

# RL 18: Self-Motivated Reinforcement Learning

Michael Herrmann

University of Edinburgh, School of Informatics

18/03/2016

## Questions

- How do we define a reward function?
- Where do rewards come from?
- Intrinsic or extrinsic rewards?
- Can an agent “learn” without rewards? What could it possibly learn?
- What actions are worth being explored?
- How can exploration be organised beyond purely random behaviour?

# Why is exploration a problem?

- Not if the number of states and actions is small and time horizon is short
- Exhaustive exploration may be impossible
- Frontier-based exploration becomes impractical in higher dimensions
- Reward signals may not reveal problem structure
- Early success may be misleading

*Intrinsic motivation is defined as the doing of an activity for its inherent satisfaction rather than for some separable consequence. When intrinsically motivated, a person is moved to act for the fun or challenge entailed rather than because of external products, pressures, or rewards.*

Ryan R. M., Deci E. L. (2000). Intrinsic and extrinsic motivations: Classic definitions and new directions. *Contemp. Educ. Psychol.* **25**, 54–67.

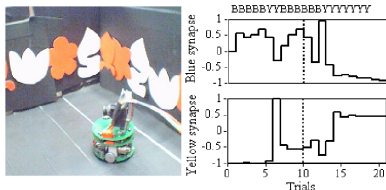
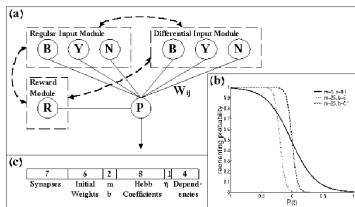
# Intrinsic Motivation: Evolutionary Perspective

- define a  $Q$ -learning agents  $\mathbf{A}$  with reward functions  $r_{\mathbf{A}}$
- do forever
  - set learning rate  $\eta$  and exploration rate  $\varepsilon$
  - for  $i = 1$  to  $N$  do
    - generate a sample  $E_i$  from the environment  $\mathcal{E}$
    - find  $Q$ -function by acting in  $E_i$
    - generate a history  $h_i$  over lifetime of the agent
    - compute fitness  $F(h_i)$
  - select and reduplicate high fitness agents and modify  $r_{\mathbf{A}}$

Result? Fitness increases over generations, but it doesn't matter whether it's a self-generated or external fitness function, see

S. Singh, R.L. Lewis, A.G. Barto, and J. Sorg (2010) Intrinsically Motivated Reinforcement Learning: An Evolutionary Perspective. IEEE Transactions on Autonomous Mental Development 2:2, 70-82.

# Evolution of reinforcement learning in a model



(a) The bee's neural network controller.  
(b) The bee's action function. Probability of reorienting direction of flight as a function of  $P(t)$  for different values of parameters  $m$ ;  $b$ . (c) The genome sequence of the simulated bee.

(a) The foraging robot. (b) Blue and yellow differential weights represent the expected rewards from the two flower colours along the trials. Top: Flower colour chosen in each trial. (blue flowers:  $1/2 \mu\text{l}$  nectar, yellow:  $1 \mu\text{l}$  in half the flowers, contingencies switched after trial 10.)

Y. Niv, D. Joel, I. Meilijson, E. Ruppin (2001) Evolution of reinforcement learning in uncertain environments: Emergence of risk aversion and matching. *Proc. ECAL*.

## Exploration: What is interesting to explore?

- “Interestingness” as a complement to utility could help shaping exploration strategies
- Agent could develop a sense of “curiosity”, e.g.
  - counter-based: states that have not been visited
  - information-based, using a novelty detector
  - “homeokinetic”: Explore as much as possible if still predictable
- or could observe secondary qualities of the learning process in a form of introspection, e.g.
  - learning time or slope of temporal reward average
  - robustness and generality
- Shape/evolve rewards signals as well as exploration strategies

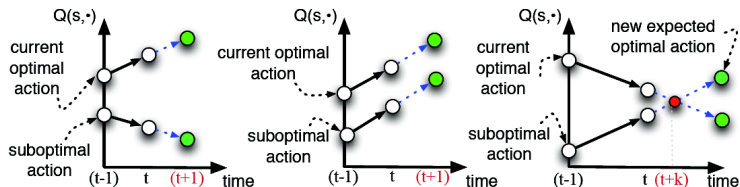
## Intrinsic motivation (Barto, NIPS)

- “Sutton & Barto point out that one should not identify this RL agent with an entire animal or robot.” (Barto, NIPS)
- Exploration of both the external environment and the internal environment
- A sophisticated system should not have to be redesigned for different problems
- Learning a collection of reusable skills in order to generate a skill knowledge base
- Skills could be options (an option is not a sequence of actions; it is a closed-loop control rule, meaning that it is responsive to on-going state changes)

Chentanez, N., Barto, A. G., & Singh, S. P. (2004). Intrinsically motivated reinforcement learning. In NIPS, pp. 1281-1288.



# A Model-Free Algorithm for Efficient Exploration



Determine time until policy change

$$d(s, a) = \frac{1}{\alpha_M} \frac{Q_t(s, a^*) - Q_t(s, a)}{\delta_{s,a}(T_{s,a}) - \delta_{s,a^*}(T_{s,a^*})}$$

$\delta_{s,a}(T_{s,a})$  is the  $\delta$  error for the last time  $(s, a)$  was updated.

$\alpha_M$  is the estimated slope of the expected reward

# A Model-Free Algorithm for Efficient Exploration

Reward based on predicted future usefulness of an action

$$\tilde{r}(s, a) = \begin{cases} \exp\left(-\frac{d^2(s, a)}{\sigma}\right) & \text{if } 0 \leq d(s, a) < \lambda \\ -p & \text{otherwise} \end{cases}$$

$-p$  is a small penalty for stabilisation

$d$  is the expected time until a policy change will occur (see previous slide)

$\sigma$  and  $\lambda$  define a prediction horizon

Da Silva, B. C., & Barto, A. G. (2012) TD- $\Delta\pi$ : A Model-Free Algorithm for Efficient Exploration. 26th Conf. on Artificial Intelligence (AAAI-2012), Toronto, Ontario.

Use generic RL algorithm  $L(s_t, a_t, \tilde{r}_t, r_t, s_{t+1}, Q^t(s', a'))$ , eg. QL.

# A Model-Free Algorithm for Efficient Exploration

For all  $(s, a)$ :  $Q_{\text{exploit}}^0(s, a) \leftarrow 0$ ,  $Q_{\text{explore}}^0(s, a) \leftarrow 0$ ,

$\delta_{s,a}(0) \leftarrow 0$ ,  $T_{s,a} \leftarrow 0$ , visited  $(s, a) \leftarrow \text{False}$

For  $t = 1, 2, \dots$  do

Choose action  $a_t = \arg \max_b Q_{\text{explore}}^t(s, b)$

observe reward  $r_t$  and next state,  $s_{t+1}$

Choose action  $a_t^* = \arg \max_b Q_{\text{exploit}}^t(s, b)$

if not visited  $(s_t, a_t)$  or not visited  $(s_t, a_t^*)$  then  $r(s_t, a_t) = 1$

else if  $|\delta_{(s_t, a_t)}(T_{s_t, a_t}) - \delta_{(s_t, a_t^*)}(T_{s_t, a_t^*})| < \kappa$  then  $r(s_t, a_t) = -p$

else use  $r(s_t, a_t)$  and  $\tilde{r}(s_t, a_t)$  (see previous slide)

$Q_{\text{exploit}}^{t+1}(s, a) \leftarrow L_{\text{exploit}}(s_t, a_t, r_t, s_{t+1}, Q_{\text{exploit}}^t(s', a'))$

$Q_{\text{explore}}^{t+1}(s, a) \leftarrow L_{\text{explore}}(s_t, a_t, \tilde{r}_t, s_{t+1}, Q_{\text{explore}}^t(s', a'))$

$T_{s_t, a_t} \leftarrow t$ , visited  $(s_t, a_t) \leftarrow \text{True}$

$\delta_{s_t, a_t}(t) \leftarrow Q_{\text{exploit}}^{t+1}(s_t, a_t) - Q_{\text{exploit}}^t(s_t, a_t)$

Da Silva, B. C., & Barto, A. G. (2012) TD- $\Delta\pi$ : A Model-Free Algorithm for Efficient Exploration. 26th Conf. on Artificial Intelligence (AAAI-2012), Toronto, Ontario.

# A Model-Free Algorithm for Efficient Exploration

## Discussion

- Does the task actually permit the agent to explore?
- Restricted to discrete state and action spaces
- Performs poorly if many crossings are expected
- Linear approximation questionable as reward often saturates exponentially
- Smoothing and function approximation will be useful
- Using two  $Q$ -functions is often good for combining stability and flexibility

Da Silva, B. C., & Barto, A. G. (2012) TD- $\Delta\pi$ : A Model-Free Algorithm for Efficient Exploration. 26th Conf. on Artificial Intelligence (AAAI-2012), Toronto, Ontario.

“Every time we teach a child something, we keep him from inventing it himself.” (Piaget)

“An AI system can create and maintain knowledge only to the extent that it can verify that knowledge itself.” (Sutton)

A. Turing (“Computing Machinery and Intelligence”, 1950)

“reckoned that it would be easier to write a program to simulate an infant’s mind, rather than an adult’s. The infant program could then be educated much like a human child, until it reached an adult level.”

“The challenge here is to find a learning program which can continuously build on what it knows, to reach increasingly sophisticated levels of knowledge.”

F. Guerin (2011) Learning Like Baby: A Survey of AI approaches. *The Knowledge Engineering Review* 26:02, 209-236.

- Additional rewards from the desire to improve the world model.
- Dynamic Curiosity and Boredom (Schmidhuber, 1991)
- Positive reward if the internal model fails to correctly predict the environment
- e.g. given a predictive model  $M(x_t) = \hat{x}_{t+1}$  we can define intrinsic reward  $r^{(2)} = 1$  if  $|x_{t+1} - \hat{x}_{t+1}| > \vartheta$  and  $r^{(2)} = 0$  otherwise, in addition to an extrinsic rewards signal  $r^{(1)}$ .
- Model is adapted in order to reduce prediction error while action are rewarded for having produced large prediction errors.

J. Schmidhuber (1991) A possibility for implementing curiosity and boredom in model-building neural controllers. In *From Animals to Animats*, 222–227, MIT Press.

J. Schmidhuber, (2010). Formal theory of creativity, fun, and intrinsic motivation (1990–2010). *IEEE Transact. Autonomous Mental Development* 2(3), 230-247.

# How can we define intrinsic motivation?

## 1 Knowledge based models

- Comparisons between the predicted flow of sensorimotor values, (internal forward model) with the actual flow of values
- Adaptive motivation: refers to mechanisms that assign different levels of interest to the same situation

## 2 Competence based models

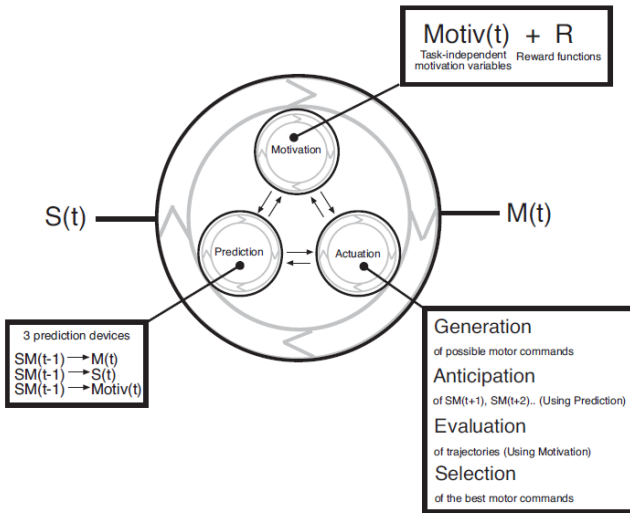
- Characterise the degree of performance/competence
- Comparisons between self-generated goals and the extent to which they are reached in practice (internal inverse model)
- Adaptive motivation

## 3 Morphological models

- Measure immediate structural relationships among multiple sensorimotor channels
- Fixed motivation

Oudeyer, P. Y., & Kaplan, F. (2008). How can we define intrinsic motivation?. Proc. 8th Int. Conf. on Epigenetic Robotics: Modeling cognitive development in robotic systems. Lund Univ. cognitive studies.

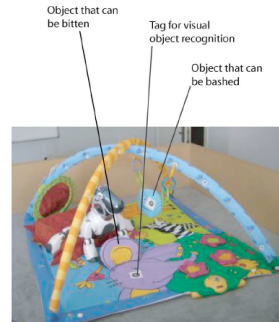
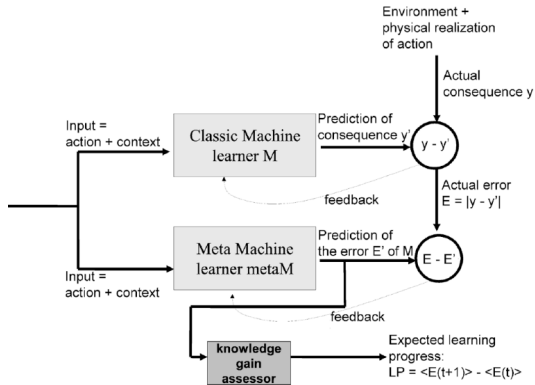
# How can we define intrinsic motivation ?



Kaplan, F., & Oudeyer, P. Y. (2003). Motivational principles for visual know-how development. In C. G. Prince et al. 3. Int. Worksh. Epigen. Robotics, 73–80, Edinburgh, Scotland, Lund Univ. Cogn. Studies.



# The playground experiment: Maximise learning progress



Oudeyer, P. Y., Kaplan, F., & Hafner, V. V. (2007). Intrinsic motivation systems for autonomous mental development. *IEEE Transact. Evol. Comput.* **11**(2), 265-286.

- Homeostasis: Maintain state in a “viable” zone (W. B. Cannon, 1926; W. R. Ashby, 1948)
- Allostasis: achieving stability (homeostasis) through physiological or behavioral change (P. Sterling and J. Eyer, 1988)
- Heterostasis: Drive away from the habitual state (H. Selye, 1973)
- Homeokinesis: Self-organised behaviour aiming at “predictable sensitivity” (R. Der, 1999)

- Aim at state transitions that are predictable  $\Rightarrow$  model with minimal prediction error
- Aim at states where actions have an effect (or at actions that affect the state)  $\Rightarrow$  sensitivity
- Playful behaviour as a compromise between these two conflicting goals
- Self-generated behaviours can be used as options for RL

LPZrobots (<http://robot.informatik.uni-leipzig.de>)

- Predictive model for state transitions  $M(s_t) = \hat{s}_{t+1}$
- Self-evaluation of the model: Sliding average of prediction error
- Choose actions that minimise the 2nd derivative of the prediction error  $\langle |\hat{s}_{t+1} - s_{t+1}|^2 \rangle$
- Result: Agent follow a behaviour as long as it improves in learning. If the rate of the error reduction decays, agent is likely to move on to other behaviours

# Actor-Critic: Heuristic balance model

Soft policies: how soft exactly? Use entropy.

Consider a game between Actor and critic:

- Actor aims at decrease  $\langle H(\pi) \rangle_{\mu(s)}$  in order to get more reward  $\Delta r$ , i.e. the actor transfers entropy into reward. For a given entropy reduction prefer actions that increase  $\Delta r$ .
- Critic aims at increase  $\langle H(\pi) \rangle_{\mu(s)}$  in order to explore, which may (or may not) result in a decrease of the reward. For given entropy reduction prefer actions that decrease  $\Delta r$  least.

Act such as to keep the balance. Balance will obviously shift.

- Intrinsic rewards can
  - speed-up learning
  - generalise beyond known tasks
  - direct exploration
- Can be obtained from
  - From demonstration by inverse reinforcement learning
  - General principles related to homeostasis
  - Successful self-generated options
- Intrinsic rewards are essential in biological and psychological systems