasons to doubt this. (1) Many aspects of animal cognition are highly sophisticated and seem to go far beyond the scope of purely associative mechanisms (e.g., Wasserman & Zentall 2006). (2) Associative theories of learning typically assume gradual modifications; yet actual behavior is roughly all-or-none (Gallistel et al. 2004), just as though the animal is adopting or rejecting a hypothesis about possible environmental contingencies. The familiar smooth learning curves arise only from data averaging. (3) Putative conditioning phenomena in animals appear to be highly sensitive to rational factors (Courville et al. 2006). So, for example, blocking (Kamin 1969) can be rationally understood in terms of "explaining away" (Pearl 1988); the slower rate of extinction from partially reinforced contingencies has a natural statistical explanation; and so on.

**The role of computational models of reinforcement learning.** There have been remarkable recent developments in computational models of reinforcement learning (Dayan & Abbott 2001) – often implicitly or explicitly viewed as capturing the computational principles of a distinct, striatal, non-cognitive, learning system (Jog et al. 1999). If Mitchell et al. are right, then such computational models should perhaps be interpreted differently: as providing an account of rational inferences that can be drawn from data concerning actions and rewards, given minimal background knowledge. But where background knowledge *is* available (e.g., about likely causal connections between actions, events, and rewards), we should expect that such knowledge will be incorporated appropriately (Gopnik & Schulz 2007). According to this perspective, computational models of reinforcement learning apply to a narrow class of situations, in which background causal knowledge is restricted, rather than describing the operation of a particular neural system that drives behavior.

**Clash of reasons, not clash of mechanisms.** One intuitive appeal of the idea of a split between associative and cognitive

systems, competing for the control of behavior, is a potential explanation for many paradoxical aspects of human behavior, both in laboratory studies of, for example, time-discounting and weakness-of-will and in real-world phenomena of addiction, depression, or phobias (Epstein 1994; McClure et al. 2004).

If, following Mitchell et al., we reject evidence for a distinct associative system, how are we to explain the origin of internal cognitive conflict? One straightforward approach (Chater, in press) is to propose that internal conflict arises from a "clash of reasons" rather than a clash of systems. In almost all nontrivial reasoning problems, different lines of argument appear to favour different conclusions. One source of reasons, among many, may be past experience (including the "reinforcement history"). Moreover, reasons are often not equally persuasive; nor are they equally easy to evaluate. When paying close attention and given sufficient time, it may become evident one reason is valid, whereas another reason is weak. But when attention is reduced, the weaker reason may nonetheless prevail. Therefore, to choose a classic example from probabilistic reasoning, the reasoner may decide that, given information about, say, Linda's intellectual and political background, it is more likely that Linda is a feminist bank teller, than that she is a feminist (Tversky & Kahneman's [1983] conjunction fallacy), because there is a better overall match with the former description (for which at least the first part matches), than the second description (which seems entirely incongruous). Considered reflection on probability may, or may not, lead the reasoner to draw the opposite conclusion.

More generally, it seems entirely possible that there will be systematic differences between responses when time and attention are limited and responses when time and attention are plentiful (see Cunningham & Zelazo [2007] for a similar perspective on apparent dissociations between two putative routes underpinning social cognition, as exemplified by, e.g., Bargh & Chartrand 1999); and concomitant differences in the degree to which brain areas are activated in the contemplation of different reasons. In summary, observing battles for control of the behavioral "steering wheel," and evidence for different behavioral and neural bases for the competitors, need not be interpreted as indicating a clash between distinct mechanisms (e.g., associative vs. cognitive), but might equally arise from a clash of reasons within a unified cognitive system.