





Looking for the right time to shift strategy in the Exploration-Exploitation Dilemma

Filipo STUDZINSKI PEROTTO



2015 February

Anecdote





XIX Century Immigrants

Summary

- Exploitation-Exploration Dilemma
- Generalities
- Standard Approaches x Proposed Approach
- Experimental Setting and Results
- Conclusions and Perspectives

Exploration – Exploitation Dilemma



Exploration-Exploitation Dilemma

Should I...



...try something new? or just... enjoy a known pleasure?



Exploration-Exploitation Dilemma

Any kind of adaptive agent...

(humans, animals, robots, societies, companies, ...)

...must perform and learn at the same time.

- how to balance it?
- more knowledge = better performance
- but learning = investment



Exploration-Exploitation Dilemma

• When to explore ?

• How to explore ?





Generalities



Markovian Decision Process





Markovian Decision Process (MDP)

$$S = \{s_1, s_2 \dots s_{|S|}\}$$
set of states $A = \{a_1, a_2 \dots a_{|A|}\}$ set of actions $T = \Pr(s' \mid s, a)$ transition function $R = \Pr(r \mid s, a, s')$ reward function

Solution : a policy of actions $f: S \rightarrow A$ which maximizes expected future rewards



Reinforcement Learning



Non-episodic setting : continuous lifetime experience



Agent Lifecycle

- Observe **s**
- Decide
- Execute a
- Receive **r**
- Observe s'
- Learn

- Choose Strategy
- Select Action (from policy)

- Update Models
- Calculate Utility
- Redefine Policy

*model-based learning

Agent Lifecycle

- Observe **s**
- Decide
- Execute a
- Receive **r**
- Observe s'
- Learn

- Choose Strategy
- Select Action (from policy)

- Update Models
- Calculate Utility
- Redefine Policy

*model-based learning

Issues



RL with Limited Resources





• How about costs?



Reinforcement Learning with Energy



New challenge : managing costs, alert under ω_{min}



Standard Approaches



Epsilon-Greedy Method

- Choose exploration with probability V
- Choose exploitation with probability 1-V

* {
$$\epsilon \in \Re \mid 0 \le \epsilon \le 1$$
 }

Optimism in face of incertitude

- Optimistic initialization
 - Optimistic-greedy
 - R-max
 - $Q_0(s, a) > Q^*(s, a)$ $R_0(s) > R(s)$
- Exploration bonus
 - UCB
 - Uncertain or infrequent states



How to explore ?

<u>Undirected</u> :
 = Random

"when exploring, do something unexpected"

 \rightarrow choose action at random

- Directed :
 - = Exploration Bonus

"when exploring, go towards unknown situations"

→ search less visited states
→ or more unpredictable states

Standard Approaches Difficulties

- Random Exploration
 - not very efficient in general
 - not like curiosity
 - can leads the agent to known bad choices
- Epsilon Methods
 - not efficient for sequential problems
- Optimistic Methods
 - efficient exploration
 - but long forced initial training time
 - costs not taken in account

Proposed Approach



- Intuition :
 - Choose the strategy (explore / exploit)
 - Stay engaged for a while (until to reach a "peak")
 - Except in critical situations



- 2 Policies :
 - $-\pi_R$: policy for exploiting
 - $-\pi_{K}$: policy for exploring

- 2 Utilities:
 - $V_{\rm R}$, $Q_{\rm R}$: Exploiting Utility based on reward
 - $V_{\rm K}$, $Q_{\rm K}$: Exploring Utility based on uncertainty

- Parameters :
 - ξ : time in the current strategy
 - ξ_{max} : maximum engagement time
 - r_{peak} : peak reward
 - ω_{crit} : critical energy level alert
 - ε : exploration rate

• Update Utilities :







discount factor $\gamma = 0.9$



discount factor $\gamma = 0.9$





Choose Strategy (Engaged-Climber)



Choose Strategy (I)



Choose Strategy (II)



Choose Strategy (III)



Experience : Chain States

- 10 states, 2 actions :
 - go forward
 - back to start



Experimental Results











Experimental Results

Table 1.: Experimental results

Method	Time to discover the goal	Minimum Score
ε -Greedy	≈ 1200	≈ -400
Optimistic-Greedy	≈ 60	≈ -30
Engaged- $Climber$	≈ 170	≈+10

*using value-iteration method for calculating utilities

Conclusions

- Curiosity
 - Policy for Exploring
 - Policy for Exploiting
- Engagement
 - More human-like behavior
- Consider limited resources

Perspectives

- Concrete search for peaks
- Factored Representations
- Partial Observation
- Multiple Policies
- Robustness







Looking for the right time to shift strategy in the Exploration-Exploitation Dilemma

Filipo STUDZINSKI PEROTTO



2015 February