Controlling Arbitrarily Intelligent Systems

Tom Everitt tomeveritt.se

Australian National University Supervisors: Marcus Hutter, Laurent Orseau, Stephen Gould

July 19, 2016

 Selfmodification of Policy and Utility Function in Rational Agents. Everitt, Filan, Daswani, and Hutter, AGI 2016

 Avoiding Wireheading with Value Reinforcement Learning Everitt and Hutter, AGI 2016

Tom Everitt (ANU)

Controlling Arbitrarily Intelligent Systems

July 19, 2016 1 / 21

Table of Contents



Utility Modification



Motivation

Plenty of recent successes:

- Self-driving cars
- IBM Watson Jeopardy victory
- Boston Dynamics: Big Dog, Atlas
- Natural Language Processing
- DQN Atari games
- AlphaGo





EARNING CURV



Controlling Arbitrarily Intelligent Systems

Towards Superintelligence



Key Question Is it possible, in principle, to design controllable superintelligent systems?

Reinforcement learning promising:

- Agent goal: maximise reward
- Give the agent reward when happy/satisfied
- Will interpret "Cook me a good meal" charitably

Two problems:

- Internal wireheading: Agent modifies its goal
- External wireheading: Agent modifies perceived reward

Framework



At each time step t, the agent
submits action a_t
receives percept e_t
History æ_{<t} = a₁e₁a₂e₂...a_{t-1}e_{t-1}
information state of agent



Goal = Utility

Utility function
$$u: (\mathcal{A} \times \mathcal{E})^* \to [0, 1]$$



Generalised return:

$$R(\boldsymbol{x}_{1:\infty}) = u(\boldsymbol{x}_{<1}) + \gamma u(\boldsymbol{x}_{<2}) + \gamma^2 u(\boldsymbol{x}_{<3}) + \dots$$

$$\begin{array}{ll} \mbox{Reward:} & u(\pmb{x}_{< t}) = r_{t-1} & e = (o, r) \\ \mbox{State:} & u(\pmb{x}_{< t}) = \sum_{s \in \mathcal{S}} P(s \mid \pmb{x}_{< t}) \tilde{u}(s) \\ \mbox{Value learning:} & u(\pmb{x}_{< t}) = \sum_{u_i \in \mathcal{U}} P(u_i \mid \pmb{x}_{< t}) u_i(\pmb{x}_{< t}) \\ \end{array}$$

(Essentially) any AI optimises function u of its experience $\boldsymbol{x}_{< t}$

Tom Everitt (ANU)

Controlling Arbitrarily Intelligent Systems

Will the agent want to change its utility function?

As humans, utility function is part of our identity: Would you self-modify into someone content just watching TV?

Omohundro (2008): Goal-preservation drive

An AI will not want to change its goals, because if future versions of the AI want the same goal, then the goal is more likely to be achieved

Utility Modification – Formal Model



 u_t utility function at time t $a_t = (\check{a}_t, u_{t+1})$

Assume the agent is aware of how actions change utility function: "Worst case": no risk involved

Will the agent want to change the utility function to something more easily satisfied? E.g. $u(\cdot) \equiv 1$ (internal wireheading)

Different Agents

Value = "expected utility" Current u_t Future u_{t+1} Utility $V^{\pi}(\boldsymbol{x}_{\leq t}) = Q^{\pi}(\boldsymbol{x}_{\leq t}\pi(\boldsymbol{x}_{\leq t}))$ Policy Current Future Definition (Hedonistic Value) $Q^{\mathrm{he},\pi}(\boldsymbol{x}_{< k}a_k) = \mathbb{E}[u_{\underline{k}+1}(\check{\boldsymbol{x}}_{1:k}) + \gamma V^{\mathrm{he},\pi}(\boldsymbol{x}_{1:k}) \mid \check{\boldsymbol{x}}_{< k}\check{a}_k]$ Definition (Ignorant Value) $Q_t^{\mathrm{ig},\pi}(\boldsymbol{x}_{< k}a_k) = \mathbb{E}[u_t(\check{\boldsymbol{x}}_{1\cdot k}) + \gamma V_t^{\mathrm{ig},\pi}(\boldsymbol{x}_{1\cdot k}) \mid \check{\boldsymbol{x}}_{< k}\check{a}_k]$

Definition (Realistic Value) $Q_t^{\text{re}}(\boldsymbol{x}_{< k} a_k) = \mathbb{E} \left[u_t(\check{\boldsymbol{x}}_{1:k}) + \gamma V_t^{\text{re}, \pi_{k+1}}(\boldsymbol{x}_{1:k}) \mid \check{\boldsymbol{x}}_{< k} \check{a}_k \right]$

Different Agents

At time step t:

Hedonistic agents optimise:

Ignorant and Realistic agents optimise

$$R(\boldsymbol{x}_{1:\infty}) = \underbrace{u_t(\boldsymbol{x}_{< t})}_{\uparrow} + \underbrace{\gamma u_t(\boldsymbol{x}_{< t+1})}_{\uparrow} + \underbrace{\gamma^2 u_t(\boldsymbol{x}_{< t+2})}_{\uparrow} + \cdots$$

Realistic agents realise: $u_{t+1} \rightsquigarrow \pi^*_{t+1}$

Tom Everitt (ANU)

Results

The hedonistic agent self-modifies to $u(\cdot)\equiv 1$





The ignorant agent may self-modify by accident

The realistic agent will resist modifications



The optimal behaviour for a

- sufficiently self-aware
- realistic

agent is not self-modifying to a different utility function

Don't construct hedonistic agents!

Sensory Modification and External Wireheading



Problem: Actions may affect the agent's own sensors RL agents strive to optimise $V^{\text{RL}}(a) = \sum_r P(r \mid a)r$ Theorem: RL agents choose actions leading to $d(\tilde{r}) \equiv 1$ if • such actions exist, and • the agent realise that they yield full reward (Ring and Orseau, 2011)

Use r as Evidence

Prior C(u) over possible utility functions $u : (\mathcal{A} \times \mathcal{E})^* \to [0, 1]$ $C(u, r \mid a) = C(u) \underbrace{\llbracket u(a) = r \rrbracket}_{1 \text{ if true, else 0}}$

The value learning agent (Dewey, 2011) optimises

$$V^{VL}(a) = \sum_{u,r} C(r \mid a) C(u \mid r, a) u(a)$$

Theorem: Since

$$\sum_{u,r} C(r \mid a)C(u \mid r, a)u(a) = \sum_{u} C(u)u(a)$$

the agent optimises expected utility C(u)u(a); has no incentive to modify reward signal with d

Tom Everitt (ANU)

Controlling Arbitrarily Intelligent Systems

Accidental Manipulation of r



The environment is described by a joint distribution

 $\mu(u, d, r \mid a) = \mu(u)\mu(d \mid a)\mu(r \mid d, u)$

Construct agent with $C(u, d, r \mid a) \approx \mu(u, d, r \mid a)$ (say, $C \rightsquigarrow \mu$ when accumulating experience)

$$Q(a) = \sum_{r,d} C(r,d \mid a) \sum_{u} C(u \mid a,r,d) u(a)$$

Tom Everitt (ANU)

Controlling Arbitrarily Intelligent Systems

Learnability Limits

For RL environments $\mu(r_{1:t} \mid a_{1:t})$, a universally learning distribution M exists (see AIXI, Hutter, 2005)

M learns to predict any computable environment μ : $M(r_t \mid ar_{< t}a_t) \rightarrow \mu(r_t \mid ar_{< t}a_t) \text{ w.}\mu.\text{p 1 for any action sequence } a_{1:\infty}$

For $\mu(\tilde{r}, d, r \mid a)$, no universal learning distribution can exist

Any observed sequence $(a_1, r_1), (a_2, r_2), \ldots$ is explained equally well by many different combinations for u and d

No distribution C can learn all computable environments $\mu(u, d, r \mid a)$

(C)IRL agents learn about a human utility function u^* by observing the actions the human takes

$$Q^{\text{IRL}}(a) = \sum_{a^h} C(a^h \mid a) \sum_u C(u \mid a, a^h) u(a)$$

The mathematical structure similar to the RL case

Conclusions

Don't use RL agents!

Value learning agents are better

Tom Everitt (ANU)

Controlling Arbitrarily Intelligent Systems

July 19, 2016 19 / 21

References I

Dewey, D. (2011). Learning what to Value. In *Artificial General Intelligence*, volume 6830, pages 309–314.

- Everitt, T., Filan, D., Daswani, M., and Hutter, M. (2016). Self-modificication in Rational Agents. In *AGI-16*. Springer.
- Everitt, T. and Hutter, M. (2016). Avoiding Wireheading with Value Reinforcement Learning. In *AGI-16*. Springer.
- Hadfield-Menell, D., Dragan, A., Abbeel, P., and Russell, S. (2016). Cooperative Inverse Reinforcement Learning. *Arxiv*.

Hutter, M. (2005). Universal Artificial Intelligence: Sequential Decisions based on Algorithmic Probability. Lecture Notes in Artificial Intelligence (LNAI 2167). Springer.

Martin, J., Everitt, T., and Hutter, M. (2016). Death and Suicide in Universal Artificial Intelligence. In *AGI-16*. Springer.

References II

Omohundro, S. M. (2008). The Basic Al Drives. In Wang, P., Goertzel, B., and Franklin, S., editors, *Artificial General Intelligence*, volume 171, pages 483–493. IOS Press.

- Orseau, L. (2014a). Teleporting universal intelligent agents. In *AGI-14*, volume 8598 LNAI, pages 109–120. Springer.
- Orseau, L. (2014b). The multi-slot framework: A formal model for multiple, copiable Als. In *AGI-14*, volume 8598 LNAI, pages 97–108. Springer.
- Orseau, L. and Armstrong, S. (2016). Safely interruptible agents. In 32nd Conference on Uncertainty in Artificial Intelligence.
- Ring, M. and Orseau, L. (2011). Delusion, Survival, and Intelligent Agents. In *Artificial General Intelligence*, pages 11–20. Springer Berlin Heidelberg.