Learning Relative Return Policies With **Upside-Down Reinforcement Learning**

Dylan R. Ashley ^{1,2,3} Kai Arulkumaran^{4,5} Jürgen Schmidhuber ^{1,2,3,6,7} Rupesh Kumar Srivastava⁷

- 1 The Swiss AI Lab IDSIA 💈
- 2 Università della Svizzera italiana 🗳
- 3 Scuola universitaria professionale della Svizzera italiana 🖾
- 4 ARAYA Inc. 💌
- 5 Imperial College London 😹
- 6 AI Initiative, King Abdullah University of Science and Technology 🜌
- 7 NNAISENSE 🖾

Abstract

Lately, there has been a resurgence of interest in using supervised learning to solve reinforcement learning problems. Recent work in this area has largely focused on learning command-conditioned policies. We investigate the potential of one such method—upside-down reinforcement learning—to work with commands that specify a desired relationship between some scalar value and the observed return. We show that upside-down reinforcement learning can learn to carry out such commands online in a tabular bandit setting and in CartPole with non-linear function approximation. By doing so, we demonstrate the power of this family of methods and open the way for their practical use under more complicated command structures.

- 1: $\pi \leftarrow$ initial policy
- 2: while not done do
- pick a command C_0 of the form (desire, horizon)
- generate experience by selecting actions according to π and C_0
- for $S_t, A_t, (D_t, H_t), G_t \in$ new experience do
- train π using input $(S_t, (G_t, H_t))$ and desired output A_t
- end for 7:
- 8: end while

While the formal UDRL paradigm is more general than the above simple version, the underlying principles are the same. The basic idea of UDRL is to have one agent learning command-conditioned policies and a second agent issuing the first agent increasingly challenging commands. The first agent is then trained with the hindsight method.



Istituto Dalle Molle 📃 di studi sull'intelligenza



Scuola universitaria professionale della Svizzera italiana





tablished by the European Commission





The desired return versus the observed return under the learned command-conditioned policy in the CartPole setting. Standard deviation is shown with shading. Note how the observed returns center around the lower-end of the valid returns.







ARAYA

