
GENERIC REINFORCEMENT LEARNING AGENTS

Marcus Hutter

Canberra, ACT, 0200, Australia

<http://www.hutter1.net/>



ANU



RSISE



NICTA

Abstract

General purpose intelligent learning agents cycle through (complex, non-MDP) sequences of observations, actions, and rewards. On the other hand, reinforcement learning is well-developed for small finite state Markov Decision Processes (MDPs). It is an art performed by human designers to extract the right state representation out of the bare observations, i.e. to reduce the agent setup to the MDP framework. Before we can think of mechanizing this search for suitable MDPs, we need a formal objective criterion. The main contribution in these slides is to develop such a criterion. I also integrate the various parts into one learning algorithm. Extensions to more realistic dynamic Bayesian networks are briefly discussed.

Contents

- UAI, AIXI, Φ MDP, ... in Perspective
- Agent-Environment Model with Reward
- Universal Artificial Intelligence
- Markov Decision Processes (MDPs)
- Learn Map Φ from Real Problem to MDP
- Optimal Action and Exploration
- Extension to Dynamic Bayesian Networks
- Outlook and Jobs

Universal AI in Perspective

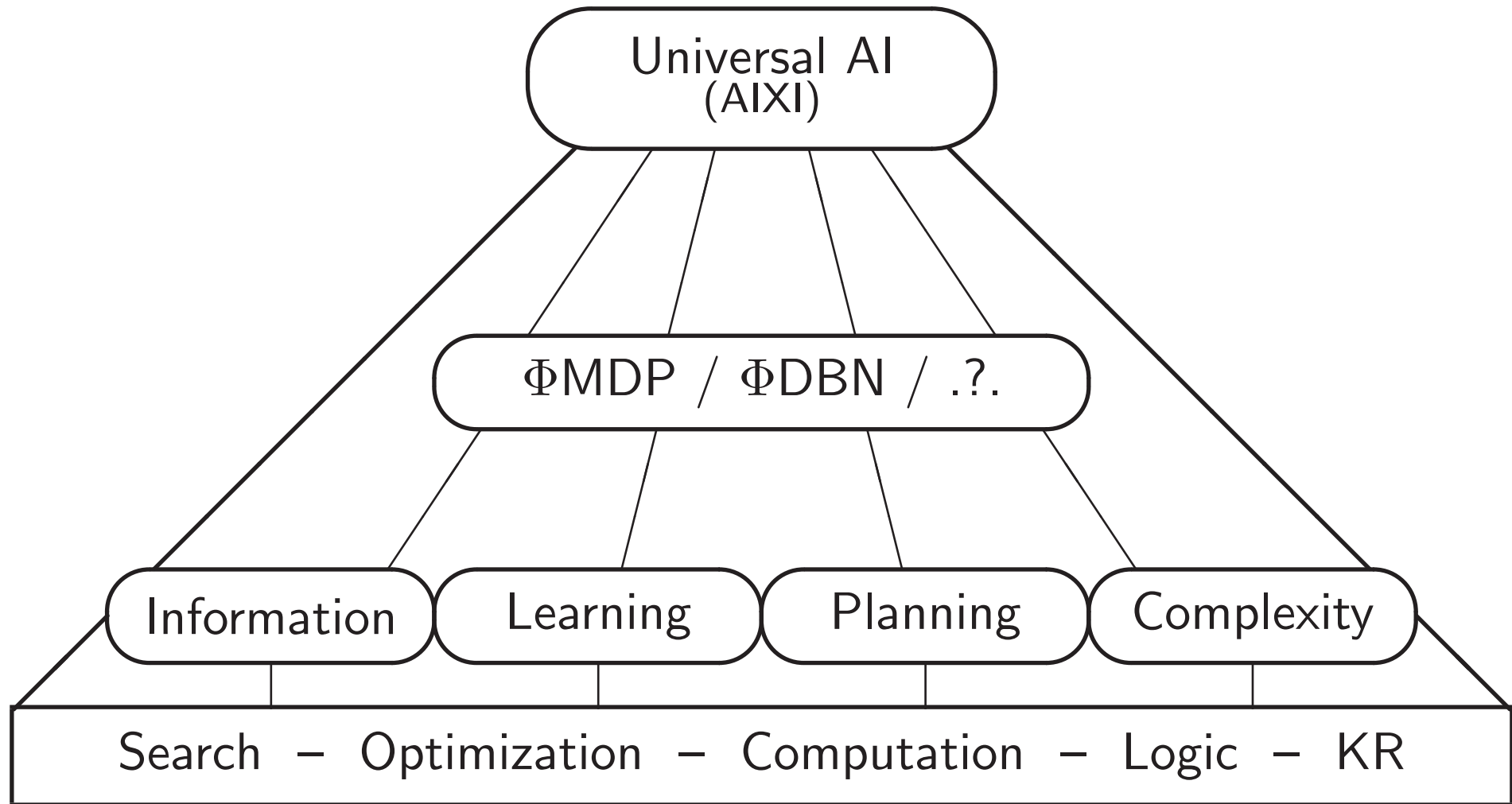
What is A(G)I?	Thinking	Acting
humanly	Cognitive Science	Turing Test
rationally	Laws of Thought	Doing the right thing

Difference matters until systems reach self-improvement threshold

- **Universal AI:** analytically analyzable generic learning systems
- **Real world is nasty:** partially unobservable, uncertain, unknown, non-ergodic, reactive, vast but luckily structured, ...
- **Dealing properly** with uncertainty and learning is crucial.
- Never trust a ~~theory~~ if it is not supported by an ~~experiment~~
experiment
theory

Progress is achieved by an interplay between theory and experiment !

Φ MDP in Perspective



Agents = General Framework, Interface = Robots, Vision, Language

Φ MDP Overview in 1 Slide

Goal: Develop efficient general purpose intelligent agent.

State-of-the-art: (a) AIXI: Incomputable theoretical solution.

(b) MDP: Efficient limited problem class.

(c) POMDP: Notoriously difficult. (d) PSRs: Underdeveloped.

Idea: Φ MDP reduces real problem to MDP automatically by learning.

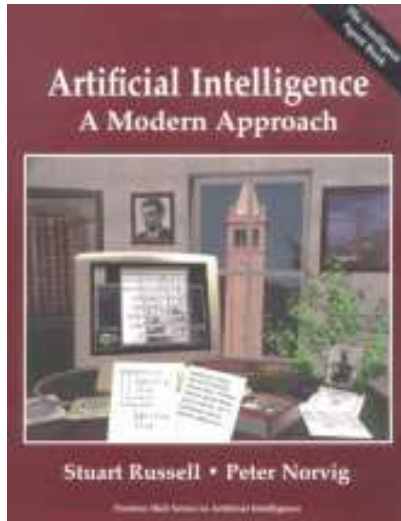
Accomplishments so far: (i) Criterion for evaluating quality of reduction.

(ii) Integration of the various parts into one learning algorithm.

(iii) Generalization to structured MDPs (DBNs)

Φ MDP is promising path towards the grand goal & alternative to (a)-(d)

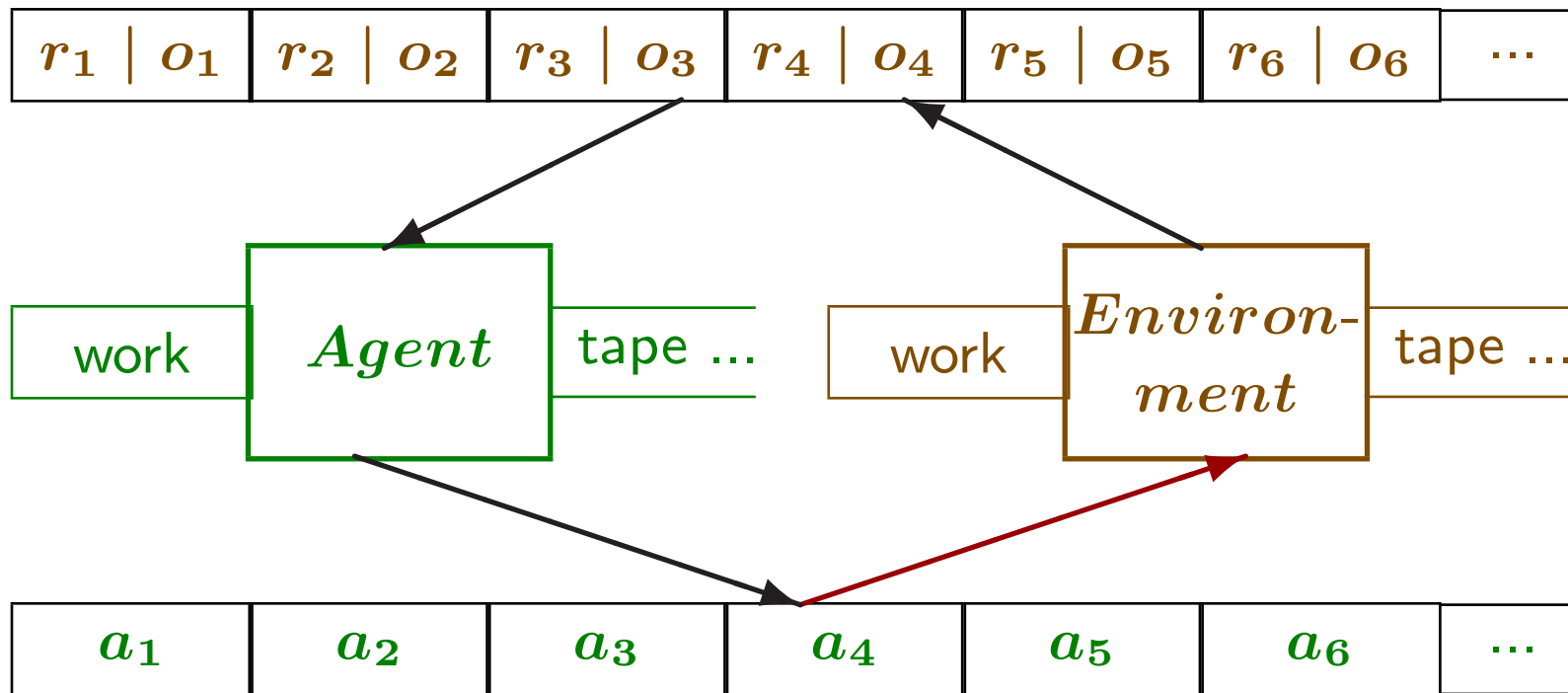
Problem: Find reduction Φ efficiently (generic optimization problem?)



Agent Model with Reward



Framework for **all** AI problems!
Is there a **universal** solution?



Types of Environments / Problems

all fit into the general Agent setup but few are MDPs

sequential (prediction)	⇔	i.i.d (classification/regression)
supervised	⇔	unsupervised
	⇔	reinforcement learning
known environment	⇔	unknown environment
planning	⇔	learning
exploitation	⇔	exploration
passive prediction	⇔	active learning
Fully Observable MDP	⇔	Partially Observed MDP
Unstructured (MDP)	⇔	Structured (DBN)
Competitive (Multi-Agents)	⇔	Stochastic Env (Single Agent)
Games	⇔	Optimization

Universal Artificial Intelligence

Key idea: Optimal action/plan/policy based on the simplest world model consistent with history. Formally ...

$$\text{AIXI: } a_k := \arg \max_{a_k} \sum_{o_k r_k} \dots \max_{a_m} \sum_{o_m r_m} [r_k + \dots + r_m] \sum_{p: U(p, a_1 \dots a_m) = o_1 r_1 \dots o_m r_m} 2^{-\ell(p)}$$

action, *reward*, *observation*, *Universal TM*, *program env.*, $k = \text{now}$

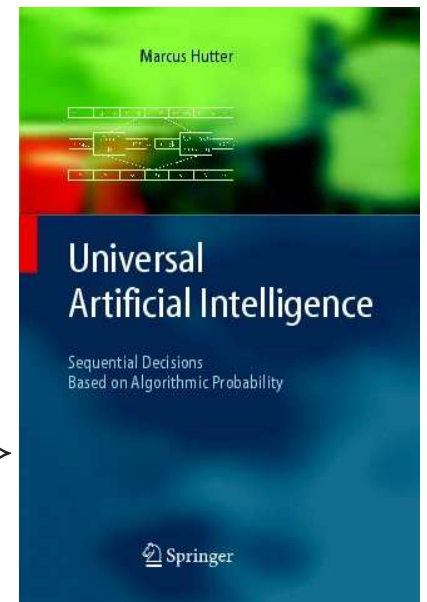
AIXI is an elegant, complete, essentially unique, and limit-computable mathematical theory of AI.

Claim: AIXI is the most intelligent environmental independent, i.e. universally optimal, agent possible.

Proof: For formalizations, quantifications, proofs see \Rightarrow

Problem: Computationally intractable.

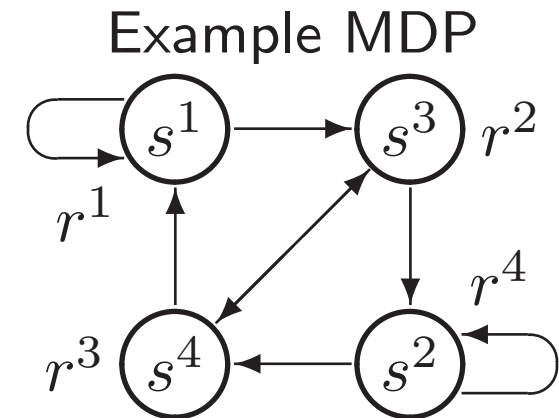
Achievement: Well-defines AI. Gold standard to aim at. Inspired practical algorithms. Cf. infeasible exact minimax.



Markov Decision Processes (MDPs)

a computationally tractable class of problems

- **MDP Assumption:** State $s_t := o_t$ and r_t are probabilistic functions of o_{t-1} and a_{t-1} only.
- **Further Assumption:** State=observation space \mathcal{S} is finite and small.
- **Goal:** Maximize long-term expected reward.
- **Learning:** Probability distribution is unknown but can be learned.
- **Exploration:** Optimal exploration is intractable but there are polynomial approximations.
- **Problem:** Real problems are not of this simple form.



Map Real Problem to MDP

Map history $h_t := o_1 a_1 r_1 \dots o_{t-1}$ to state $s_t := \Phi(h_t)$, for example:

Games: Full-information with static opponent: $\Phi(h_t) = o_t$.

Classical physics: Position+velocity of objects = position at two time-slices: $s_t = \Phi(h_t) = o_t o_{t-1}$ is (2nd order) Markov.

I.i.d. processes of unknown probability (e.g. clinical trials \simeq Bandits),
Frequency of obs. $\Phi(h_n) = (\sum_{t=1}^n \delta_{o_t o})_{o \in \mathcal{O}}$ is sufficient statistic.

Identity: $\Phi(h) = h$ is always sufficient, but not learnable.

Find/Learn Map Automatically

$$\Phi^{best} := \arg \min_{\Phi} \text{Cost}(\Phi | h_t)$$

- What is the best map/MDP? (i.e. what is the right Cost criterion?)
- Is the best MDP good enough? (i.e. is reduction always possible?)
- How to find the map Φ (i.e. minimize Cost) efficiently?

Minimum Description Length Principle

suitable for large semi/non-parametric model classes

MDL Principle: The best model \hat{P} from a class of models \mathcal{P} for data D minimizes the 2-part code length of the data D given probability P plus code length CL of the model:

$$\hat{P} = \arg \min_{P \in \mathcal{P}} \{ |\log P(D)| + CL(P) \}$$

Shortest Code for I.I.D. Data

$D \equiv x_{1:n} \equiv x_1 \dots x_n \in \{1, \dots, d\}^n$ i.i.d, i.e. $P_{\theta}(x_{1:n}) = \theta_{x_1} \dots \theta_{x_n} = \prod_{i=1}^d \theta_i^{n_i}$

Each θ_i need to be coded to precision $O(1/\sqrt{n}) \Rightarrow CL(\theta) = \frac{d}{2} \log n$.

$\Rightarrow \hat{\theta} = \arg \min_{\theta} \{ |\log P_{\theta}(x_{1:n})| + CL(\theta) \} = \dots = \mathbf{n}/n$

If $\mathcal{P} = \{\text{all } P\}$ then still $\hat{P} = P_{\hat{\theta}}$ with high probability.

$CL(x_{1:n}) = CL(\mathbf{n}) := nH(\mathbf{n}/n) + \frac{d}{2} \log n$
 is shortest code length for i.i.d. data $x_{1:n}$.

$H(\hat{\theta}) := -\sum_{i=1}^d \hat{\theta}_i \log \hat{\theta}_i$
 is the entropy of p .

MDP Coding

- Definition:** $n_{ss'}^{ar'} := \#\{t \leq n : s_{t-1} = s, a_{t-1} = a, s_t = s', r_t = r'\}$
 = number of times at which action a in state s resulted in state s' and reward r' ($s \xrightarrow{a} s' r'$).
- Notation:** $\mathbf{n}_\bullet = (n_1, n_2, \dots)$ = vector, $n_+ = n_1 + n_2 + \dots$ = sum.
- Code length of $s_{1:n}$:**
$$\text{CL}(s_{1:n} | a_{1:n}) = \sum_{s,a} \text{CL}(n_{s\bullet}^{a+}).$$
- Code length of $r_{1:n}$ given $s_{1:n}$:**
$$\text{CL}(r_{1:n} | s_{1:n}, a_{1:n}) = \sum_{s'} \text{CL}(n_{+s'}^{+\bullet}).$$

Code is optimal (w.h.p) iff $s_{1:n}$ and $r_{1:n}$ are samples from an MDP.

Φ MDP Cost Criterion

Reward \leftrightarrow State Trade-Off

- Small state space \mathcal{S} has short $\text{CL}(s_{1:n}|a_{1:n})$ but obscures structure of reward sequence $\Rightarrow \text{CL}(r_{1:n}|s_{1:n}a_{1:n})$ large.
- Large $|\mathcal{S}|$ usually makes predicting=compressing $r_{1:n}$ easier, but a large model is hard to learn, i.e. the code for $s_{1:n}$ will be large

$\text{Cost}(\Phi|h_n) := \text{CL}(s_{1:n}|a_{1:n}) + \text{CL}(r_{1:n}|s_{1:n}, a_{1:n})$
is minimized for Φ that keeps all and only relevant information for predicting rewards.

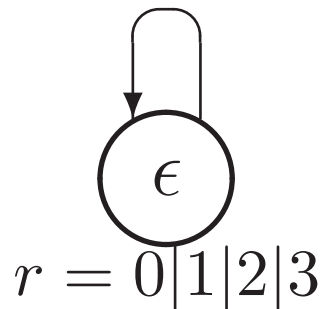
- Recall $s_t := \Phi(h_t)$ and $h_t := a_1 o_1 r_1 \dots o_t$.

A Tiny Example

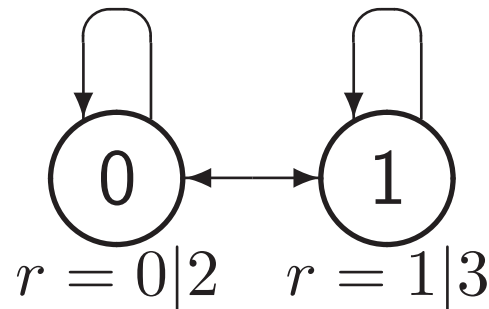
Observations $o_t \in \{0, 1\}$ are i.i.d, and reward $r_t = 2o_{t-1} + o_t$.

Considered features: $\Phi_k(h_t) = o_{t-k+1} \dots o_t =$ last k obs. ($\mathcal{S} = \{0, 1\}^k$)

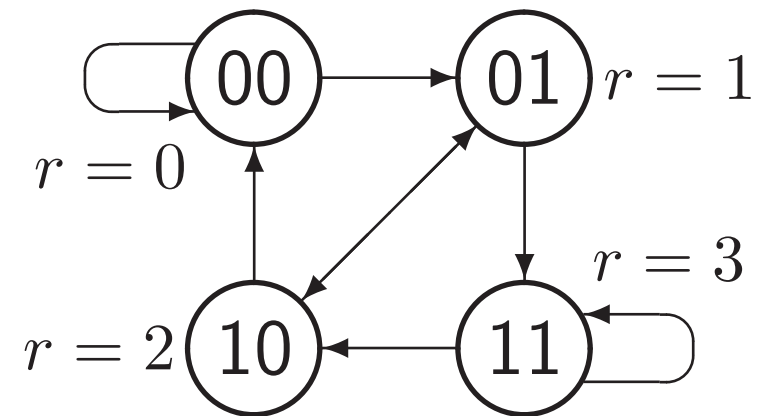
Φ_0 MDP



Φ_1 MDP



Φ_2 MDP



$$\text{Cost}(\Phi_0|h) = 2n + \frac{3}{2} \log n$$

$$\text{Cost}(\Phi_1|h) = 2n + 2 \log \frac{n}{2}$$

$$\text{Cost}(\Phi_2|h) = n + 2 \log \frac{n}{4}$$

$\text{Cost}(\Phi_k|h) = n + 2^{k-1} \log \frac{n}{2^k} (k \geq 2)$. **Cost minimal for $k=2$** (if $n \gg 2^k$).

Intuition: Φ_2 is smallest exact MDP \Rightarrow best observation summary.

Cost(Φ) Minimization

- Minimize $\text{Cost}(\Phi|h)$ by search: random, blind, informed, adaptive, local, global, population based, exhaustive, heuristic, other search.
- Most algs require a neighborhood relation between candidate Φ .
- Φ is equivalent to a partitioning of $(\mathcal{O} \times \mathcal{A} \times \mathcal{R})^*$.
- Example partitioners: Decision trees/lists/grids/etc.
- Example neighborhood: Subdivide=split or merge partitions.

Stochastic Φ -Search (Monte Carlo)

- Randomly choose a neighbor Φ' of Φ (by splitting or merging states)
- Replace Φ by Φ' for sure if Cost gets smaller or with some small probability if Cost gets larger. Repeat.

Context Tree Example

State Space \mathcal{S} = History Suffix Tree

Φ Improve($\Phi_{\mathcal{S}}, h_n$)

[Randomly choose a state $s \in \mathcal{S}$;

Let p and q be uniform random numbers in $[0, 1]$;

if ($p > 1/2$) then split s i.e. $\mathcal{S}' = \mathcal{S} \setminus \{s\} \cup \{os : o \in \mathcal{O}\}$

else if $\{os : o \in \mathcal{O}\} \subseteq \mathcal{S}$

then merge them, i.e. $\mathcal{S}' = \mathcal{S} \setminus \{os : o \in \mathcal{O}\} \cup \{s\}$;

if ($\text{Cost}(\Phi_{\mathcal{S}}|h_n) - \text{Cost}(\Phi_{\mathcal{S}'}|h_n) > \log(q)$) then $\mathcal{S} := \mathcal{S}'$;

[**return** ($\Phi_{\mathcal{S}}$);

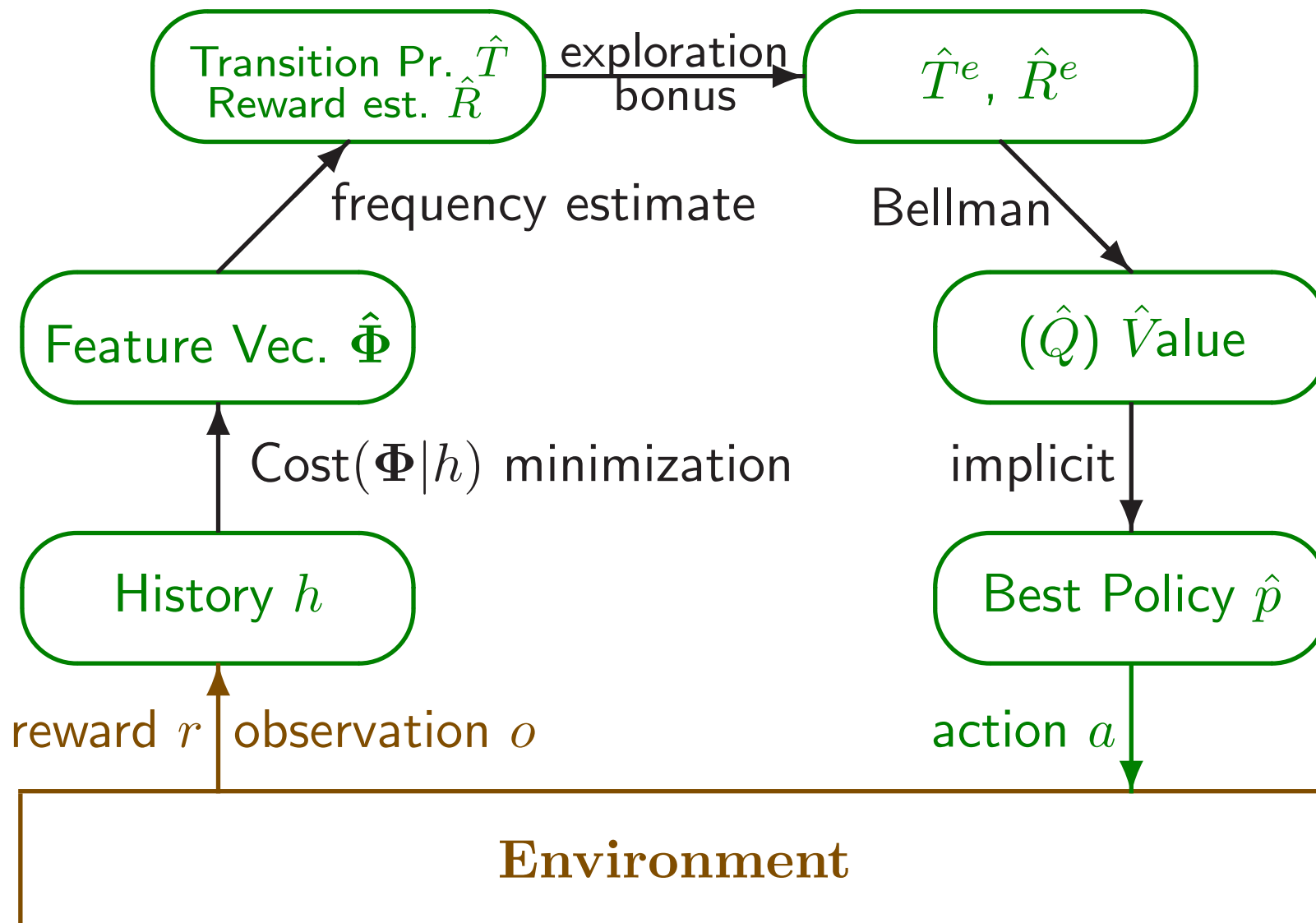
Optimal Action

- Let $\hat{\Phi}$ be a good estimate of Φ^{best} .
 \Rightarrow Compressed history: $s_1 a_1 r_1 \dots s_n a_n r_n \approx$ MDP sequence.
- For a finite MDP with known transition probabilities, optimal action a_{n+1} follows from Bellman equations.
- Use simple frequency estimate of transition probability and reward function \Rightarrow Infamous problem ...

Exploration & Exploitation

- **Polynomially optimal solutions:** Rmax, E3, OIM [KS98,SL08].
- **Main idea:** Motivate agent to explore by pretending high-reward for unexplored state-action pairs.
- Now compute the agent's action based on **modified rewards**.

Computational Flow



Overall Algorithm for Φ MDP Agent

Φ MDP-Agent(\mathcal{A}, \mathcal{R})

[Initialize $\Phi \equiv \epsilon$; $\mathcal{S} = \{\epsilon\}$; $h_0 = a_0 = r_0 = \epsilon$;
for($n = 0, 1, 2, 3, \dots$)
 [Choose e.g. $\gamma = 1 - 1/(n + 1)$;
 Set $R_{max}^e = \text{Polynomial}((1 - \gamma)^{-1}, |\mathcal{S} \times \mathcal{A}|) \cdot \max \mathcal{R}$;
 While waiting for o_{n+1} { $\Phi := \Phi\text{Improve}(\Phi, h_n)$ };
 Observe o_{n+1} ; $h_{n+1} = h_n a_n r_n o_{n+1}$;
 $s_{n+1} := \Phi(h_{n+1})$; $\mathcal{S} := \mathcal{S} \cup \{s_{n+1}\}$;
 Compute action a_{n+1} from Bellman Eq. with explore bonus;
 Output action a_{n+1} ;
 [Observe reward r_{n+1} ;

Extension to Dynamic Bayesian Networks

- Unstructured MDPs are only suitable for relatively small problems.
- Dynamic Bayesian Networks = Structured MDPs for large problems.
- $\Phi(h)$ is now vector of (loosely coupled binary) features=nodes.
- Assign global reward to local nodes by linear regression.
- $\text{Cost}(\Phi, \text{Structure}|h) = \text{sum of local node Costs.}$
- Learn optimal DBN structure in pseudo-polynomial time.
- Search for approximation $\hat{\Phi}$ of Φ^{best} as before.
Neighborhood = adding/removing features.
- Use local linear value function approximation.
- Optimal action/policy by combining [KK99,KP00,SDL07,SL09].

Incremental Updates

Incremental updates/computations can significantly speedup many computations.

- $\text{Cost}(\Phi|h)$ can be computed incrementally when updating Φ or h .
- Local rewards can be updated in linear time.
- Old value estimate can be used as initialization for next cycle.

Outlook

- Real-valued observations suggest real-valued features.
- Extra edges in DBN can improve the linear value function approximation. Needs right incentives in Cost criterion.
- Extend Φ DBN to large structured action spaces.
- Improve various polynomial-time algs. to linear-time algorithms.
- Optimal learning&exploring DBNs required a lot of assumptions.
- Developing smart Φ generation and smart stochastic search algorithms for Φ are the major open challenges.

Conclusion

Goal: Develop efficient general purpose intelligent agent.

State-of-the-art: (a) AIXI: Incomputable theoretical solution.

(b) MDP: Efficient limited problem class.

(c) POMDP: Notoriously difficult. (d) PSRs: Underdeveloped.

Idea: Φ MDP reduces real problem to MDP automatically by learning.

Accomplishments so far: (i) Criterion for evaluating quality of reduction.

(ii) Integration of the various parts into one learning algorithm.

(iii) Generalization to structured MDPs (DBNs)

Φ MDP is promising path towards the grand goal & alternative to (a)-(d)

Problem: Find reduction Φ efficiently (generic optimization problem?)

Connection to AI SubFields

- **Agents:** Φ MDP is a reinforcement learning (single) agent.
- **Search & Optimization:** Minimizing $\text{Cost}(\Phi, \text{Structure}|h)$ is a well-defined but hard (non-continuous, non-convex) optimization problem.
- **Planning:** More sophisticated planners needed for large DBNs.
- **Information Theory:** Needed for analyzing&improving Cost criterion.
- **Learning:** So far mainly reinforcement learning, but others relevant.
- **Logic/Reasoning:** For agents that reason, rule-based logical recursive partitions of domain $(\mathcal{O} \times \mathcal{A} \times \mathcal{R})^*$ are predestined.
- **Knowledge Representation (KR):** Searching for Φ^{best} is actually a search for the best KR. Restrict search space to reasonable KR Φ .
- **Complexity Theory:** We need polynomial time and ultimately linear-time approximation algorithms for all building blocks.
- **Application dependent interfaces:** Robotics, Vision, Language.

Thanks! Questions? Details:

- M.H., *Feature Markov Decision Processes*. (AGI 2009)
- M.H., *Feature Dynamic Bayesian Networks*. (AGI 2009)
- Human Knowledge Compression Contest: (50'000€)
- M. Hutter, *Universal Artificial Intelligence: Sequential Decisions based on Algorithmic Probability*.
EATCS, Springer, 300 pages, 2005.
<http://www.idsia.ch/~marcus/ai/uaibook.htm>



$$\begin{aligned}
 \text{Decision Theory} &= \text{Probability} + \text{Utility Theory} \\
 \text{Universal Induction} &= \text{Ockham} + \text{Bayes} + \text{Turing} \\
 &= \text{A Unified View of Artificial Intelligence}
 \end{aligned}$$

