# RSA for Bayes Filters and POMDPs

February 4, 2019

## 1 Definitions

- $b_i(s)$ is the listener's current belief in $s$

- $b^0(s \mid u)$ is the primitive interpretation of utterance $u$

- $b_i^d(s \mid u) = \dfrac{O^{d-1}(u \mid s)b(s)}{\displaystyle\sum_{s'} O^{d-1}(u \mid s')b(s')}$ would be the listener's updated belief

  if they interpreted $u$ with RSA of depth $d$

- $O_i^d(u \mid s) = \dfrac{e^{\alpha \ln(b_i^d(s|u))}}{\displaystyle\sum_{u'} e^{\alpha \ln(b_i^d(s|u'))}} = \dfrac{\left(b_i^d(s \mid u)\right)^\alpha}{\displaystyle\sum_{u'} \left(b_i^d(s \mid u')\right)^\alpha}$ is the probability of the

  speaker saying $u$ to communicate $s$ with RSA of depth $d$.

## 2 Desired behavior

We hope that using the listener's current belief will allow utterances to have context-dependent meaning. As an example, we would like for a single, when spoken under belief $b_i$, to be evidence for $s_0$, but when spoken under $b_1$ act as evidence against $s_0$. Equivalently, we wish to find $b_i, b_j, b^0, d, \alpha, u, s_0, s_1$ s.t.

$$(1) \quad \frac{O_i^d(u \mid s_0)}{O_i^d(u \mid s_1)} > 1, \qquad\qquad \frac{O_j^d(u \mid s_0)}{O_j^d(u \mid s_1)} < 1$$

Substituting the speaker formulas gives

$$
(2) \qquad \frac{\dfrac{\left(b_i^d(s_0 \mid u)\right)^\alpha}{\displaystyle\sum_{u'} \left(b_i^d(s_0 \mid u')\right)^\alpha}}{\dfrac{\left(b_i^d(s_1 \mid u)\right)^\alpha}{\displaystyle\sum_{u'} \left(b_i^d(s_1 \mid u')\right)^\alpha}} > 1, \qquad \frac{\dfrac{\left(b_j^d(s_0 \mid u)\right)^\alpha}{\displaystyle\sum_{u'} \left(b_j^d(s_0 \mid u')\right)^\alpha}}{\dfrac{\left(b_j^d(s_1 \mid u)\right)^\alpha}{\displaystyle\sum_{u'} \left(b_j^d(s_1 \mid u')\right)^\alpha}} < 1
$$

$$
(3) \qquad \frac{\left(b_i^d(s_0 \mid u)\right)^\alpha \displaystyle\sum_{u'} \left(b_i^d(s_1 \mid u')\right)^\alpha}{\left(b_i^d(s_1 \mid u)\right)^\alpha \displaystyle\sum_{u'} \left(b_i^d(s_0 \mid u')\right)^\alpha} > 1, \qquad \frac{\left(b_j^d(s_0 \mid u)\right)^\alpha \displaystyle\sum_{u'} \left(b_j^d(s_1 \mid u')\right)^\alpha}{\left(b_j^d(s_1 \mid u)\right)^\alpha \displaystyle\sum_{u'} \left(b_j^d(s_0 \mid u')\right)^\alpha} < 1
$$

$$
(4) \qquad \frac{\left(b_i^d(s_0 \mid u)\right)^\alpha}{\left(b_i^d(s_1 \mid u)\right)^\alpha} > \frac{\displaystyle\sum_{u'} \left(b_i^d(s_0 \mid u')\right)^\alpha}{\displaystyle\sum_{u'} \left(b_i^d(s_1 \mid u)\right)^\alpha}, \qquad \frac{\left(b_j^d(s_0 \mid u)\right)^\alpha}{\left(b_j^d(s_1 \mid u)\right)^\alpha} < \frac{\displaystyle\sum_{u'} \left(b_j^d(s_0 \mid u')\right)^\alpha}{\displaystyle\sum_{u'} \left(b_j^d(s_1 \mid u)\right)^\alpha}
$$

## 2.1 Unraveled recursion formula

$$
(5) \qquad O_i^d(u \mid s) = \frac{\left(b_i^d(s \mid u)\right)^\alpha}{\displaystyle\sum_{u'} \left(b_i^d(s \mid u')\right)^\alpha}
$$

$$
(6) \qquad = \frac{\left(\dfrac{O^{d-1}(u \mid s)b(s)}{\displaystyle\sum_{s'} O^{d-1}(u \mid s')b(s')}\right)^\alpha}{\displaystyle\sum_{u'} \left(\dfrac{O^{d-1}(u' \mid s)b(s)}{\displaystyle\sum_{s'} O^{d-1}(u' \mid s')b(s')}\right)^\alpha}
$$

# 3 Strategies

Write out desired behavior explicitly, find solution.

Uncurl recursion to see the effects of different initial belief after multiple steps.

Run simulations until a good example is found. Will need to write out theory of why that example works afterwards.