

ulm university universität

IIIm

An Information-theoretic On-line Learning Principle for Specialization in Hierarchical Decision-Making Systems

Heinke Hihn, Sebastian Gottwald, and Daniel A. Braun Ulm University Institute for Neural Information Processing

> December 12, 2019 58<sup>th</sup> IEEE Conference on Decision and Control, Nice

## Emergence of Specialized Decision-Makers Decision-Maker: optimizes a utility *U* in state *s*:

$$a_s^* = rg\max_a U(s, a)$$

#### **Central Idea:**

Limited resources such as

- Linear Decision-Makers
- Limited information processing

drive specialization. 1 2

Motivation: Linear decision-makers are easy to

#### <u>analyze.</u>

<sup>2</sup> Hihn, H., Gottwald, S., and Braun, D. A. (2018). Bounded rational decision-making with adaptive neural network priors. In IAPR Workshop on Artificial Neural Networks in Pattern Recognition.

<sup>&</sup>lt;sup>1</sup>Genewein, T., Leibfried, F., Grau-Moya, J., and Braun, D.A. Bounded rationality, abstraction, and hierarchical decision-making: An information-theoretic optimality principle. Frontiers in Robotics and AI, 2:27, 2015.

# Bounded Rationality and Specialization



Intelligent agents must invest their resources such that they optimally trade off utility versus processing costs <sup>3 4</sup>

Herbert A. Simon coined the term Bounded Rationality

Consequence: Specialization

<sup>&</sup>lt;sup>3</sup>Simon, H. A. A behavioral model of rational choice. The Quarterly Journal of Economics, 69(1):99–118, 1955.

<sup>&</sup>lt;sup>4</sup>Gershman, S. J., et al. Computational rationality: A converging paradigm for intelligence in brains, minds, and machines. Science (2015)

## Information-theoretic Bounded Rationality <sup>5</sup>



$$\max_{p(a|s)} \mathbb{E}_{p(s),p(a|s)} \left[ U(s,a) \right] \text{ s.t. } I(S;A) \le C \tag{1}$$

$$p^*(a|s) = \arg_{p(a|s)} \mathbb{E} \left[ U(s,a) \right] - \frac{1}{\beta} I(S;A) \tag{2}$$

Mutual Information:  $I(S; A) = \mathbb{E}_{p(a|s)} [D_{KL}(p(a|s)||p(a))]$ 

<sup>&</sup>lt;sup>5</sup>Ortega, P. A., and Braun, D.A.. *Thermodynamics as a theory of decision-making with information-processing costs.* Proceedings of the Royal Society of London A: Mathematical, Physical and Engineering Sciences, 469(2153), 2013.

# Hierarchical Decision-Making

Extend to two-level hierarchy with experts  $x \in X^{6}$ 

$$S \to X \to A$$
 (3)

Extended objective:

$$\max_{p(a|s,x),p(x|s)} \mathbb{E}[U(s,a)] - \frac{1}{\beta_1} I(S;X) - \frac{1}{\beta_2} I(S;A|X). \quad (4)$$

<sup>&</sup>lt;sup>6</sup>Genewein, T., Leibfried, F., Grau-Moya, J., and Braun, D.A. Bounded rationality, abstraction, and hierarchical decision-making: An information- theoretic optimality principle. Frontiers in Robotics and AI, 2:27, 2015.

#### Learning via Gradient Descent

Parametrize distributions with parameters  $\theta$  and  $\vartheta$ :

$$J(s, x, a) = U(s, a) - \frac{1}{\beta_1} \log \frac{p_\theta(x|s)}{p(x)} - \frac{1}{\beta_2} \log \frac{p_\theta(a|s, x)}{p(a|x)}$$
(5)

Approximate the prior distributions p(x) and p(a|x) by running means.

#### Utilities for Classification and Regression

1. cross-entropy loss  $\mathcal{L}(y, \hat{y}) = \sum_{i} y_i \log \frac{1}{\hat{y}_i} = -\sum_{i} y_i \log \hat{y}_i$ 

2. mean squared error  $\mathcal{L}(y, \hat{y}) = \sum_{i} (\hat{y}_{i} - y_{i})^{2}$ 

$$\max_{\theta} \qquad \mathbb{E}_{p_{\theta}(x|s)} \left[ \hat{f}(x,s) - \frac{1}{\beta_{1}} \log \frac{p_{\theta}(x|s)}{p(x)} \right] \qquad (8)$$

$$\hat{f}(x,s) = \underbrace{\mathbb{E}_{p_{\vartheta}(\hat{y}|x,s)} \left[ -\mathcal{L}(\hat{y},y) - \frac{1}{\beta_{2}} \log \frac{p_{\vartheta}(\hat{y}|s,x)}{p(\hat{y}|x)} \right]}_{\text{Expert Objective}} \qquad (9)$$

## Classification



## Reinforcement Learning: Setup

Markov Decision Process as a tuple (S, A, P, r), where

- S is the set of states
- A the set of actions
- $P: S \times A \times S \rightarrow [0, 1]$  is the transition probability
- $r : S \times A \rightarrow \mathbb{R}$  is a reward function

Find policy  $\pi_{\theta}$  maximizing expected reward:

$$\theta^* = \arg\max_{\theta} \mathbb{E}_{\tau \sim \pi_{\theta}} \left[ \sum_{\substack{t=0 \ J(\pi_{\theta})}}^{\infty} r(s_t, a_t) \right].$$
(10)

## **RL** Objective

Penalize deviation from a prior policy:

$$\arg\max_{\pi} \mathbb{E}_{\pi} \left[ \sum_{t=0}^{\infty} \gamma^{t} \left( r(s_{t}, a_{t}) - \frac{1}{\beta} \log \frac{\pi(a_{t}|s_{t})}{\pi(a)} \right) \right]. \quad (11)$$

Similar to MaxEnt RL<sup>7</sup>, Trust Region Policy Optimization <sup>8</sup>, Mutual Information Regularized RL<sup>9</sup>

<sup>&</sup>lt;sup>7</sup> Eysenbach, B. and Levine, S. If MaxEnt RL is the Answer, What is the Question?. arXiv preprint (2019).

<sup>&</sup>lt;sup>8</sup>Schulman, J., et al. *Trust region policy optimization*. In International Conference on Machine Learning (2015)

<sup>&</sup>lt;sup>9</sup> Leibfried, F., and Grau-Moya, J. *Mutual-information regularization in markov decision processes and actor-critic learning.* Conference on Robot Learning (2019).

## **RL** Objectives

# Advantage-Actor-Critic <sup>10</sup> Selection Stage Objective:

$$\max_{\theta} \mathbb{E}_{\pi_{\theta}(x|s)} \left[ \hat{f}(s, x) - \frac{1}{\beta_{1}} \log \frac{\pi_{\theta}(x|s)}{\pi(x)} \right], \quad (12)$$

where

$$\hat{f}(s,x) = \underbrace{\mathbb{E}_{\pi_{\vartheta}(a|s,x)} \left[ r(s,a) - \frac{1}{\beta_2} \log \frac{\pi_{\vartheta}(a|s,x)}{\pi(a|x)} \right]}_{\text{Expert Objective}} \quad (13)$$

<sup>&</sup>lt;sup>10</sup> Schulman, J., et al. High-dimensional continuous control using generalized advantage estimation. International Conference on Learning Representations (2015)





#### Reinforcement Learning - Continuous Control Problems



<sup>&</sup>lt;sup>11</sup> Schulman, J., et al. *Trust region policy optimization*. In International Conference on Machine Learning (2015)

## Gain Scheduling

$$\dot{x} = A_i x + B_i u + \epsilon, \quad \text{for } x \in X_i$$
$$B_i = \begin{cases} 1 & \text{if } x \ge 0\\ -1 & \text{if } x < 0 \end{cases}$$
(14)



## Conclusion

- Principled method applicable to a variety of tasks
- Resource limitation drives specialization
- No prior task information required: utility driven partitioning
- Normative framework to analyze hierarchical structures
- System build only by linear decision-makers
- Open Questions
  - High dimensional tasks
  - Sample efficiency in RL