

Formal Approaches to Decision-Making under Uncertainty

Lecture 4-1: Dealing with Large MDPs

Arnd Hartmanns
Formal Methods and Tools
UNIVERSITY OF TWENTE

Very Large MDPs

Example from ML: Atari games

States: all possible screen images

Actions: joystick movement, button

Reward: winning game points

(+ state of RAM)

128 colors (210 × 160 pixels) =

state space explosion
Bang!
> 10⁷⁰⁸⁰² states



reinforcement learning instead of model checking
using deep neural network as function approximators
→ get good strategy, but lose all optimality guarantees



Q-Learning

Large MDP given implicitly, e.g. via interface

- *init()* → initial state
- *acts(s)* → actions of state s
- *sample(s, a)* → randomly select successor s' and get reward r
- *term(s)* → is s a goal or S_0 state?

but not

- *distr(s, a)* → get full distribution info for action a



Q-Learning

Model checking algorithms:

maintain vectors $x[s]$ or $x_i[s]$

Q-learning:

maintain function $Q: S \times A \rightarrow \mathbb{R}$

with $Q(s, a)$ indicating the "quality" of action a from s

– i.e. an approx. of the goal probability or expected reward



do simulation runs up to goal or S_0 state,
updating the Q -function with
the newly found rewards as you go

Q-Learning

$R(\square G) \quad r_1 + r_2 + r_3 + \dots + r_n \leftarrow \text{goal}$
 discounted: $r_1 + \gamma(r_2 + \gamma(r_3 + \dots + \gamma r_n))$

Algorithm for one episode (simulation run):

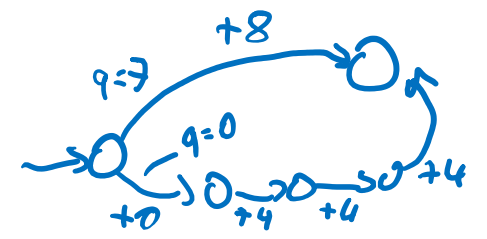
1. $s := \text{init}()$ "ε-greedy strategy"; often start with ε near 1, decrease over more episodes "exploration"
2. with probability ϵ , select uniformly random a from $\text{acts}(s)$;
 with probability $1 - \epsilon$, select $a := \arg \max_{a' \in \text{acts}(s)} Q(s, a')$ "exploitation"
3. $r, s' := \text{sample}(s, a)$

4. $Q(s, a) := Q(s, a) + \alpha \left(r + \gamma \cdot \max_{a' \in \text{acts}(s')} Q(s', a') - Q(s, a) \right)$

(learning rate) "value" of state s' ($x[s']$)
 (discounting factor)

$(1-\alpha)Q(s,a) + \alpha($ $)$

5. $s := s'$
6. if $\neg \text{term}(s')$ then go to 2



Q-Learning

Q-learning algorithm:

1. Perform n learning episodes

2. Return $\max_{a \in \text{acts}(s_I)} Q(s_I, a)$ as approx. for $R_{\max}(\diamond G)$ if $\gamma=1$

and $\mathcal{S}_{\max} = \{ s \mapsto \arg \max_{a \in \text{acts}(s)} Q(s, a) \}$ as optimal scheduler

Fact: $\lim_{n \rightarrow \infty} \max_{a \in \text{acts}(s_I)} Q(s_I, a) = R_{\max}(\diamond G)$ if we play every action infinitely often

But when to stop?

SMC vs. PMC vs. Learning

Memory usage

Runtime

Statistical
Model Checking

constant
(in size of DTMC) } for
DTMC
only

depends
on desired error & confidence

Probabilistic
Model Checking

$O(\text{size of DTMC/MDP})$
to store MDP/DTMC
+ iteration vectors $O(|S|)$

depends
on size of DTMC/MDP
and probabilities & cycles

Reinforcement
Learning (via
Q-learning)

+ do not need DTMC/MDP
fully in memory, just sample
(like in SMC)
- Q-function is $O(|S| \cdot |A|)$
but only for states we
actually visit

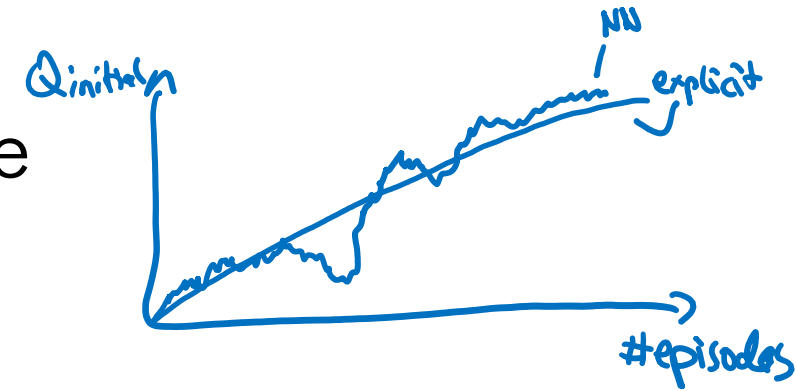
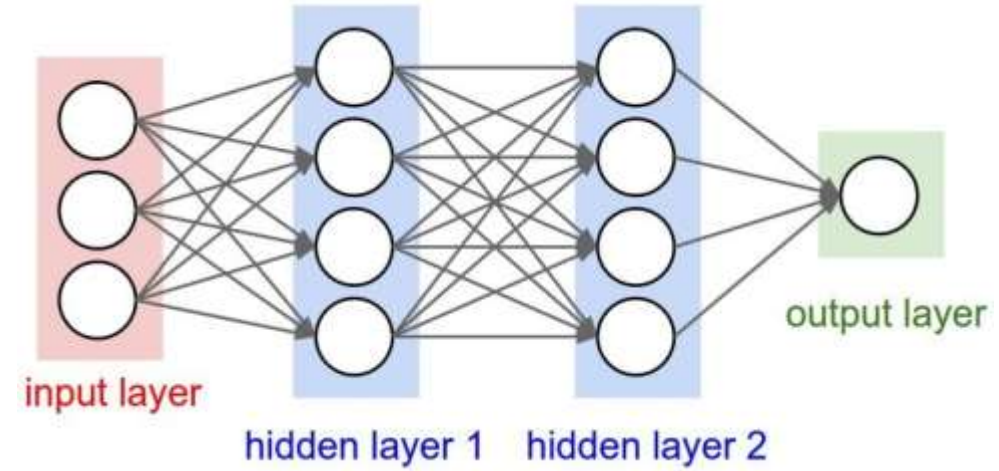
? depends
on # episodes (runs)
(how to determine?)

Deep Learning

Neural networks
are function approximators.

💡 store Q -function
in (deep) neural network

- + fixed memory usage independent of MDP size
- no more $\lim_{n \rightarrow \infty}$ convergence,
learning behaviour very unpredictable



SMC for MDPs

DTMCs: $P(\sigma \models G) = ?$
MDP: $P_{\max}(\sigma \models G) = ?$

What about SMC for MDPs?

→ LSS: lightweight scheduler sampling

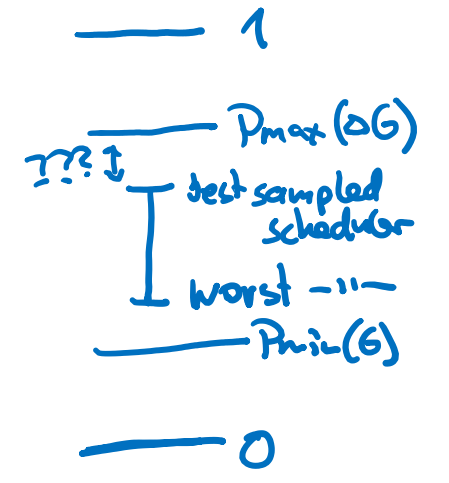
⚡ nondeterminism: must optimise, not estimate

💡 Perform SMC for M randomly chosen schedulers, return max/min

💡 Identify scheduler by single integer σ

(step counter. for step-bounded props.)

$a := (\mathcal{H}(\sigma.s) \bmod |\text{acts}(s)|)$ -th element of $\text{acts}(s)$



O(M)

+ $O(1)$ memory usage like original SMC

- distance from best scheduler found to optimal scheduler unknown