# Chapter 1

# Unification

## 1.1 The Emotion / Reason Dichotomy

Morals, preferences, goals, utilities, rewards/punishments, and desires all have something in common: the behavior that characterizes them is an optimization. They answer a question about *why* something was done, in a way that is compatible with planning and rational, well-thought-out behavior. If one sees a person repeatedly doing something, say walking their dog, it is reasonable to conclude that either they hold that this is a morally good thing to do, have a preference or goal / sub-goal of walking their dog, enjoy it, or have times when they want to do it. Moreover, these are considered explanations of *why* behavior happens. They are descriptions of the things that agents optimize.

The second feature they share is a subjectivity: anyone can have any preferences or utilities or rewards (perhaps subject to certain constraints if you don't want to be manipulated, want to behave robustly in the presence of adversaries, and so on). Once the space has been constrained, having different preferences is merely... a matter of preference. The theories we use for modeling agents do not take into account the mechanisms by which one might obtain preferences, which

More egregiously still, standard utility / reward maximization picture has nothing to say about the possibility of these quantities changing over time.

All of this is to be contrasted with two other concepts:

1. Empirical analysis which determines that some behavior (as opposed to another one) is occurring.

2. Theories of belief, rationality, and *how* to optimize.

All of these have been postulated as reasons why people do things, and more generally, as theories about ways to shape behavior of agents.

### 1.1.1 Utilities and Preferences

To simplify things, economists and decision theorists start with the case when the set of possibilities is small and finite. If $O$ is the set of things one could choose between, i.e., the set of outcomes, then we can formalize a preference as a partial order $(O, \preccurlyeq)$ on $O$.

However, there is often a lot more structure on outcomes

### 1.1.2   Utilities and Rewards

Both utilities and rewards are real-valued functions from something in the world to a one-dimensional notion of 'good-ness' represented by $\mathbb{R}$, and hence are sometimes thought of as equivalent; here we will do some of the work to explore in what sense they might be equivalent, and the sort of thing that would have to matter for it not to be equivalent.

**Determinism**

If $\mathcal{W}$ is the set of all possible world states, and $\tau : \mathcal{W} \to \mathcal{W}$ is a deterministic function describing the evolution of the world, then having a preference over futures starting at $w_0$ and consistent with $\tau$ is meaningless, as there is only one such sequence of worlds. Therefore, pure determinism only makes sense with cognitively and

## 1.2   The Gödelian Knot

### 1.2.1   Selective Pressure

### 1.2.2   Art History

### 1.2.3   Pedagogy

### 1.2.4

## 1.3   Applications

### 1.3.1   Content Recommendation

The way that these assumptions of static preferences manifest themselves in content recommendation systems is the

### 1.3.2   AI Safety

### 1.3.3   Better Inverse Reinforcement Learning

### 1.3.4   Robotics: Life-long Learning

### 1.3.5   Meta Ethics

## 1.4   As a Learning Problem