

CONTENTS

1	Introduction	1
2	A Brief History of Rationality	16
3	The Rational Use of Cognitive Resources	36
4	(Ir)rationality Revisited	63
5	Strategy Selection, Metareasoning, and Learning to Be Rational	89
6	Strategy Discovery	111
7	Representations and Architectures	140
8	Improving Decisions	151
9	Conclusion	195
10	Appendix: Mathematical Details	205
	<i>Acknowledgments</i>	233
	<i>Bibliography</i>	235
	<i>Index</i>	259

1

Introduction

What does it mean to be rational? In the middle of the twentieth century, it looked like there was a definitive answer to this question. Mathematicians, economists, and statisticians all agreed that rational agents choose actions that maximize their **expected utility**, weighing the utility of each outcome by its probability (von Neumann and Morgenstern, 1944; Savage, 1954). But when psychologists looked carefully at how people choose actions, it became clear that this is an imperfect description of human behavior: people systematically violate the axioms of rationality and act in ways that are inconsistent with maximizing expected utility (Tversky and Kahneman, 1974; Kahneman et al., 1982).

One response to these findings is simply to decide that people are irrational—that the mathematicians, economists and statisticians had characterized rationality correctly and humans are faulty actors. However, there is another possible response, which is to ask whether we are holding humans to the right standard. Perhaps the fault lies not in human behavior, but in the definition of rational action. One clue that this might be the case is that, despite not following the prescriptions of classical rationality, human beings are the best example we have of a general-purpose intelligent system. If we want to understand rationality not in abstract, but in the context of what produces intelligent behavior, there might be lessons we can learn by starting from the choices that humans make.

The perspective of trying to understand how an intelligent agent should behave led to an alternative notion of rationality, one that came not from mathematics, economics, and statistics, but from computer science. One of the things that is missing from the principle of maximizing expected utility is any notion of how difficult it might be to do so. To make a classically rational decision, we would have to consider all the outcomes that might arise from each action we take and assign each outcome a utility and a probability. In any realistic setting, the number of possible outcomes will likely be very large, and assigning a precise utility and probability to each outcome could be very

difficult. In particular, it could be extremely **computationally costly** to run an algorithm that does this—the algorithm might take so long to execute that the original decision becomes irrelevant.

Computer scientists are used to trading off the quality of an output with the cost of computing it. Formalizing this trade-off is crucial for designing artificial intelligence (AI) systems that need to take informed actions in the real world. This is especially true for designing autonomous robots that interact with unstructured real-world environments in real time and compute all of their decisions on onboard hardware, such as autonomous vehicles. By considering how such AI systems should act, computer scientists developed the theory of **bounded optimality**—an alternative definition of rationality that focuses not on whether an agent selects the best *action*, but on whether the agent follows the best *algorithm* for making that decision taking into account both utility and computational cost (Horvitz, 1987; Russell and Wefald, 1991a; Russell and Subramanian, 1995; Russell, 1997). The bounded-optimal agent settles for a good enough action when time is limited. This poses a new kind of decision problem: deciding what internal computations to execute in order to select an external action (Russell and Wefald, 1991b).

We (the authors, Falk, Fred, and Tom) are interested in considering alternative ways of defining rationality not because we want to know whether or not humans are rational, but because knowing how people *should* be solving a problem is often a key part of understanding how they *do* solve that problem and how they could solve it even better. In cognitive science, **rational analysis** is a method for understanding human behavior in which we consider the optimal solution to the problems that human minds face (Anderson, 1990). Traditionally, this kind of analysis was done within the framework of classical rationality. However, this omits an important part of actually understanding human behavior: understanding the constraints under which human minds operate. To address this omission, we and our colleagues developed the approach of **resource-rational analysis**, which uses bounded optimality and associated ideas to ask how human minds *should make the best use of their limited cognitive resources* to solve the problems they face (Lieder et al., 2012; Griffiths et al., 2015; Lieder and Griffiths, 2020b).

In this book, we present a formal framework for applying resource-rational analysis to understand and improve human behavior, a set of tools we have developed to make this easier, and examples of how we have used this approach to revisit classic questions about human cognition, pose new ones, and enhance human rationality. Our target audience is not limited to the people who have traditionally been concerned about rationality: psychologists studying judgment and decision-making, economists, and philosophers. We

also address the broader group of psychologists, cognitive scientists, and neuroscientists trying to understand human minds and brains, as well as computer scientists interested in reproducing such systems in machines. We see these ideas as being relevant to this broader audience because understanding the rational use of cognitive resources offers not just a new way to think about the evidence for human irrationality in decision-making, but also a new perspective on many long-standing questions in cognitive science. Many theories in the social sciences and humanities are (implicitly or explicitly) built on assumptions about human rationality. This book provides constructive suggestions for how those assumptions can be improved. Moreover, the new approach to improving people's decisions that we present is relevant for anyone who wants certain decisions to be made well, including researchers, educators, and practitioners who care about improving people's lives or the functioning of groups, organizations, or society.

To understand the potential we see in the resource-rational approach, it is worth observing that the fundamental limitation of the classical notion of rationality is not just that it does not engage with how difficult it might be to make a decision, but that it does not engage with cognitive processes at all. Maximizing expected utility is fundamentally a **behaviorist** (Skinner, 1953) notion of rationality: it defines rational action purely in terms of the agent's response to its environment. By contrast, considering the question of how an agent should best deploy its cognitive resources to solve a problem places resource-rational analysis in firm contact with those cognitive processes. For the first time, we can say not just what action an agent should take but also how it should think about what to do. This is a powerful idea that extends across the psychologist's entire toolbox of cognitive processes—attention, planning, and memory all involve the use of limited resources, and resource-rational analysis gives us a new way to understand them.

Even more significantly, resource-rational analysis provides a new way to think about how to support and improve human cognition. When we use classical rationality to assess people's performance, we risk assuming that performance can be improved by teaching people the classically rational solutions to the problems they face. However, because these solutions are computationally costly, they may be impractical for many real-world decisions. Resource-rational analysis suggests an alternative approach: teaching people practical cognitive strategies that make better use of their limited cognitive resources. Thinking about computational costs also highlights another way that we can support human decision-making: building more computation into human environments. Using this approach, we can design **cognitive prostheses** to support longer-term decision-making (Lieder et al., 2019) and understand the

circumstances under which behavioral **nudges** are most likely to be effective (Callaway et al., 2023).

In addition to a framework, tools, and examples, we hope to share with you our excitement about this approach and our vision for the many ways it can be used to understand and improve human minds. The cognitive revolution provided an alternative to behaviorism by showing how mathematics can be used to express rigorous hypotheses about cognitive processes (Miller, 2003). In the same way, we see resource-rational analysis as offering a path to a cognitive revolution in how we use and think about rationality when developing theories to explain, predict, and improve human behavior. The result is a blueprint for a new bridge between ideal solutions and cognitive processes that is directly informed by empirical data and supports interventions that meaningfully improve human lives. We hope this book will help you join us in building that bridge.

1.1 A simple example

Before we dive into any more details, we will give a simple example that illustrates what a resource-rational analysis looks like, how it differs from classical notions of rationality, and how this difference gives us new tools for understanding human behavior. This example takes the simplest kind of decision an agent could make, considers how it becomes more complicated when resources are limited, and shows a surprising connection to something people do that has long been considered irrational. It is also the first resource-rational analysis that any of us worked on (Vul et al., 2009).

Imagine you have a choice between two options. For simplicity, we will call them Option A and Option B. If you select Option A, you receive a reward with probability p . If you select Option B, you receive the same reward with a probability of $1 - p$. Which option should you choose?

From the perspective of classical rationality, this problem is trivial. If p is greater than 0.5, then you should choose Option A, otherwise choose Option B. But it becomes much more challenging if you don't know the value of p . In many real human decisions, assessing the probability of different outcomes is the hardest part. Doing so may require imagining future possibilities, considering different sources of evidence, recalling similar situations from memory, or engaging in various other cognitive processes.

To simplify things, we are going to model all of those different cognitive processes as methods of **sampling** a possible outcome. Whether we imagine, consider, or recall, we are getting a piece of information about the value of p . We will assume that obtaining this piece of information is equivalent to

sampling a possible outcome for this decision—getting a result that favors Option A with probability p or Option B with probability $1 - p$.

Assuming these cognitive processes are costly—they take time and mental effort—making rational use of our cognitive resources in this setting is a matter of deciding how many samples to generate before choosing one of the options. Making this decision requires considering the cost of sampling relative to the cost of making an error. Finding the optimal trade-off between these two costs can be formulated as an optimization problem.

The solution to this optimization problem might be surprising: for a wide range of relative costs, it is optimal to take a single sample. Intuitively, there are few values of p where more samples are informative. If p is close to 1 or 0, then a single sample indicates the right answer with high probability. If p is close to 0.5, then more samples yield little benefit over choosing at random. If the cost of making an error is large relative to the cost of sampling, then the few cases where p falls in an intermediate range may be enough for it to make sense to draw more samples. But in many situations, one is enough.

This resource-rational analysis explains an aspect of human behavior that, on its surface, appears irrational. Given choices like this, where options are rewarded probabilistically, people tend to select those options with the probability that matches the probability of reward. This **probability matching** behavior is inconsistent with the classically rational solution, which is to always choose the option that is best according to all available information (Vulkan, 2000). However, probability matching makes sense if people are using a cognitive mechanism that behaves like sampling and have to pay a cost for each sample they generate.

While this is a simple example—perhaps the simplest we can imagine—it illustrates some themes that will reappear many times in this book. First, while the classically rational solution is easy to define abstractly, computing a concrete rational action is much more cumbersome and often impractical due to limited cognitive resources. Second, by taking people's limited cognitive resources into account, we can understand aspects of human behavior that differ systematically from the prescriptions of classical rationality. Finally, we can reach that understanding while abstracting away from the specific cognitive mechanisms involved. Throughout the book, we are going to define models at a fairly high level of generality, using processes such as sampling that are intended to stand in as abstract proxies for a variety of different cognitive mechanisms.

While we have not specifically explored ways of improving human decision-making in the simple kinds of decisions we are considering in this example, our resource-rational analysis of such decisions suggests some ways

we could improve them. If we can identify situations where the cost of error is high relative to the cost of sampling, these would be good situations in which to signal that people might want to think carefully before responding. Likewise, if we have information about what range the value of p is likely to fall into, we might be able to tell people whether generating more than one sample is worthwhile. This is a way to support people's decisions by putting some of the computation they might otherwise have to perform into their decision environment.

In the rest of the book, we explore how making the move that we made here—going from a task to a decision about how to deploy our cognitive resources to best solve that task—can be used to shed light not just on decision-making, but also on planning, attention, and memory. The optimization problems involved in allocating these resources become increasingly complex, and we will present a framework for formulating and solving those problems. We will also show how the resulting models connect to traditional questions in cognitive psychology, and how they can be used to design effective interventions for agents with bounded cognitive resources.

1.2 Answers to skeptical questions

The work presented in this book is the culmination of a decade-long research program exploring applications of resource-rational analysis in cognitive science. Consequently, we have had the opportunity to present these ideas to various audiences and, hence, to field various skeptical questions from those audiences. To help you decide whether reading further is worth the cognitive resources involved, this section provides our answers to those skeptical questions. Of course, if you are already convinced, you can skip this section.

1.2.1 *Why should I care?*

The answer to this question depends on the person asking it. Depending on your interests, you may find different aspects of resource-rational analysis appealing.

If you are a **psychologist**, resource-rational analysis provides a different way to think about cognitive processes. Many of the phenomena that interest psychologists, such as problem-solving and deductive reasoning, are trivial from the perspective of classical rationality: the goal is simply to solve the problem. In these cases, the journey—the sequence of cognitive processes we engage in when solving the problem—is more important than the destination. For a long time, psychologists have developed theories about how people perform such tasks by manually assembling different cognitive

mechanisms, strategies, or heuristics guided by empirical results. If you use resource-rational analysis, you can have computers do this tedious work for you: having defined a task and a set of elementary cognitive processes that could be involved, the optimal way to combine these processes can be derived automatically. This provides a way to answer questions about how people should think and what they should be thinking about, in the same way that classical rational analysis answers the question of how people should act. Consequently, resource-rational analysis makes deeper contact with the traditional preoccupations of cognitive psychology than rational analysis. It also provides a different way to design interventions to support human decision-making, focusing on guiding people toward better strategies and offloading some of the computation involved into the environment in which decisions are made.

If you are an **economist**, you are likely to be familiar with the idea of modeling human behavior as the solution to an optimization problem, but possibly also concerned that people systematically deviate from classical rationality. Resource-rational analysis is a way to create models of human behavior that are just as general—and use many of the same tools—as classical rational models, but also capture some systematic deviations from classical economic rationality. As we will discuss in the next chapter, this approach is aligned with recent work in economics, particularly in the **rational inattention** literature (Sims, 2003; Caplin and Dean, 2015), but also provides a way to more explicitly connect those ideas to cognitive processes, and thereby make the models' predictions even more accurate.

If you are a **neuroscientist**, having a theory of rational action that makes closer contact with cognitive mechanisms means that theory also makes closer contact with the brain. The successes of classical rational analysis led to a flurry of research on how the brain might implement probabilistic inference, typically by searching for one kind of computation or neural circuit that might solve this problem. The resource-rational perspective suggests that different instances of probabilistic inference might be solved in different ways, depending on the structure of the problem and the elementary cognitive processes available. For example, we will discuss a variety of sampling and planning algorithms that make sense to use under different circumstances. This provides a more nuanced set of targets for neuroscientists to search for in the brain. In addition, the core problem of selecting which cognitive processes to engage in has a structure that parallels that of **reinforcement learning**, creating the possibility that the extensive work on the neuroscience of reinforcement learning (Dolan and Dayan, 2013; O'Doherty et al., 2015; Lee et al., 2012; Mattar and Lengyel, 2022; Miller and Venditto, 2021) might also help us understand this metacognitive process. Moreover, resource-rational

analysis is a promising vehicle for translating neuroscientific insights into the basic computations neurons and neural circuits perform and how costly those computations are into computational models of cognition and their resulting applications.

If you are a **computer scientist**, resource-rational analysis provides an opportunity to understand how human minds are so efficient in their use of limited cognitive resources to solve a wide range of problems and to define a formal framework for emulating this in machines. Bounded optimality is a notion that originated in the AI literature, and there are still a lot of opportunities to improve on the computational methods that are used to solve the problem of discovering efficient algorithms. Research in machine learning has explored a variety of ad hoc methods for deciding how much computation to perform (Graves, 2016; Banino et al., 2021). Formulating these problems more explicitly in terms of **rational metareasoning** (Russell and Wefald, 1991a) may be a way to improve on these methods. Finally, a formal framework that captures some of the ways that people systematically deviate from classical rationality is a valuable tool for making sense of human behavior. If you are interested in building systems that interact with people, then better understanding the mapping from people's preferences to their behavior is essential for being able to infer people's preferences from their behavior (Ho and Griffiths, 2022).

If you are a **curious person** who is primarily interested in how to make better decisions in your own life, this book offers a different perspective on what it actually means to be rational. Many books about decision-making emphasize human irrationality, focusing on deviations from classical rationality. Our perspective is more nuanced: when considering the resource constraints that human minds operate under, many of the seemingly irrational things people do make more sense. It is reasonable to want to be more rational in your own decisions, but the path to getting there has to respect those constraints. We discuss improving human decisions in detail in Chapter 8, which might be of particular interest to you. You can find more pragmatic advice in Tom's book *Algorithms to Live By* (Christian and Griffiths, 2016), which is not explicitly framed in terms of resource-rational analysis but shares with this book the idea that computer science has an important contribution to make to understanding rational action.

1.2.2 *Isn't this just bounded rationality?*

Immediately after everyone started to agree on how to define rationality, Herbert Simon pointed out that it was an unreasonable standard for human behavior (Simon, 1955). His concern was one we share: that humans have limited time and cognitive resources, making maximization of expected utility

impractical. Simon wanted researchers to acknowledge the idea of **bounded rationality**—that the constraints under which human minds operate make classical rationality unachievable. This idea is very much a part of resource-rational analysis.

However, Simon viewed bounded rationality as an idea, rather than a formal framework. In a letter to the psychologist Gerd Gigerenzer, he wrote, “I have never thought of either bounded rationality or satisficing as precisely defined technical terms, but rather as signals to economists that they needed to pay attention to reality, and a suggestion of some ways in which they might. But I do agree that I have used bounded rationality as the generic term, to refer to all the limits that make a human being’s problem spaces something quite different from the corresponding task environments: knowledge limits, computational limits, compatibility of component goals” (Augier and March, 2004, p. 406). This is where we differ: we view resource-rational analysis as a formal framework that can be useful for understanding the behavior of bounded agents, turning the idea of bounded rationality into a generalizable approach that can be used to derive precise predictions about cognition and behavior given a specification of the problem to be solved and the cognitive processes available for doing so.

In the same letter, Simon objected to the interpretation of bounded rationality in terms of constrained optimization—assuming that people solve an optimization problem similar to the one solved by a classically rational approach but with bounded resources providing constraints on the solution. In fact, he threatened to sue the next person who interpreted bounded rationality in this way. To anticipate any such litigation, we want to state that this is not how we think about resource-rational analysis. Rather, we see this approach as fundamentally changing the optimization problem implicit in the definition of rationality from one that is focused on external actions to one that is focused on internal computations. The simple example given above illustrates this: thinking about cognitive processes in terms of sampling leads to a completely new optimization problem, one that is focused on using those cognitive processes effectively.

It is notable that the one formal model that Simon published in his papers on bounded rationality has a very similar flavor to the models we consider in this book. Simon’s model of **satisficing** (Simon, 1955) turns a choice problem into a stopping problem, where options are considered until one is found that is above a derived threshold. This kind of structure arises repeatedly in resource-rational models, where limited cognitive resources require us to consider our options (or gather information, or sample from distributions) sequentially. Many of the models we present later in the book focus on the problem of making intelligent decisions about when to stop thinking in a way that is very similar to satisficing.

1.2.3 *How is this different from other perspectives on rationality?*

Understanding what it means to be rational, and how the definition of rationality relates to what people actually do, has been a preoccupation of many researchers from many different disciplines. We provide a detailed treatment of these ideas in Chapter 2, but want to briefly sketch how we see the relationship of resource-rational analysis to previous approaches.

The **heuristics and biases** research program carried out by Daniel Kahneman and Amos Tversky in the 1970s provided clear evidence that people's judgments and decisions deviate from classical rationality (Tversky and Kahneman, 1974; Kahneman et al., 1982). This research program was originally framed in terms of the computational costs of rationality: given that probabilistic inference and calculating expected utility are intractable, people must find some way to approximate them. Heuristics are the natural solution, and biases are the clues that we can use to identify those heuristics. We view resource-rational analysis as being very aligned with this original goal: considering how the costs of using cognitive resources trade off with performance gives us a way to characterize what makes a heuristic worth using. As a consequence, we can rederive some of the heuristics that Kahneman and Tversky discovered (see Chapter 4) and derive new heuristics directly from the formulation of the problem people have to solve (see Chapter 6).

The notion of **ecological rationality** was offered by Gerd Gigerenzer and colleagues as an alternative way of understanding heuristics (Gigerenzer and Todd, 1999a; Gigerenzer, 2000). Focusing on the interaction between cognition and the environment, ecological rationality explores the idea that particular heuristics might be adaptive solutions in particular environments. In this way, using heuristics should not be considered a deviation from rationality, but rather a reasonable strategy in itself. More recently, this perspective has been argued for from the perspective of the **bias-variance trade-off**, which shows that simple strategies can be best when information is limited (Gigerenzer and Brighton, 2009). Again, we see this approach as being well aligned with resource-rational analysis, and some of the heuristics identified by this research program fall out of our framework (see Chapter 6). However, one important difference is that we consider computational costs as an important additional factor that pushes people toward simpler heuristics, beyond the effect of the bias-variance trade-off.

Work on **rational inattention** in economics has explored how classical rationality can be brought into closer alignment with human behavior by assuming that gathering information is costly (Sims, 2003; Caplin and Dean, 2015). The resulting models consider how much investment an agent

should make in gathering information—allocating attention to it—in order to find the best compromise between the costs it incurs and the benefits it provides for decision-making. This approach shares with the models we will discuss in this book the idea of introducing resource constraints into the formulation of the decisions that humans face. One difference is that the models we present here include cognitive processes—computations—as well as information gathering. However, we see the class of models developed under the flag of rational inattention as instances of resource-rational analysis.

Other researchers have used the term **computational rationality** to refer to applications of bounded optimality to cognitive science (Lewis et al., 2014; Gershman et al., 2015). We prefer to refer to these as resource-rational models (and have used this term since 2012; Lieder et al., 2012) as a way to emphasize the fact that they are about the rational use of cognitive resources and to avoid making a commitment to those resources being computational (something that allows this approach to make direct contact with rational inattention, for example).

1.2.4 Why are you proposing this now?

While the mathematical ideas behind bounded optimality were first proposed in the 1980s, it has been hard to use this approach to identify concrete cognitive strategies until the last decade or so. It has now become possible to actually define meaningful cognitive models based on bounded optimality because of innovations in the way we formulate and solve the underlying computational problems, and because of significant increases in the speed of the computers on which we solve those problems. Using the resulting models as the basis for designing effective behavioral interventions is potentially even more computationally challenging, which is one reason why it is a topic we have only been able to explore recently. One of the important things we do in this book is to provide you with the details on how we draw on ideas in artificial intelligence and statistics to set these problems up in a way that makes them (at least approximately) solvable.

1.2.5 So aren't you also creating a harder problem for human minds to solve?

Yes, we are defining a harder problem. Deciding how to think is certainly harder than deciding how to act. However, that harder problem does not necessarily need to be solved by human minds, and definitely does not need to be solved on the same timescale.

The problem of optimizing how we use our cognitive resources could be solved in different ways. One way would be simply considering all of the ways we could use our resources and choosing the best strategy. This is extremely computationally challenging—it involves presenting the harder version of the problem directly to human minds. But, importantly, there are other ways to solve this problem. One is to make use of the fact that we have repeated opportunities to perform particular tasks, and even new tasks often share components with tasks we have performed before. This means that we can *learn* better cognitive strategies over our lifetimes. Another is to recognize that variants of those same tasks were also faced by our ancestors, meaning that we can *evolve* better cognitive strategies across generations. In some cases, we do this via cultural evolution, coming up with intuitive strategies or algorithms we share with other people (e.g., Thompson et al., 2022). In other cases, biological evolution can act as the optimizing force—something that seems particularly plausible for making effective use of resources such as attention and memory that are shared with many other species (e.g., Tomlin et al., 2015).

In this way, although deciding how to think is harder than deciding how to act, humanity has had much more time to solve the problem. Individual human decisions are typically made on a timescale ranging from a few hundred milliseconds to a few weeks. Cognitive strategies are learned over a lifetime or tuned over many generations. And the way we use particular cognitive resources may pre-date all human minds.

We also expect that the mind solves the problem of deciding how to allocate its cognitive resources imperfectly. Resource-rational analysis is useful to the extent that people’s cognitive strategies approximate the optimal solutions to these problems; it does not require that they be perfect. There is room for slippage, suboptimality, and satisficing in reasoning, learning, and evolution.

1.2.6 *Doesn’t this lead to an infinite regress?*

Yes, we have introduced another level of abstraction into decision-making, and the same move can potentially be iterated. Thus, while our focus here is on deciding how to decide, we can imagine deciding how to decide how to decide, and so on. However, we do not see this kind of regress as a major problem. The key reason is that each level of abstraction has significantly diminishing returns. Deciding how to decide helps us explain a wide range of otherwise puzzling aspects of human behavior. Adding another level of abstraction—deciding how to decide how to decide—moves the focus from cognitive processes to cognitive architectures, which we explore in Chapter 7 (e.g., Milli et al., 2021). At this level, the question is how we should design a mind for it to be most effectively used by an agent that decides how to best

allocate cognitive resources. While we can still make progress in exploring this question, the answers are fairly abstract and harder to relate directly to human behavior. Adding further levels of abstraction gets us even further from behavior and begins to engage us with questions that start to push us beyond the limits of cognitive science (e.g., how should you design an evolutionary process to select a cognitive architecture that allows an agent that optimally allocates cognitive resources to solve problems?). While these questions are theoretically interesting, we see the biggest payoffs of these ideas lying at the lower levels of the potentially infinite regress.¹

1.2.7 *Are you really not interested in whether people are actually rational?*

Yes, really. Or, to be more precise, we are not interested in testing the hypothesis that people are actually rational under different criteria for rationality. To us, what matters more is whether a particular definition of rationality is useful for understanding human behavior. We use resource rationality as a **methodological assumption** that allows us to make progress in the broader project of making sense of people's thoughts and actions (Godfrey-Smith, 2001). Classical rationality is arguably useful in this respect, even if people deviate from it in systematic ways. Most Bayesian models of cognition are implicitly formulated within this framework. These models help us understand the aspects of people's behavior that we can explain in terms of the optimal solutions to the problems that human minds face (Anderson, 1990; Griffiths et al., 2010). Resource-rational analysis allows us to go a step beyond that, introducing assumptions about constrained cognitive resources and using them to explain some of the systematic deviations from classical rationality.

Given the vast space of possible Bayesian models and cognitive mechanisms, we anticipate (and our critics have suggested; Jones and Love, 2011; Bowers and Davis, 2012) that it will be hard to falsify the assumption of either of these forms of rationality. For that reason, we do not think it is productive to engage in debates about whether or not people are rational. Rationality, of whatever denomination, provides a *framework* in which we can develop *models* of human behavior (Lakatos, 1970). Individual models are falsifiable—they generate specific predictions that we can test in further experiments.

1. It is worth noting that there are also precedents of not being too concerned about infinite regress in the context of understanding decision-making. The classic example is in game theory, which considers the asymptotic equilibria of processes where players deliberate about the consequences of their actions. Whether similar asymptotic results can be obtained in the case of resource-rational analysis is an interesting research question.

Establishing whether people's behavior is consistent with specific models allows us to evaluate whether the factors considered in those models—hypothesis spaces and prior distributions in Bayesian models, augmented by cognitive mechanisms and constraints in resource-rational models—are candidates for explaining different aspects of that behavior.

Ultimately, we take the success of specific models based on resource-rational analysis as evidence that the framework is *useful*, not as evidence that the assumption that people always act in a resource-rational way is *true*. The approach we present in this book will be successful if it helps us develop theories that help us predict and explain why people act in the way they do and design interventions that are effective in improving that behavior.

1.3 The rest of the book

If you have gotten this far, you are hopefully satisfied with our answers to your skeptical questions—or are at least reserving judgment until you get more details. In the same spirit of letting you make an informed decision about how to use your time, we will use the remainder of this chapter to briefly summarize the rest of the book. Our focus in the book is on presenting the theoretical ideas behind resource-rational analysis and illustrating these ideas through some of the applications and experimental results that demonstrate the value of this approach. However, this is still a new perspective and we anticipate that some of the strongest evidence for its utility is still to come. We aim to equip readers to be able to gather that evidence for themselves.

Chapter 2 provides a brief history of different perspectives on rationality in psychology, economics, and computer science, expanding on the description that appears above and providing more details on the background behind resource-rational analysis and the context in which it arose.

Chapter 3 introduces the key ideas behind the formal framework of resource-rational analysis, drawing on bounded optimality and rational metareasoning to define a set of mathematical tools that we will use to define problems and state their solutions in the rest of the book.

Chapter 4 revisits the evidence for human irrationality, providing an analysis of some of the classic heuristics proposed by Kahneman and Tversky from the perspective of resource-rational analysis. The heuristics include **anchoring-and-adjustment** (Tversky and Kahneman, 1974) and the **availability heuristic** (Tversky and Kahneman, 1973), both of which are related to ideas about how to best use sampling algorithms to estimate probability distributions.

Chapter 5 applies resource-rational analysis to strategy selection. According to explanations of human decision-making postulating heuristics, the

mind is equipped with a toolbox full of different decision strategies. This raises the question of how people know when to use which strategy. This chapter explains how the problem of selecting a strategy can be expressed as a form of rational metareasoning—a problem that has been studied in the artificial intelligence literature (Russell and Wefald, 1991a). We then present a theory according to which people solve this problem by learning to predict how well different strategies work in different situations. We conclude the chapter with some empirical findings demonstrating that people can learn to become more rational.

Chapter 6 turns from strategy selection to strategy discovery. While strategy selection is about choosing between known strategies, strategy discovery is much harder: it requires us to put together new strategies out of elementary cognitive operations. In this chapter, we show how this problem can be solved using **metalevel Markov decision processes** (Hay et al., 2012), with applications to decision-making, planning, and memory.

Chapter 7 considers the implications of resource-rational analysis for two other key questions in cognitive science: how should we represent the world and how should human minds be structured? In answering the first question, we show how thinking about representations in terms of the computational costs they incur can account for how people form abstractions (Ho et al., 2022; Correa et al., 2023). In answering the second question, we show how resource rationality favors minds composed of a small number of cognitive systems (justifying the idea of having a “fast” and a “slow” system) (Milli et al., 2021).

Chapter 8 considers how resource-rational analysis can be used as a guide to constructing effective interventions for improving people’s decisions. In this chapter, we show how having a more accurate normative standard for people’s behavior can allow us to better identify real opportunities for improvement and to derive effective interventions to achieve those improvements. In particular, we discuss the approaches we have taken to teach people resource-rational decision strategies (Callaway et al., 2022a), construct cognitive prostheses (Lieder et al., 2019), and nudge people toward making better decisions with less effort (Callaway et al., 2023).

Chapter 9 concludes the book, summarizing the results from the previous chapters and considering their implications, as well as highlighting future directions for this research program.

We have constructed the book with the intent that it be read from start to finish, with later chapters building on ideas presented in earlier chapters and assuming some of their content as context. Hence, perhaps appropriately, the main decision you have to make about how to best use your cognitive resources is when to stop reading and start acting on the things you have learned.

INDEX

Page numbers in italics indicate figures and tables.

- ACT-R, cognitive architecture, 111, 137, 168
- Adaptive Character of Thought, The* (Anderson), 27, 29
- AI. *See* artificial intelligence (AI)
- algorithm(s), 2, 8, 12, 48, 68, 140; AI-Interpret, 162, 162–63, 163; importance sampling, 75–76; learning, 26, 206n1; Metropolis-Hastings algorithm, 67; planning, 142–44; sampling, 31, 65–66; SARSA, 95
- algorithmic level, rational process models, 30
- Allais paradox, 22, 24, 25, 64, 83
- anchoring-and-adjustment, 71, 72; heuristic, 14, 64; process, 24; resource-rational, 68, 69, 70
- anchoring bias, 71, 72, 73; experiments, 73, 74; resource-rational perspective on, 64–74
- Anderson, John: *The Adaptive Character of Thought*, 27, 29; rational analysis, 60
- architecture(s), 140; definition, 61; resource-rational cognitive architectures, 146–48
- artificial intelligence (AI), 2, 19, 67, 69; implications of resource-rational analysis for AI, 200–201
- attention, 31–32, 75–83, 125–28
- availability biases, resource-rational perspective, 74–83
- availability heuristic, 14, 24
- axiom(s) of rational choice, 17, 18, 84
- backward induction, 224–26; recursive implementation, 225
- Bayesian approach to cognitive science, 215
- Bayesian belief updating, 127, 216–18, 221
- Bayesian decision theory, 65, 66
- Bayesian inference, 28, 30–31, 122, 126, 135, 214
- Bayesian learning, 94
- Bayesian metalevel MDP(s), 215–19; belief states, 215–16; belief updating, 216; partially observable MDPs (POMDP), 218–19
- Bayesian metalevel policy search (BMPS), 228–31; learning to select computations, 230–31; value of information (VOI), 229, 230
- Bayesian model(s) of cognition, 13, 14, 30, 94, 215
- Bayes' rule, 28, 29
- behavioral economics, 172, 197; rise of, 203–4
- behavioral law, 198
- behaviorism, 3, 4, 42
- Bellman equations, 224
- Bernoulli, Daniel: expected value, 16; subjective value, 16
- Bernoulli, Nicolas: expected value, 16; subjective probability, 17
- Bernoulli distribution, 217; metalevel probability model, 217–18
- bias(es), 23, 151; rationality, 37. *See also* heuristics and biases
- bias-variance trade-off, 10, 26

- blinkered approximation, 227–28
- BMPS. *See* Bayesian metalevel policy search (BMPS)
- boost(s), 157, 159, 167, 169, 185, 186, 189, 194
- boosting: decision-making, 156, 157; decision-making with resource-rational heuristics, 159–69; future directions, 167–69; optimal, 160–61; theory of resource-rational, 167
- bounded optimality, 2, 11, 14; accounting for cognitive constraints, 43–46; cognitive system, 39; interpolating between value of cognition (VOC) and, 53–57; resource-rational analysis and, 32–33
- bounded rationality, 19–21, 33–34; Simon on, 8–9
- brain points, 181
- ChatGPT, 170
- choice architecture, 152–59, 172–85, optimal nudging as framework for choice architecture, 178–80
- classification, 168, 200
- cognition: Bayesian models of, 30; benefit of, 47; cost of, 49–53; human, 2, 3, 19, 22, 27, 29, 30, 33, 36, 59, 63, 88, 156, 185, 233; as resource, 59; social, 170; value of, 39
- cognitive architecture(s), 43, 54, 61; ACT-R, 61, 111, 137, 168; metalevel MDP, 115; SOAR, 61, 111, 137
- cognitive augmentation, 188; approach, 169; cognitive prostheses, 170; decision-making, 155, 157, 169–72
- cognitive control, plasticity of, 108–9
- cognitive costs, metalevel MDP, 129, 209
- cognitive mechanism, 5; rationality, 13–14
- cognitive operation(s), 48; metalevel MDP, 115, 116, 122, 127, 130, 135
- cognitive policies, 44
- cognitive process(es), 5, 43, 55; decision-making strategies, 111–12; general framework for characterizing optimal, 136–37; optimal sequential, 213
- cognitive prostheses, 3–4, 170, 171
- cognitive psychology, 6, 141, 149
- cognitive resources, 2, 3, 5, 15, 50, 63; bounded, 6
- cognitive science, 11, 13
- cognitive strategies: discovering, 201–2; learning better, 12
- cognitive systems, resource rational cognitive systems, 56
- cognitive training, decision-making, 164–67
- cognitive tutor(s), decision-making, 164–67
- collective decision-making, 203
- compensatory, 106
- computation(s), 48, 200; metalevel MDP, 209, 210–11; views of, 214–15
- computational architecture C, 43, 69
- computational cost. *See* cognitive costs
- computational level, 30, 61
- computational rationality, term, 11, 33
- computer science, 8, 14, 140
- conjunction fallacy, representativeness heuristic, 23
- constrained optimization, bounded rationality, 9
- construal(s): definition, 146; efficient planning, 144–46; navigation problem, 144–46, 145
- continuity, axiom, 18
- cost: metareasoning, 148, 149; opportunity, 50
- cost-benefit trade-offs, 46
- cost function, metalevel MDP, 115, 116, 122, 127, 130
- cost of cognition, 49–53, 58
- Cox's theorem, 28
- credit assignment problem, 228
- debiasing, decision-making, 155
- decision-making, 2, 151–52; approaches to improving, 152–56; beyond cognitive

- strategies, 193–94; beyond individual, 190–93; beyond rational self-interest, 186–87, 189–90; boosting, 156; boosting, with resource-rational heuristics, 159–69; cognitive augmentation, 155, 169–72; cognitive tutors and training, 164–67; collective, 203; debiasing, 155; evidence for human irrationality, 3; formalizing rational, 17–19; future work on resource-rational approaches, 188; game theory, 13n1; generating decision aids, 162–63; heuristics, 14–15; human, 5–6; human behavior, 12–13; human rationality and, 109–10; incentives, 154; learning how to decide, 99–102; learning to become more rational, 102–9; metacognitive feedback, 164, 165; neuroscience, 198–200; nudging, 154–55, 183–84; optimal boosting, 160–61; optimal incentives, 175; optimal incentive structures, 172–74; prospect theory, 24; providing decision support, 155; resource rationality, 156–59; resource-rational nudging, 176–83; restricting choice to good options, 153–54; risk aversion, 85–86; teaching people better decision strategies, 155–56
- decision mechanisms: model-based, 48; model-free, 48
- decision problem(s), 26, 38, 89, 99, 101, 106, 111, 112, 113
- decision procedures, 21, 76
- decision process(es), 191
- decision science, resource-rational analysis and, 196–97
- decision support, 153, 155
- decision tree: search, 130n4; thinking ahead, 128
- de Montmort, Pierre, on expected value, 16
- directed cognition model, 227
- Donders, Franciscus, work on “mental chronometry,” 111
- “drift rate,” 48
- “Dutch book” argument, 28
- EBA. *See* Elimination-By-Aspects (EBA) heuristic
- ecological rationality, 25–27, 34, 106; Gigerenzer’s theory of, 109; notion of, 10
- economics, 1, 7, 14, 196; optimal incentive design, 172; rational inattention in, 10–11; resource-rational analysis, 197–98
- elementary information processes (EIPs), 111; concept, 210
- Elimination-By-Aspects (EBA) heuristic, 99–101
- environment, 40, 45
- episode, Markov decision process (MDP), 114
- error cost, expected value of, 70, 71
- evidence accumulation, 125, 137n6
- expected utility, 1, 18, 70; action, 74–75, 76; expected opportunity cost and, 77–81; maximizing, 3
- expected utility theory (EUT), 38, 39, 196, 197; cognitive system, 39; Markov decision process (MDP), 119, 120; rational system, 38–42
- expected value, 16
- expected value of control (EVC) theory, 48
- expected value theory, 197
- fast-and-frugal heuristic(s), 93, 99–101; decision-making, 168; resource-rational analysis, 167
- feature vectors (f), 91–92
- feedforward neural network, 95, 96
- framing, 32, 188, 193, 194
- frequency judgments, 25
- game theory, 13n1
- gamification app, to-do list gamification, 175
- Gigerenzer, Gerd: adaptive heuristics, 25; on bounded rationality, 9; ecological rationality, 26; on ecological rationality, 10

- global architecture, 54
- Google Maps, 170
- Google Search, 170
- groups: behavior of, 202–3; decision-making, 190–93
- groupthink, 190, 191

- hash function, 226
- heuristic(s), 14, 23; decision-making, 14–15; Elimination-By-Aspects (EBA), 99–101; fast-and-frugal, 99–101; lexicographic heuristic (LEX), 99–101; optimal boosting teaching people, 160–61; Take-The-Best (TTB), 105–8; weighted additive strategy (WADD), 105–8
- heuristics and biases: anchoring-and-adjustment, 24; availability heuristic, 24; conjunction fallacy, 23; framework, 34; Kahneman and Tversky’s program, 10, 87–88; representativeness heuristic, 23; research program, 23; resource-rational models, 83–87
- hierarchical reinforcement learning (HRL), 139, 202
- homo economicus model, 173
- homunculus, 137–39; term, 138
- Horvitz, Eric: building AI systems, 32; term bounded optimality, 43
- Howes, Andrew: cognition, 33; resource rationality, 37
- human behavior, 2, 7; decision-making, 12–13; model predictions, 62; models of, 13–14; realistic standard, 63. *See also* decision-making
- human cognition. *See* cognition
- human irrationality, decision-making, 3
- human rationality, 3; debate about, 87–88; improving decision-making, 109–10

- Icard, Thomas: bounded optimality, 33; resource rationality, 37
- implementation level, rational process, 30
- importance sampling, 75–76, 78

- incentive(s): decision-making, 154, 157; optimal, 188, 190, 192; optimal structures for decision-makers, 172–74
- independence of irrelevant alternatives, axiom, 18
- induction, backward, 224–26
- inductive bias, rationality, 36
- inductive problems, 28
- intelligent agent(s). *See* artificial intelligence
- intelligent system. *See* artificial intelligence
- intelligent tutoring (systems), 166, 168, 169
- irrationality, 3, 8, 14, 29, 63–64, 74, 87, 88, 198

- Jeffrey, Richard, utility function, 18

- Kahneman, Daniel: heuristics and biases research, 10; rationality, 22
- Keynes, John Maynard, on making decisions, 17

- large language models (LLMs), 168–69, 170
- law and policy, human behavior, 198
- learning mechanism(s), 91–93, 95–97, 109–10, 165
- learning process, modeling, 138
- levels of analysis, Marr’s, 30
- Lewis, Richard: cognition, 33; resource rationality, 37
- LEX. *See* lexicographic heuristic (LEX)
- lexicographic heuristic (LEX), 99–101
- loss aversion: explaining, 85; phenomenon, 84

- machine learning, 67, 69
- Markov chain Monte Carlo (MCMC), 31, 67, 69–70, 141
- Markov decision process(es) (MDPs): contextualizing, 119–20; illustration, 206; Markov property, 222–23; mathematical details, 205–8; metalevel, 114–19, 208–12; navigation problem, 145–46;

- neuroscience, 199; optimal policies and value functions, 207–8; sequential decision problems, 113, 113–14
- Marr, David, levels of analysis, 30
- maximum entropy optimal policy, definition, 208n3
- MCMC. *See* Markov chain Monte Carlo (MCMC)
- MDP(s). *See* Markov decision process(es) (MDPs)
- memoization, 226, 224
- memory, 111, feeling of knowing in memory recall, 132–36, 134
- mental chronometry, 111
- mental state(s), 57; metalevel MDP, 115, 116, 122, 126, 130, 129, 134, 135, 209, 210
- metacognition, 201
- metacognitive feedback, 164, 165
- metacognitive learning, 102, 200
- meta-greedy policy, 226–27
- metalevel Markov decision processes (MDPs), 15, 113, 114–19, 201; attention in preferential choices, 125–28, 126; backward induction, 224–26; Bayesian, 215–19; blinkered approximation, 227–28; cognitive costs, 209; components of, 115; computations, 209, 210–11, 214–15; contextualizing, 119–20; external action, 209; external reward, 209; feeling of knowing in memory recall, 132–36, 134; identifying good policies, 223–28; illustration, 209; mathematical details, 208–12; mental states, 209, 210; multistep lookahead, 227–28; myopic policy, 226–27; optimal policies for, 213–14, 214; policies, 213–14; reward function, 211–12; termination operation, 209; thinking ahead, 128–32, 129; transition function, 211; value of computation (VOC), 223–24; world states, 209, 210
- metalevel policy, Markov decision policy (MDP), 117
- metamemory, 133
- metareasoning; cost of, 148, 149; rational, 140; sequential problem, 118; testing the predictions of rational, model, 97–102
- metareasoning model: learning how to decide, 99–102; learning how to sort, 97–99; testing the predictions of rational, 97–102
- methodological assumption, resource rationality as, 13
- methods: computational, 8, 34, 161, 163, 186; sampling, 66, 75; strategy discovery, 161, 163, 168
- Metropolis-Hastings algorithm, 67
- Monte Carlo. *See* Markov Chain Monte Carlo (MCMC)
- Morgenstern, Oskar, on making decisions, 17
- Mouselab-MDP paradigm, 164
- Mouselab paradigm, 124–25, 131; decision-making, 179; gambling game, 102–3
- multiple-cued recall task, 134, 136
- myopic policy, 226–27; definition, 226
- navigation problem, construals, 144–46, 145
- neural circuits, 138
- neural network(s), feedforward, 95–96, 96
- neuroscience, 7–8, 196; human decision-making, 198–200
- Newell, Allen, human cognition, 111
- non-compensatory, 105–6; environments, 26
- nudge(s). *See* nudging
- nudging, 4, 15, 176; decision-making, 154–55, 157; ethics of, 184; future directions, 183–84; optimal, 176–78, 182; optimal, as framework for choice architectures, 178–80; real-world applications of optimal, 180–83; resource-rational perspective, 179. *See also* optimal nudging

- objective function, 38
opportunity cost, 50, 51, 52
optimal boosting, 188
optimal cognitive system, 40
optimal incentives, 188
optimality, 38, 42
optimal nudging, 176–78, 188; definition, 176; as framework for choice architectures, 178–80; real-world applications of, 180–83. *See also* nudging
optimal policy: Markov decision process (MDP), 114; metalevel MDP, 213–14, 214
optimal rationality enhancement, 185
optimal stopping, decision-making, 20
optimization, 37, 43; Bayesian, 231; framework, 32–33
optimization problem, 5, 7, 9, 49, 54, 60, 171, 184–85, 201
organizations: behavior of, 202–3; decision-making, 190–93

partially observable Markov decision processes (POMDPs), 218–19
perfect rationality, 41
planning: construals, 144–46; subgoals, 142–44
policy, 41; cognitive system, 41; Markov decision process (MDP), 114; policies as programs, 202
POMDP(s). *See* partially observable Markov decision processes (POMDPs)
posterior distribution, 28
predictably irrational, 25
preference reversals, 84
preferences, 17
present bias, resource-rational analysis, 86
prior distribution, 28
probabilistic inference, 67
probability, 5, 25
probability matching, 5, 66
problem domain, 51; metalevel MDP, 115, 116
process models, 67, 70, 190, 192; psychological, 76, 88; rational, 30–32, 34, 65, 86, 96
process-tracing paradigm(s), 62, 121
program π , definition, 43
prosocial, 187, 188, 189–91, 194
prospect theory, 197; decision-making, 24
psychology, 6–7, 14

Q-learning, 231

Ramsey, Frank, on making decisions, 17
rational, learning to become more, 102–9
rational action, new view of the mind, 203–4
rational altruism, 187–93
rational analysis, 2, 27–30, 34; Anderson's, 27, 30, 60; Bayes' rule, 28; definition, 60; inductive problems, 28; method, 2
rational decision-making, formalizing, 17–19
rational inattention, 7, 10–11, 31–32, 34
rationality, 1; alternative definition of, 2; bounded, 8–9, 19–21, 33–34, 90, 202; bounded optimality, 43–46; classical, 1, 5, 6, 33, 63, 64; cognitive systems, 59–60; computational costs of, 10; cost of cognition, 49–53; expected utility theory (EUT), 38–42; formalizing, 38–60; human, 3, 87–88, 109–10; models of human behavior, 13–14; notions for cognitive systems, 39; rational choice of mental action, 46–49; resource rationality, 13, 36, 37, 39, 53–57, 202, 233; sequential case, 57–59; value of cognition (VOC), 46–49
rationality enhancement, 158, 184–85; decision process, 158, 158
rational metareasoning, 8, 14, 48, 140; model, 97–102
rational process model(s), 30–31, 34; strategy selection learning, 96

- reasoning, 6, 12, 24, 28, 29, 31, 48, 63, 66, 80; causal, 70; model-based, 228; model-free, 228; resources, 43
- reference-dependent preferences, 84
- reinforcement learning, 7, 93; hierarchical, 139, 202; model-based, 95; model-free, 78
- representation(s), 140; construals, 144–46; improving mental representations, 193–94; navigation problem, 145–46, 147; resource-rational, 141–46; subgoals, 141–44, 143
- representativeness heuristic, 23
- resource-constrained cognition, 53, 59n6, 136, 139
- resource-rational: altruism, 187–93; anchoring-and-adjustment, 68, 69, 70, 73; architectures, 146–48; cognitive policy, 54–55, 118; cognitive tutor teaching planning strategy, 166; illustration of best intervention, 157; learning from experience, 103–8; model of numerical estimation, 73; nudging, 176–78; perspective on anchoring bias, 64–74; representation, 141–46; teaching, heuristics by demonstration, 161–62
- resource-rational analysis, 2, 3, 4, 12, 14, 36, 60–62, 140; answering questions, 33–35; artificial intelligence (AI), 200–201; assumptions behind, 63; bounded optimality and, 32–33; cognitive science, 15; collective decision-making, 203; computer scientist, 8; curious person, 8; economics, 197–98; economist, 7; expanding the scope of, 148–50; foundations of decision science, 196–97; interventions for decisions, 15; metalevel Markov decision processes (MDPs), 121–36; neuroscientist, 7–8; optimal incentive structures, 172–74; planning algorithms, 142–44; present bias, 86; psychologist, 6–7; simple example of, 4–6; steps of, 60–61; strategy selection, 14–15; subgoals, 142–44, 143
- resource-rational approach, 3
- resource-rational cognitive policy, 54
- resource-rational framing, 188, 194
- resource rationality, 36, 196; cognitive system, 39; definition, 57–58, 58; environment, 56; framework for decision-making enhancements, 156–59; mathematical formalism, 194; methodological assumption, 13; rationality enhancements, 158, 158; theory of, 158–59, 168
- resource-rational model(s), 11, 85; understanding heuristics and biases, 83–87
- resource-rational perspective: anchoring-and-adjustment, 68, 69, 70, 73; availability biases in judgment and decision-making, 74–83
- reward, Markov decision process (MDP), 113
- reward function(s): marginal, 220–21; metalevel MDP, 211–12
- risk aversion, human decision-making, 85–86
- Russell, Stuart: building AI systems, 32; optimization of computational utility, 43
- St. Petersburg paradox, 16, 24
- sampling, 65–66, 67; importance sampling, 75–76, 78; methods, 4–5; Thompson sampling, 94
- SARSA algorithm, 95
- satisficing, 20, 124; heuristic, 86; Simon's model of, 9
- SAT-TTB, 124
- Savage, Leonard, utility function, 18
- sequential case, resource rationality, 57–59
- sequential decision problem(s), 112; cognition as, 112–20; external environments, 113; Markov decision process (MDP), 113, 113–14; metalevel MDPs, 114–19
- Simon, Herbert: bounded rationality, 8–9, 19–21, 202; human cognition, 111; model of satisficing, 9

- Sims, Christopher, rationality in economics, 31
- SOAR, cognitive architecture model, 111, 137
- social choice theory, 192
- social dilemmas, 191, 191–92
- society, decision-making, 190–93
- “softmax” rule, 19, 32
- speed-accuracy trade-off, 74
- state, Markov decision process (MDP), 113, 113
- state space(s), 201, 224–26
- statistics, 1, 67, 69, 71
- stimulus-response learning, 95
- strategies, 15; cognitive, 12, 98–99; model-based, 228; model-free, 228; optimal, 26, 90–92, 104–5, 162–64, 168; resource-rational, 122–23, 161–62, 164, 167–68, 177; sorting, 97–99; weighted additive strategy (WADD), 99–100
- strategy discovery, 108–9, 111, 202; AI-Interpret, 162; attention in preferential choices, 125–28, 126; automatic, 162, 163, 168, 198; feeling of knowing in memory recall, 132–36, 134; metalevel MDP, 123; planning, 128–32, 129; risky choice, 122–25
- strategy selection, 15; computational models, 102; problem, 90–92; rational process model of, 96, 97
- strategy selection learning, 103–8
- strategy selection mechanism(s), 96
- Stroop task, 108–9
- subjective probability, 17
- subjective value, 16, 17
- Take-The-Best (TTB), 26, 87, 105–8, 111, 124; strategy discovery, 122, 123
- tallying heuristic, 217
- teleological explanations, rational model, 37
- termination operation, metalevel MDP, 115, 117
- Theory of Games and Economic Behavior* (von Neumann and Morgenstern), 17
- Thompson sampling, 94
- transition function(s): marginal, 219–20; Markov decision process (MDP), 113; metalevel MDP, 115, 116, 122, 127, 130, 135, 211
- transitivity, axiom, 18
- TTB. *See* Take-The-Best (TTB)
- Tversky, Amos: heuristics and biases research, 10; rationality, 22
- “Type II” rationality, 41
- unbounded optimality, 60
- utility, metalevel MDP, 122
- utility function(s), 18, 38; metalevel MDP, 115, 126, 130, 133
- utility-weighted sampling (UWS), 78, 79; model, 80, 81–83
- UWS. *See also* utility-weighted sampling (UWS)
- value function(s), 24, 207–8, 223–24, 226
- value of cognition (VOC), 39, 51; bounded optimality and, 53–57; cognitive system, 39; metalevel MDPs, 119, 120; rational choice of mental action, 46–49; strategy selection problem, 90–92
- value of computation (VOC), 223–24
- value of information (VOI), 228, 229; bounding, 229–30; illustration of features, 230
- VOC. *See* value of cognition (VOC); value of computation (VOC)
- VOI. *See* value of information (VOI)
- von Neumann, John, on making decisions, 17
- WADD. *See* weighted additive strategy (WADD)
- weighted additive strategy (WADD), 99–100, 105–8, 111, 122, 123, 124
- working memory, 138
- world state(s), 57; metalevel MDP, 115, 116, 122, 126, 130, 133, 209, 210