

Progress Report: Social Feedback for Robotic Collaboration

Emily Wu, Brown University

February 18, 2016

1 Introduction

In my proposal, I introduced a project to enable robots to interact socially with humans in human robot collaboration tasks. In order to establish this social feedback loop, robots must be able to both perceive and convey their mental state to the human partner. Presented below is what I have accomplished so far.

2 Implementing the Domain

For this project, I built a domain describing an object handover interaction on the BURLAP framework. In order to build a domain, the states, actions, observations, reward function, and their interactions must be described within the framework. I have implemented three different variants of these domains, and begun evaluation on them in a simulated environment.

2.1 Base implementation

In the basic implementation, we define the states, actions, transition function, reward function, and observations and observation models as follows:

States are given by a tuple of $\langle \omega, \mathcal{O}, b_H \rangle$. The object the human wishes the robot to handover is w , the set of all objects and their locations and names are given as \mathcal{O} , and the robot’s belief about which object the human believes they will hand over is represented by a distribution b_H over \mathcal{O} .

The **actions** available for the robot to take are to **WAIT** (take no action) **PICK(x)** some object x , **ASK(x)** about some object(s) x , **POINT(x)** to some object, or **LOOK(x)** at some object x .

The **transition function** describes how state changes in response to the robot’s actions. The desired object ω only changes once the robot has **PICKed** the correct object. Similarly, \mathcal{O} only changes if an object was picked up. b_H , the belief of the human about which object the robot hands over, changes with every action the robot takes. Specifically, b_H is updated in a bayesian manner as if the human had observed the robot’s action. Specifically:

$$b_{Ht}(\omega_t) = p(a_t|\omega_t) \sum_{\omega_{t-1}} p(\omega_t|\omega_{t-1}) b_{Ht-1}(\omega_{t-1})$$

$p(a_t|\omega_t)$ gives the probability of observing action a_t (the robot’s last action) given an object ω_t , which is the object the human believes the robot is trying to communicate. this directly mirrors the observation functions described below.

The **observations** that the robot can make of its environment (i.e., the human’s actions) are speech and gesture made by humans. Speech is given by a NLP provided by google, while gesture is provided by a kinect.

The **observation function** describes how observations are produced by the current state. In a POMDP, the current state is hidden, so it is from these observations that the agent (robot) must determine the current state. Our observations pertain mostly to the desired object ω . The details of the observation function are unchanged from the proposal, but I use the same models to update the state variable b_H , except from the perspective receiving observations from the robot’s actions.

The **reward function** describes how the robot receives reward according to its actions and the hidden states. This essentially incentivizes and disincentivizes certain behaviors. In this domain, we provide a large positive reward for **PICKing** the correct object, and a large negative reward for **PICKing** the incorrect object. In addition, several smaller negative rewards are given for taking actions such as **ASKing** the user questions or as a penalty for “annoying” the human. In addition, we give the robot a small reward for how closely b_H matches the distribution representing the robot’s belief about which object the human desires, which we label b_R .

2.2 Incremental Picks

On top of the base implementation, I have experimented with variations that improve the robot’s ability to communicate. The first of these variations is the inclusion of an incremental pick action. In this variation, the robot must deliberately choose the pick action until it observes the pick is complete in order to receive the positive (or negative) reward. This allows the robot to continue taking in information from the human as it is picking, as attempting to pick up the object is an indisputable sign of the robot’s internal state, allowing the human to correct the robot if the robot is picking the incorrect object.

This experiment is running in simulation, where the desired effect of the robot cancelling its pick if additional correcting information is observed. I have not yet run this domain with real robots as solvers (described later) for the POMDP behave poorly with regards to the delayed reward. This suggests approach that abstract out the actions the robot can take to allow for higher level planning, which should be more robust to delayed reward. In the coming semester, I will investigate Abstract MDPs (AMDPs) as a solution to this problem.

2.3 Alternative Reward Function

The reward function of the domain described above is incompatible with the POSS solver used to plan actions (see below), so an alternative had to be devised.

In the new version of the domain, we expanded the state from $\langle \omega, \mathcal{O}, b_H \rangle$ to $\langle \omega_R, \omega_H, \mathcal{O}, b_H \rangle$. In this new state variable, ω_R is the same as the old ω : the object the human desires. The new variable ω_H represents the object the robot has communicated to the human through its actions. This is a hidden variable that the robot must maintain an estimate over, as it is not certain of the interpretation of its actions. This is a parallel interpretation of b_H . The transition function was adjusted such that the underlying distribution over ω_H matched b_H . Actions, observations and the observation function remained unchanged. On top of the existing reward function defined above, we also add an addition reward that operates over ω_R and ω_H .

$$R(\langle \omega_R, \omega_H \rangle) = \begin{cases} 3 & \text{if } \omega_H == \omega_R \\ 0 & \text{otherwise} \end{cases}$$

The reward returned by this function is added to the reward defined in the base domain.

This reward function is designed to offset the cost of taking more expensive actions: if you have correctly communicated the object you believe the human has to the human, the cost of communicative actions is offset. As a result, running the domain in simulation results in more communicative actions which reflect which object the human desires.