



Optimising actions for control objectives

Concept. The idea here is

13.1 States, actions and formulating rewards

Up to this point, we have only considered actions which were either scheduled up front through some fixed process or through user interaction via a game interface. In order to start creating algorithms to act on the system state for us, we now need to develop a formalism which ‘closes the loop’ by feeding information back from the stochastic process to another decision-making process. Note that in most cases, the state of real-world phenomena cannot be measured perfectly. So in order to enable any agent trained on simulated phenomena to potentially act in the real world, we will need to model this measurement process as part of the information retrieval step.

Let’s now define the concept of a ‘measured state’ \mathcal{S}_t of the system (some vector that doesn’t have to share the same length as X_t) at timestep t . We can then say generally that this measured state is ‘observed’ using the following measurement function

$$\mathcal{S}_t^i = M_t^i(X', Z_t, t), \quad (13.1)$$

where we have also extended the definition of Z_t to include parameters which control how this measurement is performed.

In a Markov Decision Process (MDP), based on this observation alone, we would then take actions $X_{t+1} = \mathcal{F}_{t+1}(X', Z_t, \mathcal{A}_t, t)$, for which we would later attribute reward \mathcal{R}_t .

Bibliography