# Extreme State Aggregation Beyond MDPs

### Marcus Hutter

Research School of Computer Science Australian National University Canberra, ACT, 0200, Australia

http://www.hutter1.net/

12 July 2014

### Abstract

We consider a Reinforcement Learning setup where an agent interacts with an environment in observation-reward-action cycles without any (esp. MDP) assumptions on the environment. State aggregation and more generally feature reinforcement learning is concerned with mapping histories/raw-states to reduced/aggregated states. The idea behind both is that the resulting reduced process (approximately) forms a small stationary finite-state MDP, which can then be efficiently solved or learnt. We considerably generalize existing aggregation results by showing that even if the reduced process is not an MDP, the (q-)value functions and (optimal) policies of an associated MDP with same state-space size solve the original problem, as long as the solution can approximately be represented as a function of the reduced states. This implies an upper bound on the required state space size that holds uniformly for all RL problems. It may also explain why RL algorithms designed for MDPs sometimes perform well beyond MDPs.

#### Contents

1	Introduction	2
<b>2</b>	Feature Markov Decision Processes (ΦMDP)	3
3	Exact Aggregation for $P_{\phi} \in MDP$	6
4	Approximate Aggregation for General P	8
<b>5</b>	Approximate Aggregation Results	10
6	Extreme Aggregation	17
7	Reinforcement Learning	18
8	Feature Reinforcement Learning	22
9	Miscellaneous	25
<b>10</b>	Discussion	25
Re	eferences	26
$\mathbf{A}$	List of Notation	28

### Keywords

state aggregation, reinforcement learning, non-MDP.

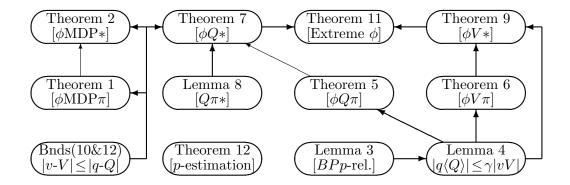
## 1 Introduction

In Reinforcement Learning (RL) [SB98], an agent Π takes actions in some environment P and observes its consequences and is rewarded for them. A well-understood and efficiently solvable [Put94] and efficiently learnable [SLL09, LH12] case is where the environment is (modelled as) a finite-state stationary Markov Decision Process (MDP). Unfortunately most interesting real-world problems P are neither finite-state, nor stationary, nor Markov. One way of dealing with this mismatch is to somehow transform the real-world problem into a small MDP: Feature Reinforcement Learning (FRL) [Hut09c] and U-tree [McC96] deal with the case of arbitrary unknown environments, while state aggregation assumes the environment is a large known stationary MDP [GDG03, FPP04]. The former maps histories into states (Section 2), the latter groups raw states into aggregated states.

Here we follow the FRL approach and terminology, since it is arguably most general: It subsumes the cases where the original process P is an MDP, a k-order MDP, a POMDP, and others (Section 3). Thinking in terms of histories also naturally stifles any temptation of a naive frequency estimate of P (no history ever repeats). Finally we find the history vs state terminologically somewhat neater than raw state vs aggregated state.

More importantly, we consider maps  $\phi$  from histories to states for which the reduced process  $P_{\phi}$  is not (even approximately) an MDP (Section 4). At first this seems to defeat the original purpose, namely of reducing P to a well-understood and efficiently solvable problem class, namely small MDPs. The main novel contribution of this paper is to show that there is still an associated finite-state stationary MDP p whose solution (approximately) solves the original problem P, as long as the solution can still be represented (Section 5). Indeed, we provide an upper bound on the required state space size that holds uniformly for all P (Section 6). While these are interesting theoretical insights, it is a-priori not clear whether they could by utilized to design (better) RL algorithms. We also show how to learn p from experience (Section 7), and sketch an overall learning algorithm and regret/PAC analysis based on our main theorems (Section 8). We briefly discuss how to relax one of the conditions in our main theorems by permuting actions (Section 9). We conclude with an outlook on future work and open problems (Section 10). A list of notation can be found in Appendix A.

The diagram below depicts the dependencies between our results:



# 2 Feature Markov Decision Processes (ΦMDP)

This section formally describes the setup of [Hut09c]. It consists of the agent-environment framework and maps  $\phi$  from observation-reward-action histories to MDP states. This arrangement is called "Feature MDP" or short  $\Phi$ MDP. We use upper-case letters P, Q, V, and  $\Pi$  for the Probability, (Q-)Value, and Policy of the original (agent-environment interactive) Process, and lower-case letters p, q, v, and  $\pi$  for the probability, (q-)value, and policy of the (reduced/aggregated) MDP.

**Agent-environment setup** [Hut09c]. We start with the standard agent-environment setup [RN10] in which an agent  $\Pi$  interacts with an environment P. The agent can choose from actions  $a \in \mathcal{A}$  and the environment provides observations  $o \in \mathcal{O}$  and real-valued rewards  $r \in \mathcal{R} \subseteq [0;1]$  to the agent. This happens in cycles t = 1,2,3,...: At time t, after observing  $o_t$  and receiving reward  $r_t$ , the agent takes action  $a_t$  based on history

$$h_t := o_1 r_1 a_1 \dots o_{t-1} r_{t-1} a_{t-1} o_t r_t \in \mathcal{H}_t := (\mathcal{O} \times \mathcal{R} \times \mathcal{A})^{t-1} \times \mathcal{O} \times \mathcal{R}$$

Then the next cycle t+1 starts. The agent's objective is to maximize its long-term reward. To avoid integrals and densities, we assume spaces  $\mathcal{O}$  and  $\mathcal{R}$  are finite. They may be huge, so this is not really restrictive. Indeed, the  $\Phi$ MDP framework has been specifically developed for huge observation spaces. Generalization to continuous  $\mathcal{O}$  and  $\mathcal{R}$  is routine [Hut09a]. Furthermore we assume that  $\mathcal{A}$  is finite and smallish, which is restrictive. Potential extensions to continuous  $\mathcal{A}$  are discussed in Section 10.

The agent and environment may be viewed as a pair of interlocking functions of the history  $\mathcal{H} := (\mathcal{O} \times \mathcal{R} \times \mathcal{A})^* \times \mathcal{O} \times \mathcal{R}$ :

Env. 
$$P: \mathcal{H} \times \mathcal{A} \leadsto \mathcal{O} \times \mathcal{R}$$
,  $P(o_{t+1}r_{t+1}|h_ta_t)$ , Agent  $\Pi: \mathcal{H} \leadsto \mathcal{A}$ ,  $\Pi(a_t|h_t)$  or  $a_t = \Pi(h_t)$ ,  $action$ 

where  $\rightsquigarrow$  indicates that mappings  $\rightarrow$  are in general stochastic. We make no (stationarity or Markov or other) assumption on environment P. For most parts, environment P is assumed to be fixed, so dependencies on P will be suppressed. For convenience and since optimal policies can be chosen to be deterministic, we consider deterministic policies  $a_t = \Pi(h_t)$  only.

Value functions, optimal Policies, and history Bellman equations. We measure the performance of a policy  $\Pi$  in terms of the P-expected  $\gamma$ -discounted reward sum  $(0 \le \gamma < 1)$ , called (Q-)Value of Policy  $\Pi$  at history  $h_t$  (and action  $a_t$ )

$$V^{\Pi}(h_t) := \mathbb{E}^{\Pi}[R_{t+1}|h_t] \text{ and } Q^{\Pi}(h_t, a_t) := \mathbb{E}^{\Pi}[R_{t+1}|h_t a_t], \text{ where } R_t := \sum_{\tau=t}^{\infty} \gamma^{\tau-t} r_{\tau}$$

The optimal Policy and (Q-)Value functions are

$$V^*(h_t) := \max_{\Pi} V^{\Pi}(h_t) \quad \text{and} \quad Q^*(h_t, a_t) := \max_{\Pi} Q^{\Pi}(h_t, a_t),$$
where  $\Pi^* :\in \arg\max_{\Pi} V^{\Pi}(\epsilon)$  (1)

The maximum over all policies  $\Pi$  always exists [LH14] but may not be unique, in which case argmax denotes the set of optimal policies and  $\Pi^*$  denotes a representative or the whole set of optimal policies. Despite being history-based we can write down (pseudo)recursive Bellman (optimality) equations for the (optimal) (Q-)Values [Hut05, Sec.4.2]:

$$Q^{\Pi}(h_{t}, a_{t}) = \sum_{o_{t+1}r_{t+1}} P(o_{t+1}r_{t+1}|h_{t}a_{t})[r_{t+1} + \gamma V^{\Pi}(h_{t+1})], V^{\Pi}(h_{t}) = Q^{\Pi}(h_{t}, \Pi(h_{t}))$$
(2)  

$$Q^{*}(h_{t}, a_{t}) = \sum_{o_{t+1}r_{t+1}} P(o_{t+1}r_{t+1}|h_{t}a_{t})[r_{t+1} + \gamma V^{*}(h_{t+1})], V^{*}(h_{t}) = \max_{a_{t} \in \mathcal{A}} Q^{*}(h_{t}, a_{t})$$
(3)  

$$\Pi^{*}(h_{t}) \in \arg\max_{a_{t} \in \mathcal{A}} Q^{*}(h_{t}, a_{t})$$
(4)

$$Q^*(h_t, a_t) = \sum_{o_{t+1}r_{t+1}} P(o_{t+1}r_{t+1}|h_t a_t)[r_{t+1} + \gamma V^*(h_{t+1})], \quad V^*(h_t) = \max_{a_t \in \mathcal{A}} Q^*(h_t, a_t) \quad (3)$$

$$\Pi^*(h_t) \in \arg\max_{a_t \in \mathcal{A}} Q^*(h_t, a_t) \tag{4}$$

Unlike their classical state-space cousins (see below), they are not self-consistency equations: The r.h.s. refers to a longer history  $h_{t+1}$  which is always different from the history  $h_t$  on the l.h.s, which precludes any learning algorithm based on estimating the frequency of state/history visits. Still the recursions will be convenient for the mathematical development.

From histories to states  $(\phi)$ . The space of histories is huge and unwieldy and no history ever repeats. Standard ways of dealing with this are to define a similarity metric on histories [McC96] or to aggregate histories [Hut09c]. We pursue the latter via a feature map  $\phi: \mathcal{H} \to \mathcal{S}$  which reduces histories  $h_t \in \mathcal{H}$  to states  $s_t := \phi(h_t) \in \mathcal{S}$ . W.l.g. we assume that  $\phi$  is surjective. We also assume that state space  $\mathcal{S}$  is finite; indeed we are interested in small  $\mathcal{S}$ . This corresponds and indeed is equivalent to a partitioning of histories  $\{\phi^{-1}(s): s \in \mathcal{S}\}$ . Classical state aggregation usually uses the partitioning view [GDG03, Ort07], but the map notation is a bit more convenient here.

The state  $s_t$  is supposed to summarize all relevant information in history  $h_t$ , which lower bounds the size of S. We pass from the complete history  $o_1r_1a_1...o_nr_n$ to a 'reduced' history  $s_1r_1a_1...s_nr_n$ . Traditionally, 'relevant' means that the future is predictable from  $s_t$  (and  $a_t$ ) alone, or technically that the reduced history forms a Markov decision process. This is precisely the condition this paper intends to lift (later).

From histories to MDPs. The probability of the successor states and rewards can be obtained by marginalization

$$P_{\phi}(s_{t+1}r_{t+1}|h_{t}a_{t}) := \sum_{\tilde{o}_{t+1}:\phi(h_{t}a_{t}\tilde{o}_{t+1}r_{t+1})=s_{t+1}} P(\tilde{o}_{t+1}r_{t+1}|h_{t}a_{t})$$
(5)

The reduced process  $P_{\phi}$  is a Markov Decision Process, or Markov for short, if  $P_{\phi}$  only depends on  $h_t$  through  $s_t$ , i.e. is the same for all histories mapped to the same state. Formally

$$P_{\phi} \in \text{MDP} : \iff \exists p : P_{\phi}(s_{t+1}r_{t+1}|\tilde{h}_{t}a_{t}) = p(s_{t+1}r_{t+1}|s_{t}a_{t}) \ \forall \phi(\tilde{h}_{t}) = s_{t}$$
 (6)

Here and elsewhere a quantifier such as  $\forall \phi(\tilde{h}_t) = s_t$  shall mean: for all values of all involved variables consistent with the constraint  $\phi(\tilde{h}_t) = s_t$ . The MDP  $P_{\phi}$  is assumed to be stationary, i.e. independent of t; another condition to be lifted later. Condition (6) is essentially the stochastic bisimulation condition generalized to histories and being somewhat more restrictive regarding rewards [GDG03]: It is a condition on the reward distribution, while [GDG03] constrains its expectation only. This could easily be rectified but is besides the point of this paper. The bisimulation metric [FPP04] is an approximate version of (6), which measures the deviation of  $P_{\phi}$  from being an MDP.

Many problems P can be reduced (approximately) to stationary MDPs [Hut09c]: Full-information games such as chess with static opponent are already Markov, classical physics is approximately 2nd-order Markov, (conditional) i.i.d. processes such as Bandits have counting sufficient statistics, and for a POMDP planning problem, the belief vector is Markov.

Markov decision processes (MDP). We have used and continue to use uppercase letters V, Q,  $\Pi$  for the general process P. We will use lower-case letters v, q,  $\pi$  for (stationary) MDPs p. We use s and a for the current state and action, and s' and r' for successor state and reward. Consider a stationary finite-state MDP  $p: \mathcal{S} \times \mathcal{A} \leadsto \mathcal{S} \times \mathcal{R}$  and stationary deterministic policy  $\pi: \mathcal{S} \to \mathcal{A}$ . Only in Section 3 will this p be given by (6), but in general p will be different from (6). In any case, the p-expected  $\gamma$ -discounted reward sum, called (q-)value of (optimal) policy  $\pi^{(*)}$  in MDP p, are given by the Bellman (optimality) equations

$$q^{\pi}(s,a) = \sum_{s'r'} p(s'r'|sa)[r'+\gamma v^{\pi}(s')] \text{ and } v^{\pi}(s) = q^{\pi}(s,\pi(s))$$
 (7)

$$q^*(s,a) = \sum_{s'r'}^{sr} p(s'r'|sa)[r'+\gamma v^*(s')] \quad \text{and} \quad v^*(s) = \max_a q^*(s,a)$$
 (8)

$$\pi^*(s) \in \arg\max_{a} q^*(s, a).$$
 Note:  $v^{\pi}(s) \le v^*(s), \ q^{\pi}(s, a) \le q^*(s, a)$  (9)

Using p(s'r'|sa) = p(r'|sas')p(s'|sa) we could also rewrite them in terms of transition matrix p(s'|sa) and expected reward  $\mathbb{E}[r'|sa]$  [SB98].

**More notation.** While our equations often assume or imply  $s = s_t$ ,  $a = a_t$ ,  $s' = s_{t+1}$ ,  $r' = r_{t+1}$ , (and  $h_{t+1} = h_t a o' r'$ ) for some t, technically s, a, s', r' are different variables from all variables in history  $h_n = o_1 r_1 a_1 \dots o_t r_t a_t o_{t+1} r_{t+1} \dots a_n r_n$ . Less prone to confusion are  $o = o_t$ ,  $o' = o_{t+1}$ ,  $h = h_t$ , h' = h a o' r'.

We call a function f(h), piecewise constant or  $\phi$ -uniform iff  $f(h) = f(\tilde{h})$  for all  $\phi(h) = \phi(\tilde{h})$ . Here and elsewhere  $\forall \phi(h) = \phi(\tilde{h})$  is short for  $\forall h, \tilde{h}: \phi(h) = \phi(\tilde{h})$ . Similarly  $\forall s = \phi(h)$  is short for  $\forall s, h: s = \phi(h)$ . Etc.

The Iverson bracket,  $[\![R]\!] := 1$  if R = true and  $[\![R]\!] := 0$  if R = false, denotes the indicator function. Throughout,  $\varepsilon, \delta \ge 0$  denote approximation accuracy. Note that this includes the exact = 0 case.

We now show that if P reduces via  $\phi$  to an MDP p, the solution of these equations yields (Q-)Values and optimal Policy of the original process P. This is not surprising and just a history-based versions of classical state-aggregation results [GDG03]. We state and prove them here, since notation and setup are somewhat different, and proof ideas and fragments will be reused later.

# 3 Exact Aggregation for $P_{\phi} \in MDP$

The following two theorems show that if  $\phi$  reduces P to a stationary MDP via (5) and (6), then V and Q (and  $\Pi^*$ ) essentially coincide with v and q (and  $\pi^*$ ), where policy  $\Pi$  ( $\Pi^*$ ) has to be assumed (will be shown) constant within each partition  $\phi^{-1}(s)$ . This allows to efficiently solve for (and learn in the case of unknown P) V and Q (and  $\Pi^*$ ) in time polynomial in S by solving/learning (7) (or (8) and (9)) instead of (2) (or (3) and (4)).

**Theorem 1** ( $\phi$ MDP $\pi$ ) Let  $\phi$  be a reduction such that  $P_{\phi} \in MDP$  reduces to MDP p defined in (6), and let  $\Pi$  be some policy such that  $\Pi(h) = \Pi(\tilde{h})$  for all  $\phi(h) = \phi(\tilde{h})$ . Then for all a and h it holds:

$$V^\Pi(h) = v^\pi(s) \quad and \quad Q^\Pi(h,a) = q^\pi(s,a), \quad where \quad \pi(s) := \Pi(h) \quad and \quad s = \phi(h)$$

Note that  $\pi(s)$  is well-defined, since  $\phi$  is surjective and  $\Pi(h)$  is the same for all  $h \in \phi^{-1}(s)$ . The standard proof considers an m-horizon truncated MDP and induction on m and  $m \to \infty$ . Besides the adaptation to histories, the proof below is a slight variation that avoids such truncation and limit. This style will be useful later. We explain all steps in detail here, since variations will be utilize later.

**Proof.** Let  $\delta := \sup_{s=\phi(h),a} |q^{\pi}(s,a) - Q^{\Pi}(h,a)|$ . Using  $a' := \pi(s') = \Pi(h')$  for  $s' = \phi(h')$  and (2) and (7) lets us bound the value difference

$$|v^{\pi}(s') - V^{\Pi}(h')| = |q^{\pi}(s', a') - Q^{\Pi}(h', a')| \le \delta \quad \forall s' = \phi(h')$$
 (10)

For any a and h, this implies

$$Q^{\Pi}(h,a) \stackrel{(a)}{=} \sum_{o'r'} P(o'r'|ha)[r' + \gamma V^{\Pi}(h')] \qquad [h' = hao'r']$$

$$\stackrel{(b)}{\leq} \sum_{s'r'} \sum_{o':\phi(h')=s'} P(o'r'|ha)[r' + \gamma(v^{\pi}(s') \pm \delta)] \qquad (11)$$

$$\stackrel{(c)}{=} \sum_{s'r'} P_{\phi}(s'r'|ha)[r' + \gamma v^{\pi}(s')] \pm \gamma \delta$$

$$\stackrel{(d)}{=} \sum_{s'r'} p(s'r'|sa)[r' + \gamma v^{\pi}(s')] \pm \gamma \delta \qquad [s := \phi(h)]$$

$$\stackrel{(e)}{=} q^{\pi}(s,a) \pm \gamma \delta$$

(a) is just (2). In (b) we sum over all o' by first summing over all o' such that  $\phi(hao'r')=s'$  and then summing over all s'. We have also upper/lower bounded  $V^{\Pi}(h')$  via (10). (c) is the definition (5) of  $P_{\phi}$  and pulls out  $\gamma\delta$  using that probability  $P_{\phi}$  sums to 1. (d) is the definition (6) of p. (e) is simply (7). The chain (11a-e) holds for all  $s=\phi(h)$  and a, hence

$$\delta = \sup_{s = \phi(h), a} |q^{\pi}(s, a) - Q^{\Pi}(h, a)| \leq \gamma \delta \quad \Rightarrow \quad \delta \leq 0$$

Hence 
$$v^{\pi}(s) = V^{\Pi}(h)$$
 and  $q^{\pi}(s,a) = Q^{\Pi}(h,a)$  for all  $s = \phi(h)$  and  $a$ .

**Theorem 2** ( $\phi$ MDP\*) Let  $\phi$  be a reduction such that  $P_{\phi} \in MDP$  reduces to MDP p defined in (6), Then for all a and h it holds:

$$\Pi^*(h) = \pi^*(s)$$
 and  $V^*(h) = v^*(s)$  and  $Q^*(h, a) = q^*(s, a)$ , where  $s = \phi(h)$ 

The core of the proof follows the same steps (11a-e) as for the previous theorem, but the rest is slightly different. Additionally we have to show that  $\Pi^*$  is piecewise constant (in Theorem 1 we assumed  $\Pi$  was).

**Proof.** Let  $\delta := \sup_{s=\phi(h),a} |q^*(s,a) - Q^*(h,a)|$ . We can bound the value difference

$$|v^*(s) - V^*(h)| \stackrel{(a)}{=} |\max_{a} q^*(s, a) - \max_{a} Q^*(h, a)| \stackrel{(b)}{\leq} \max_{a} |q^*(s, a) - Q^*(h, a)| \stackrel{(c)}{\leq} \delta \, \forall s = \phi(h)$$
(12)

(a) follows from the definitions (3) and (8). (b) follows from the following general elementary frequently used bound

$$\left| \max_{x} f(x) - \max_{x} g(x) \right| \le \max_{x} \left| f(x) - g(x) \right| \tag{13}$$

(c) follows from the definition of  $\delta$ .

One now can show that  $Q^*(h,a) \leq q^*(s,a) \pm \gamma \delta$  for  $s = \phi(h)$  by following exactly the same steps as (10a-e) just with  $\Pi$  and  $\pi$  replaced by \* and using (12) instead of (10), and using the Bellman optimality equations (3) and (8) instead of the Bellman equations (2) and (7). Also as before, this implies  $\delta \leq \gamma \delta$ , hence  $\delta \leq 0$ , hence  $v^*(s) = V^*(h)$  and  $q^*(s,a) = Q^*(h,a)$  for all  $s = \phi(h)$  and a. Finally, the latter implies  $\pi^*(s) = \operatorname{argmax}_a q^*(s,a) = \operatorname{argmax}_a Q^*(h,a) = \Pi^*(h)$ .

Approximate aggregation results if  $P_{\phi}$  is approximately MDP can also be derived [FPP04]. The core results in the next section show that aggregation is possible far beyond  $P_{\phi}$  being approximately MDP.

# 4 Approximate Aggregation for General P

This section prepares for the main technical contribution of the paper in the next section. The key quantity to relate original and reduced Bellman equations is a form of stochastic inverse of  $\phi$ , whose choice and analysis will be deferred to Section 7.

**Dispersion probability** B. Let  $B_{\phi}: \mathcal{S} \times \mathcal{A} \leadsto \mathcal{H}$  be a probability distribution on finite histories for each state-action pair such that  $B_{\phi}(h|sa) = 0$  if  $s \neq \phi(h)$ .  $B \equiv B_{\phi}$  may be viewed as a stochastic inverse of  $\phi$  that assigns non-zero probability only to  $h \in \phi^{-1}(s)$ . The formal constraints we pose on B are

$$B(h|sa) \ge 0$$
 and  $\sum_{h \in \mathcal{H}} B(h|sa) = \sum_{h:\phi(h)=s} B(h|sa) = 1 \quad \forall s, a$  (14)

This implicitly requires  $\phi$  to be surjective, i.e.  $\phi(\mathcal{H}) = \mathcal{S}$ , which can always be made true by defining  $\mathcal{S} := \mathcal{S}_{\phi} := \phi(\mathcal{H})$ . Note that the sum is taken over histories of any/mixed length. In general, B is a somewhat weird distribution, since it assigns probabilities to past and future observations given the current state and action. The interpretation and choice of B does not need to concern us, except later when we want to learn p.

The MDP requirement (6) will be replaced by the following definition:

$$p(s'r'|sa) := \sum_{h \in \mathcal{H}} P_{\phi}(s'r'|ha)B(h|sa)$$

$$\equiv \sum_{t=1}^{\infty} \sum_{h_{t} \in \mathcal{H}_{t}} P_{\phi}(s_{t+1} = s', r_{t+1} = r'|h_{t}, a_{t} = a)B(h_{t}|sa)$$
(15)

That is, the finite-state stationary MDP p is built from feature map  $\phi$ , dispersion probability B, and environment P: The p-probability of observing state-reward pair (s',r') from state-action pair (s,a) is defined as the B-average over all histories h consistent with (s,a) of the  $P_{\phi}$ -probability of observing (s',r') (obtained from P by  $\phi$ -marginalizing) given history h and action a. The r.h.s. of the first line is merely shorthand for the second line. Note that sas'r' are fixed and do not appear in h

which ranges over histories  $\mathcal{H}$  of all lengths. It is easy to see that p is a probability distribution, and it is Markov by definition. If  $P_{\phi} \in \text{MDP}$ , then definition (15) coincides with p defined in (6). In general, the MDP p, depending on arbitrary B, is not the state distribution induced by P (and  $\Pi$ ), which in general is non-Markov. Note that p is a stationary MDP for any B satisfying (14) and  $any \phi$  and P. We need the following lemmas:

**Some lemmas.** The first lemma establishes the key relation between P and p via B used later to relate original history Bellman (optimality) equations (2–4) with reduced state Bellman (optimality) equations (7–9).

**Lemma 3 (B-P-p relation)** For any function  $f: S \times R \to \mathbb{R}$  and p defined in (15) in terms of P via (5), and  $s' := \phi(h')$  and h' := hao'r' it holds

$$\sum_{h \in \mathcal{H}} B(h|sa) \sum_{o'r'} P(o'r'|ha) f(s',r') = \sum_{s'r'} p(s'r'|sa) f(s',r')$$

$$\underset{depends\ on\ hao'r'}{\uparrow} p(s'r'|sa) f(s',r')$$

Proof.

$$\sum_{h \in \mathcal{H}} B(h|sa) \sum_{o'r'} P(o'r'|ha) f(s',r')$$

$$\stackrel{(a)}{=} \sum_{h \in \mathcal{H}} B(h|sa) \sum_{s'r'} \sum_{o':\phi(h')=s'} P(o'r'|ha) f(s',r')$$

$$\stackrel{(b)}{=} \sum_{h \in \mathcal{H}} B(h|sa) \sum_{s'r'} P_{\phi}(s'r'|ha) f(s',r')$$

$$\stackrel{(c)}{=} \sum_{s'r'} p(s'r'|sa) f(s',r')$$

In (a) we sum over all o' by first summing over all o' such that  $\phi(hao'r') = s'$  and then summing over all s'. In (b) we used the definition (5) of  $P_{\phi}$ . In (c) we used the definition (15) of p.

Inequalities (10) and (12) trivially bound v-V differences in terms of q-Q differences:  $|v-V| \le \max_a |q-Q|$ . The following lemma shows that a reverse holds in expectation, i.e.  $|q-\langle Q\rangle_B| \le \gamma |v-V|$ . The expectation can (only) be dropped if Q is constant for  $h \in \phi^{-1}(s)$ . Formally define

$$\langle f(h,a)\rangle_B := \sum_{\tilde{h}\in\mathcal{H}} B(\tilde{h}|sa)f(\tilde{h},a), \text{ where } s := \phi(h)$$
 (16)

That is,  $\langle f(h,a)\rangle_B$  takes a *B*-average over all  $\tilde{h}$  that  $\phi$  maps to the same state as h. For convenience we will drop the tilde, which we can do if we declare  $s := \phi(h)$  to refer to the 'global' h in  $\langle f(h,a)\rangle_B$  and not to the 'local' variable in the  $h \in \mathcal{H}$  sum.

Lemma 4 ( $|q-\langle Q\rangle| \leq \gamma |v-V|$ ) For any P,  $\phi$ , B, define p via (15) and (5). (i) If  $|v^{\pi}(s)-V^{\Pi}(h)| \leq \delta \ \forall s = \phi(h)$ then  $|q^{\pi}(s,a)-\langle Q^{\Pi}(h,a)\rangle_{B}| \leq \gamma \delta \ \forall s = \phi(h) \ \forall a$ . (ii) If  $|v^{*}(s)-V^{*}(h)| \leq \delta \ \forall s = \phi(h)$ then  $|q^{*}(s,a)-\langle Q^{*}(h,a)\rangle_{B}| \leq \gamma \delta \ \forall s = \phi(h) \ \forall a$ .

**Proof.** (i) Let  $s := \phi(h)$  and h' := hao'r' and  $s' := \phi(h')$ . Then

$$\langle Q^{\Pi}(h,a)\rangle_{B} \stackrel{\text{(16)}}{\equiv} \sum_{h\in\mathcal{H}} B(h|sa)Q^{\Pi}(h,a)$$

$$\stackrel{\text{(2)}}{\equiv} \sum_{h\in\mathcal{H}} B(h|sa) \sum_{o'r'} P(o'r'|ha)[r' + \gamma V^{\Pi}(h')]$$

$$\stackrel{\text{(a)}}{\leq} \sum_{h\in\mathcal{H}} B(h|sa) \sum_{o'r'} P(o'r'|ha)[r' + \gamma (v^{\pi}(s') \pm \delta)]$$

$$\stackrel{Lem.3}{\equiv} \sum_{s'r'} p(s'r'|sa)[r' + \gamma v^{\pi}(s')] \pm \gamma \delta$$

$$\stackrel{\text{(7)}}{\equiv} q^{\pi}(s,a) \pm \gamma \delta$$

In (a) we used the assumption (i) of the Lemma. The derived upper and lower bounds imply  $|q^{\pi}(s,a) - \langle Q^{\Pi}(h,a) \rangle_B| \leq \gamma \delta$  (for all  $s = \phi(h)$  and a).

(ii) follows the same steps except with  $\Pi$  and  $\pi$  replaced by  $\Pi^*$  and  $\pi^*$ , and using (3) and (8) instead of (2) and (7) to justify the steps. Note that in general  $\Pi^* \neq \pi^*$ !

# 5 Approximate Aggregation Results

This section contains the main technical contribution of the paper. We show that histories (or raw states) can be aggregated and modeled by an MDP even if the true aggregated process is actually not an MDP. A necessary condition for successful aggregation is of course that the quantities of interest, namely (Q-)Value functions and Policies can be represented as functions of the aggregated states. The results in this section roughly show that this necessary condition, which is significantly weaker than the MDP requirement, is also sufficient. All but one result also holds for approximate aggregation, i.e. approximate conditions lead to approximate reductions. We also lift the stationarity assumption.

- Theorem 5 shows how (approximately)  $\phi$ -uniform  $Q^{\Pi}$  and  $\Pi$  can be obtained from the reduced Bellman equations (7).
- Theorem 6 weakens the assumptions and conclusions to (approximately)  $\phi$ -uniform  $V^{\Pi}$  and  $\Pi$ .

- Theorem 7 shows that for (approximately)  $\phi$ -uniform  $Q^*$ , the optimal policy is (approximately)  $\phi$ -uniform, and (an approximation of it) can be obtained via the reduced Bellman optimality equations (8).
- Theorem 9 shows that for (approximately)  $\phi$ -uniform  $V^*$  and  $\Pi^*$  we can obtain similar but somewhat weaker results. The proof of the latter involves extra complications not present in the other three proofs. Indeed, whether the arguably most desirable bound holds is Open Problem 10.

Note that all theorems crucially differ in their conditions and conclusions.

**Theorem 5**  $(\phi Q\pi)$  For any P,  $\phi$ , and B, define p via (15) and (5). Let  $\Pi$  be some policy such that  $\Pi(h) = \Pi(\tilde{h})$  and  $|Q^{\Pi}(h,a) - Q^{\Pi}(\tilde{h},a)| \le \varepsilon$  for all  $\phi(h) = \phi(\tilde{h})$  and all a. Then for all a and h it holds:

$$\begin{split} |Q^{\Pi}(h,a) - q^{\pi}(s,a)| &\leq \frac{\varepsilon}{1-\gamma} \quad and \quad |V^{\Pi}(h) - v^{\pi}(s)| \leq \frac{\varepsilon}{1-\gamma}, \\ where \quad \pi(s) := \Pi(h) \quad and \quad s = \phi(h) \end{split}$$

**Proof.** Let  $\delta := \sup_{s=\phi(h),a} |q^{\pi}(s,a) - Q^{\Pi}(h,a)|$ . Then  $|v^{\pi}(s) - V^{\Pi}(h)| \le \delta \ \forall s = \phi(h)$  by (10),

hence 
$$|q^{\pi}(s, a) - \langle Q^{\Pi}(h, a) \rangle_B| \leq \gamma \delta \ \forall s = \phi(h), a$$

by Lemma 4i. By assumption on  $Q^{\Pi}$  and B, for  $s = \phi(h)$  we have

$$\langle Q^{\Pi}(h,a)\rangle_{B} \equiv \sum_{\tilde{h}\in\mathcal{H}:\phi(\tilde{h})=s} B(\tilde{h}|sa)Q^{\Pi}(\tilde{h},a) \leq \sum_{\tilde{h}\in\mathcal{H}:\phi(\tilde{h})=s} B(\tilde{h}|sa)[Q^{\Pi}(h,a)\pm\varepsilon] = Q^{\Pi}(h,a)\pm\varepsilon$$

Together this implies  $|q^{\pi}(s,a) - Q^{\Pi}(h,a)| \leq \gamma \delta + \varepsilon$ , hence  $\delta \leq \gamma \delta + \varepsilon$ , hence  $\delta \leq \frac{\varepsilon}{1-\gamma}$ .

**Theorem 6**  $(\phi V\pi)$  For any P,  $\phi$ , and B, define p via (15) and (5). Let  $\Pi$  be some policy such that  $\Pi(h) = \Pi(\tilde{h})$  and  $|V^{\Pi}(h) - V^{\Pi}(\tilde{h})| \leq \varepsilon$  for all  $\phi(h) = \phi(\tilde{h})$ . Then for all a and h it holds:

$$|V^{\Pi}(h) - v^{\pi}(s)| \le \frac{\varepsilon}{1 - \gamma} \quad and \quad |q^{\pi}(s, a) - \langle Q^{\Pi}(h, a) \rangle_{B}| \le \frac{\varepsilon \gamma}{1 - \gamma}$$

$$where \quad \pi(s) := \Pi(h) \quad and \quad s = \phi(h)$$

**Proof.** Let  $\delta := \sup_{s=\phi(h),a} |v^{\pi}(s) - V^{\Pi}(h)|$ , fix some  $s = \phi(h)$ , and let  $a^{\pi} := \Pi(h)$ . Now

$$\langle Q^{\Pi}(h, a^{\pi}) \rangle_{B} \equiv \sum_{\tilde{h} \in \mathcal{H}: \phi(\tilde{h}) = s} B(\tilde{h} | s a^{\pi}) Q^{\Pi}(\tilde{h}, a^{\pi}) \stackrel{(a)}{=} \sum_{\tilde{h} \in \mathcal{H}: \phi(\tilde{h}) = s} B(\tilde{h} | s a^{\pi}) V^{\Pi}(\tilde{h})$$

$$\leq \sum_{\tilde{h} \in \mathcal{H}: \phi(\tilde{h}) = s} B(\tilde{h} | s a^{\pi}) [V^{\Pi}(h) \pm \varepsilon] = V^{\Pi}(h) \pm \varepsilon$$
(17)

where (a) follows from  $a^{\pi} = \Pi(h) = \Pi(\tilde{h})$  and  $Q^{\Pi}(\tilde{h}, \Pi(\tilde{h})) = V^{\Pi}(\tilde{h})$ . By Lemma 4i we have

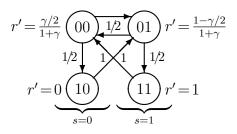
$$|q^{\pi}(s,a) - \langle Q^{\Pi}(h,a) \rangle_B| \le \gamma \delta \quad \forall s = \phi(h), a \tag{18}$$

We also have  $q^{\pi}(s,a^{\pi}) = q^{\pi}(s,\pi(s)) = v^{\pi}(s)$  from (7). Together with (17) and (18) for  $a = a^{\pi}$  this yields

$$|v^\pi(s) - V^\Pi(h)| \ \leq \ |v^\pi(s) - \langle Q^\Pi(h,a^\pi) \rangle_B| + |\langle Q^\Pi(h,a^\pi) \rangle_B - V^\Pi(h)| \ \leq \ \gamma \delta + \varepsilon$$

hence  $\delta \leq \gamma \delta + \varepsilon$  by the definition of  $\delta$ , hence  $\delta \leq \frac{\varepsilon}{1-\gamma}$ . Note that while  $|q^{\pi}(s,a^{\pi}) - Q^{\Pi}(h,a^{\pi})| \leq \frac{\varepsilon \gamma}{1-\gamma}$ , in general  $|q^{\pi}(s,a) - Q^{\Pi}(h,a)| \not \leq \frac{\varepsilon \gamma}{1-\gamma}$  for  $a \neq a^{\pi}$ .

**Example.** Consider a process P which itself is an MDP in the observations with transition matrix T and reward function R, i.e.  $P(o'r'|ha) = T^a_{oo'}R^{ar'}_{oo'}$ . The example on the right has the special form  $P(o'r'|ha) = T_{oo'} \cdot \llbracket r' = R(o) \rrbracket$ . It is an action-independent Markov process T with deterministic reward function R, which can be read off from the diagram. Observation space is  $\mathcal{O} = \{00,01,10,11\}$ . Consider reduction



$$s_t := \phi(h_t) := \begin{cases} 0 & \text{if } o_t = 00 \text{ or } 10 \\ 1 & \text{if } o_t = 01 \text{ or } 11 \end{cases} \in \mathcal{S} := \{0, 1\}$$

The reduced process  $P_{\phi}$  is not (even approximately) Markov:

$$P_{\phi}(s'=0|o=00) = T_{00,00} + T_{00,10} = 0 + 1/2 = 1/2$$
  
 $P_{\phi}(s'=0|o=10) = T_{10,00} + T_{10,10} = 0 + 0 = 0$ 

That is, P violates the bisimulation condition [GDG03], and raw states 00 and 10 have a large bisimulation distance [FPP04, Ort07]. On the other hand, the (Q-)Value function  $V(o_t) := V^{\pi}(h_t) = Q^{\pi}(h_t, a_t) \forall a_t$  can easily be verified to be

$$V(00) = V(10) = \frac{\gamma}{1 - \gamma^2}$$
 and  $V(01) = V(11) = \frac{1}{1 - \gamma^2}$ 

That is, V and Q are  $\phi$ -uniform. The conditions of Theorems 5 and 6 are satisfied exactly ( $\varepsilon = 0$ ), and hence the four raw states  $\mathcal{O}$  can be aggregated into two states  $\mathcal{S}$  despite  $P_{\phi} \notin \text{MDP}$  (the policy is irrelevant and can be chosen constant).  $\diamondsuit$ 

We now turn from the fixed policy case to similar theorems for optimal policies.

**Theorem 7**  $(\phi Q*)$  For any P,  $\phi$ , and B, define p via (15) and (5). Assume  $|Q^*(h,a)-Q^*(\tilde{h},a)| \leq \varepsilon$  for all  $\phi(h)=\phi(\tilde{h})$  and all a. Then for all a and h and  $s=\phi(h)$  it holds:

$$(i) \quad |Q^*(h,a) - q^*(s,a)| \leq \frac{\varepsilon}{1-\gamma} \quad and \quad |V^*(h) - v^*(s)| \leq \frac{\varepsilon}{1-\gamma},$$

(ii) 
$$0 \leq V^*(h) - V^{\tilde{\Pi}}(h) \leq \frac{2\varepsilon}{(1-\gamma)^2}, \quad where \quad \tilde{\Pi}(h) := \pi^*(s)$$

(iii) If 
$$\varepsilon = 0$$
 then  $\Pi^*(h) = \pi^*(s)$ 

**Proof.** (i) The proof follows the same steps as the proof of Theorem 5, replacing all  $\Pi$  and  $\pi$  by \* and using (12) instead of (10) and Lemma 4ii instead of Lemma 4i to justify the steps.

(iii) If  $\varepsilon = 0$ , then  $Q^*(h,a) = q^*(s,a)$  by (i) implies  $\Pi^*(h) = \pi^*(s)$ , where it is worthwhile to carefully check that the latter has actually not been used inadvertently in proving the former. Cf. the next theorem and proof.

(ii) For  $s = \phi(h)$  and  $\tilde{a} := \tilde{\Pi}(h) = \pi^*(s)$ ,

$$V^*(h) - \frac{\varepsilon}{1 - \gamma} \overset{(i)}{\leq} v^*(s) \overset{(8)}{=} q^*(s, \tilde{a}) \overset{(i)}{\leq} Q^*(h, \tilde{a}) + \frac{\varepsilon}{1 - \gamma}$$

which implies  $Q^*(h,\tilde{\Pi}(h)) \geq V^*(h) - \frac{2\varepsilon}{1-\gamma}$ . The claim now follows from the next Lemma 8 below.

The following lemma shows that if replacing the first action after h of the optimal policy  $\Pi^*$  by the action provided by  $\Pi$  thereafter following  $\Pi^*$  is at most  $\varepsilon$ -suboptimal, then always using  $\Pi$  is at most  $\frac{\varepsilon}{1-\gamma}$ -suboptimal.

**Lemma 8** ( $Q\pi*$ ) If  $Q^*(h,\Pi(h)) \ge V^*(h) - \varepsilon$  for all h for some policy  $\Pi$ , then for all h and a

$$0 \leq Q^*(h,a) - Q^{\Pi}(h,a) \leq \frac{\varepsilon \gamma}{1-\gamma} \quad and \quad 0 \leq V^*(h) - V^{\Pi}(h) \leq \frac{\varepsilon}{1-\gamma}$$

**Proof.** Let  $\delta := \sup_{h,a} [Q^*(h,a) - Q^{\Pi}(h,a)]$ . This implies

$$0 \stackrel{(a)}{\leq} V^*(h) - V^{\Pi}(h) \stackrel{(b)}{\leq} \varepsilon + Q^*(h, \Pi(h)) - Q^{\Pi}(h, \Pi(h)) \stackrel{(c)}{\leq} \varepsilon + \delta$$
 (19)

(a) follows from (1); (b) by assumption; and (c) by definition of  $\delta$  for  $a = \Pi(h)$ . Now for any a and h, this implies

$$Q^{\Pi}(h,a) \stackrel{(1)}{\leq} Q^*(h,a) \stackrel{(3)}{=} \sum_{o'r'} P(o'r'|ha)[r' + \gamma V^*(h')] \qquad [h' = hao'r']$$

$$\stackrel{(19)}{\leq} \sum_{o'r'} P(o'r'|ha)[r' + \gamma(V^{\Pi}(h') + \varepsilon + \delta)] \stackrel{(2)}{=} Q^{\Pi}(h, a) + \gamma(\varepsilon + \delta)$$

Hence  $\delta \leq \gamma(\varepsilon + \delta)$ , hence  $\delta \leq \frac{\varepsilon \gamma}{1 - \gamma}$ .

**Theorem 9**  $(\phi V*)$  For any P,  $\phi$ , and B, define p via (15) and (5). Assume  $\Pi^*(h) = \Pi^*(\tilde{h})$  and  $|V^*(h) - V^*(\tilde{h})| \le \varepsilon$  for all  $\phi(h) = \phi(\tilde{h})$ . Then for all a and h and  $s = \phi(h)$  it holds:

(i) 
$$|V^*(h) - v^*(s)| \le \frac{3\varepsilon}{(1-\gamma)^2}$$
 and  $|q^*(s,a) - \langle Q^*(h,a) \rangle_B| \le \frac{3\varepsilon\gamma}{(1-\gamma)^2}$ ,

(ii) If 
$$\varepsilon = 0$$
 then  $\Pi^*(h) = \pi^*(s)$ 

The proof actually implies the stronger lower bound  $V^*(h) - v^*(s) \ge \frac{3\varepsilon}{1-\gamma}$  and similarly for  $Q^*$ , but we do not know whether the upper bound can be improved.

**Proof.** While proofs start to get routine, here is a warning that care is in order when recycling similar proofs. Theorem 6 relies on the assumption that  $\pi(s) = \Pi(h)$  for  $s = \phi(h)$ , while we were lucky that the proof of Theorem 7 worked without knowing  $\pi^*(s) = \Pi^*(h)$  in advance. Here we have to work a bit harder.

Let us define  $a^0 := \pi^0(s) := \Pi^*(h)$  for  $s = \phi(h)$ . The Bellman equation for policy  $\pi^0$  is

$$q^{\pi^0}(s,a) = \sum_{s'r'} p(s'r'|sa)[r' + \gamma v^{\pi^0}(s')] \quad \text{and} \quad v^{\pi^0}(s) = q^{\pi^0}(s,\pi^0(s))$$
 (20)

At this stage  $\pi^0$  may well be different from  $\pi^*$ , since  $\pi^*$  satisfies (8), not (20), but we will now show that it actually does. First note that

$$q^{\pi^0}(s, a^0) = v^{\pi^0}(s) \le V^{\Pi^*}(h) \pm \frac{\varepsilon}{1 - \gamma} = V^*(h) \pm \frac{\varepsilon}{1 - \gamma}$$
 (21)

where the bounds follow from Theorem 6 applied to  $\Pi := \Pi^*$  (with  $\pi = \pi^0$ ). For general a we only get an upper bound:

$$q^{\pi^{0}}(s,a) - \frac{\varepsilon\gamma}{1-\gamma} \stackrel{Thm.6}{\leq} \langle Q^{\Pi^{*}}(h,a)\rangle_{B} \stackrel{(16)}{=} \sum_{h\in\mathcal{H}} B(h|sa)Q^{*}(h,a) \tag{22}$$

$$\stackrel{(4)}{\leq} \sum_{h\in\mathcal{H}} B(h|sa)Q^{*}(h,\Pi^{*}(h)) \stackrel{(14)}{=} \sum_{\tilde{h}\in\mathcal{H}:\phi(\tilde{h})=s} B(\tilde{h}|sa)V^{*}(\tilde{h})$$

$$\stackrel{(a)}{\leq} \sum_{\tilde{h}\in\mathcal{H}:\phi(\tilde{h})=s} B(\tilde{h}|sa)[V^{*}(h)+\varepsilon] \stackrel{(14)}{=} V^{*}(h)+\varepsilon$$

(a) uses the theorem's assumption on  $V^*(h)$ . Together, (21) and (22) imply

$$v^{\pi^0}(s) \stackrel{(20)}{=} q^{\pi^0}(s, a^0) \le \max_{a} q^{\pi^0}(s, a) \stackrel{(22)}{\le} V^*(h) + \frac{\varepsilon}{1 - \gamma} \stackrel{(21)}{\le} v^{\pi^0}(s) + \frac{2\varepsilon}{1 - \gamma}$$
 (23)

(ii) For  $\varepsilon = 0$ , the previous equation implies  $v^{\pi^0}(s) = \max_a q^{\pi^0}(s,a)$ , hence (20) can be rewritten as

$$q^{\pi^0}(s, a) = \sum_{s'r'} p(s'r'|sa)[r' + \gamma v^{\pi^0}(s')]$$
 and  $v^{\pi^0}(s) = \max_a q^{\pi^0}(s, a)$ 

This shows that  $(q^{\pi^0}, v^{\pi^0})$  satisfies the same Bellman *optimality* equation as  $(q^*, v^*)$  does. Since it has a unique solution, we must have  $q^{\pi^0} \equiv q^*$  and  $v^{\pi^0} \equiv v^*$  and  $\pi^* \equiv \pi^0$ , which for  $s = \phi(h)$  implies  $\Pi^*(h) = \pi^*(s)$  by definition of  $\pi^0$ . It also implies  $V^*(h) = v^*(s)$  by (21), and  $q^*(s,a) = \langle Q^*(h,a) \rangle_B$  by Lemma 4ii, i.e. the  $\varepsilon = 0$  version of (i).

(i) We now continue with the general  $\varepsilon > 0$  case. For all s and a we have

$$0 \overset{(9)}{\leq} q^*(s,a) - q^{\pi^0}(s,a) \overset{(7)}{\underset{(8)}{=}} \sum_{s'r'} p(s'r'|sa) \gamma(v^*(s') - v^{\pi^0}(s')) \overset{(a)}{\leq} \gamma \max_{s'} \{v^*(s') - v^{\pi^0}(s')\}$$

$$0 \overset{(9)}{\leq} v^*(s) - v^{\pi^0}(s) \overset{(23)}{\underset{(8)}{\leq}} \max_a q^*(s,a) - \max_a q^{\pi^0}(s,a) + \frac{2\varepsilon}{1-\gamma} \overset{(13)}{\leq} \max_a \{q^*(s,a) - q^{\pi^0}(s,a)\} + \frac{2\varepsilon}{1-\gamma} \overset{(13)}{\leq} \underset{(13)}{\leq} \underset{($$

In (a) we have upper bounded the p-expectation by the maximum. Together this gives

$$\max_{s} \{v^{*}(s) - v^{\pi^{0}}(s)\} \leq \gamma \max_{s} \{v^{*}(s) - v^{\pi^{0}}(s)\} + \frac{2\varepsilon}{1 - \gamma}$$

$$\Rightarrow \max_{s} \{v^{*}(s) - v^{\pi^{0}}(s)\} \leq \frac{2\varepsilon}{(1 - \gamma)^{2}}$$
Hence for  $s = \phi(h)$ :  $V^{*}(h) - \frac{\varepsilon}{1 - \gamma} \stackrel{(21)}{\leq} v^{\pi^{0}}(s) \stackrel{(9)}{\leq} v^{*}(s) \stackrel{\sim}{\leq} v^{\pi^{0}}(s) + \frac{2\varepsilon}{(1 - \gamma)^{2}}$ 

$$\stackrel{(21)}{\leq} V^{*}(h) + \frac{\varepsilon}{1 - \gamma} + \frac{2\varepsilon}{(1 - \gamma)^{2}} \leq V^{*}(h) + \frac{3\varepsilon}{(1 - \gamma)^{2}}$$

Together with Lemma 4ii this implies (i).

We are primarily interested in the optimal policy  $\Pi^*(h)$ ; to correctly represent the value  $V^*(h)$  is only of indirect interest. If  $\Pi^*$  is  $\phi$ -uniform, it can be represented as  $\Pi^*(h) = \pi^0(s)$  for some  $\pi^0$ , but if the  $\phi$ -uniformity condition on  $V^*$  in Theorem 9 is dropped, the conclusion  $\Pi^*(h) = \pi^*(s)$  can fail as the following example shows.

Counter Example. Let P be the MDP  $P(o'r'|ha) := T^a_{oo'} \cdot \llbracket r' = R^a_o \rrbracket$  with two raw states  $o \in \{0,1\}$  and two actions  $a \in \{\alpha,\beta\}$  formally defined on the left and depicted on the right:

$$T^{\alpha} := \begin{pmatrix} 1 & 0 \\ 1 & 0 \end{pmatrix}, \qquad R^{\alpha} := \begin{pmatrix} 1/6 \\ 1 \end{pmatrix}, \qquad \qquad \beta, \ r' = 0, \ p = \frac{1}{2}, \qquad \beta, \ r' = 1, \qquad \beta, \ r' = 1, \qquad \beta, \quad r' = 1/2, \qquad \beta, \qquad r' = 1/2$$

The value of policy  $\pi$  in vector notation is  $V^{\pi} = R^{\pi} + \gamma T^{\pi} V^{\pi}$ , where  $V_{o_t}^{\pi} := V^{\pi}(h_t)$ . The 4 stationary policies are denoted by  $\pi = a^0 a^1$ , where  $a^o$  is the action taken in raw state o. For  $\gamma = 0$ , their values are

Policy  $\pi = \alpha \alpha$  has the highest value, therefore  $\Pi^*(h) \equiv \alpha$ . Let us now aggregate raw states  $o \in \{0,1\}$  to a 1-state MDP. Its value is  $v = \frac{1}{1-\gamma} \rho^{\mathsf{T}} R = \rho^{\mathsf{T}} R$  for  $\gamma = 0$ , where  $\rho$  is the stationary distribution  $\rho^{\mathsf{T}} = \rho^{\mathsf{T}} T$  of T, in particular  $\rho^{\alpha} = \binom{1}{0}$  for  $T^{\alpha}$  and  $\rho^{\beta} = \binom{1/2}{1/2}$  for  $T^{\beta}$ . Since there is only 1 aggregated state, there are only 2 stationary policies, one for each action. This leads to  $v^{\alpha} = \frac{1}{6} < \frac{1}{4} = v^{\beta}$ , hence  $\pi^*(s) \equiv \beta \neq \alpha \equiv \Pi^*(h) \ \forall s,h$ . That is, despite  $\Pi^*$  being constant,  $\pi^* \neq \Pi^*$ , which shows that the condition on  $V^*$  in Theorem 9 cannot be dropped. Note that  $V^* = V^{\alpha\alpha} = \binom{1/6}{1}$  is far from constant. By continuity, the policy reversal also holds for  $\gamma > 0$ . Indeed, this example works for all  $\gamma < \frac{2}{5}$  and other examples work for all  $0 \le \gamma < 1$ .

Open Problem 10 ( $\phi V*$ ) Under the same conditions as Theorem 9, is

$$V^*(h) - V^{\tilde{\Pi}}(h) \stackrel{??}{=} O\left(\frac{\varepsilon}{(1-\gamma)^?}\right) \quad where \quad \tilde{\Pi}(h) := \pi^*(s)$$
 (25)

**Arguments.** Here are some arguments why it might be true (or false):

- (1) For  $\varepsilon = 0$  it immediately follows from Theorem 9, since in this case  $\Pi(h) = \Pi^*(h)$ . Some continuity argument might allow to establish a bound for small  $\varepsilon > 0$ .
- (2) Theorem 7i&iii mostly carried over to Theorem 9, so a-priori it is not too implausible that Theorem 7ii carries over to (25). On the other hand, the proofs of (i) and (iii) of both theorems were sufficiently different, so the analogy argument is weak.
  - (3) Let  $\tilde{a} := \tilde{\Pi}(h) := \pi^*(s)$  for  $s = \phi(h)$ . Then

$$\langle Q^*(h,\tilde{a})\rangle_B \overset{Thm.9i}{\geq} q^*(s,\tilde{a}) - \frac{3\varepsilon\gamma}{(1-\gamma)^2} \overset{(8)}{=} v^*(s) - \frac{3\varepsilon\gamma}{(1-\gamma)^2} \overset{(24)}{\geq} V^*(h) - \frac{3\varepsilon}{(1-\gamma)^2}$$

$$Q^*(h,\tilde{a}) \overset{(3)}{\leq} V^*(h)$$

For  $\varepsilon = 0$  this pair of inequalities implies that  $Q^*(h,\tilde{a})$  lower bounds its own expectation, therefore it must be constant and equal to  $V^*(h)$  on each  $\phi^{-1}(s)$ -partition. For  $\varepsilon > 0$ , with high probability  $Q^*(h,\tilde{a})$  cannot be much smaller than  $V^*(h)$ . If it weren't for the probability qualifier we could now apply Lemma 8 to establish (25) (as in the proof of Theorem 7ii). Low probability events could invalidate this argument.

**Discussion.** Open Problem 10 would be the main result if we had a proof for  $\varepsilon > 0$ . Absent of it we have to be content with Theorem 7ii. Both statements imply that we can aggregate histories as much as we wish, as long as the optimal value function and policy are still approximately representable as functions of aggregated states. Whether the reduced process  $P_{\phi}$  is Markov or not is immaterial. We can use surrogate MDP p to find an  $\varepsilon$ -optimal policy for P.

Most RL work, including on state aggregation, is formulated in terms of MDPs, i.e. the original process P is already an MDP. Let us call this the original or raw MDP. We could interpret the whole history as a raw state, which formally makes every P an MDP, but normally only observations are identified with raw states, i.e. P is a raw MDP iff P(o'r'|ha) = P(o'r'|oa). In this case,  $V^*(h_t) = V^*(o_t)$  etc. depends on raw states only (which is well known or follows from Theorem 2 with  $\phi(h_t) = o_t$ ). Since our results hold for all P, they clearly hold if P is a raw MDP and if  $\phi(h_t) := \phi(o_t)$  maps raw states to aggregated states.

The remainder of this paper shows how much we can aggregate and how to develop RL algorithms exploiting these insights.

# 6 Extreme Aggregation

The results of Section 5 showed that histories can be aggregated and modeled by an MDP even if the true aggregated process is not an MDP. The only restrictions were that the (Q-)Value functions and Policies could still be (approximately) represented as functions of the aggregated states. We will see in this section that in theory this allows to represent any process P as a small finite-state MDP.

Extreme aggregation based on Theorem 7. Consider  $\phi$  that maps each history to the vector-over-actions of optimal Q-values  $Q^*(h,\cdot)$  discretized to some finite  $\varepsilon$ -grid:

$$\phi(h) := \left( \lfloor Q^*(h, a) / \varepsilon \rfloor \right)_{a \in \mathcal{A}} \in \{0, 1, ..., \lfloor \frac{1}{\varepsilon(1 - \gamma)} \rfloor \}^{\mathcal{A}} =: \mathcal{S}$$
 (26)

That is, all histories with  $\varepsilon$ -close  $Q^*$ -values are mapped to the same state:

$$|Q^*(h,a) - Q^*(\tilde{h},a)| \le \varepsilon \quad \forall \phi(h) = \phi(\tilde{h}) \ \forall a$$

Now choose some B and determine p from P via (15) and (5). Find the optimal policy  $\pi^*$  of MDP p of size |S|. Define  $\tilde{\Pi}(h) := \pi^*(\phi(h))$ . By Theorem 7ii,  $\tilde{\Pi}$  is an  $\varepsilon'$ -optimal policy of original process P in the sense that

$$|V^{\tilde{\Pi}}(h) - V^*(h)| \le \frac{2\varepsilon}{(1-\gamma)^2} =: \varepsilon'$$

Extreme aggregation based on Open Problem 10. If (25) holds, we can aggregate even better: Consider  $\phi$  that maps each history to the optimal Value  $V^*(h)$  discretized to some finite  $\varepsilon$ -grid and to the optimal action  $\Pi^*(h)$ :

$$\phi(h) := \left( \lfloor V^*(h)/\varepsilon \rfloor, \Pi^*(h) \right) \in \left\{ 0, 1, ..., \lfloor \frac{1}{\varepsilon(1-\gamma)} \rfloor \right\} \times \mathcal{A} =: \mathcal{S}$$
 (27)

That is, all histories with  $\varepsilon$ -close  $V^*$ -Values and same optimal action are mapped to the same state:

$$|V^*(h) - V^*(\tilde{h})| \ \leq \ \varepsilon \quad \text{and} \quad \Pi^*(h) = \Pi^*(\tilde{h}) \qquad \forall \phi(h) = \phi(\tilde{h})$$

As before, determine p, find its optimal policy  $\pi^*$ , and define  $\tilde{\Pi}(h) := \pi^*(\phi(h))$ . If (25) holds, then  $\tilde{\Pi}$  is an  $\varepsilon'$ -optimal policy of original process P in the sense that

$$|V^{\tilde{\Pi}(h)} - V^*(h)| = O\left(\frac{\varepsilon}{(1-\gamma)^?}\right) =: \varepsilon'$$

The following theorem summarizes the considerations for the two choices of  $\phi$  above:

**Theorem 11 (Extreme \phi)** For every process P there exists a reduction  $\phi$  ((26) or (27) will do) and MDP p defined via (15) and (5) whose optimal policy  $\pi^*$  is an  $\varepsilon'$ -optimal policy  $\tilde{\Pi}(h) := \pi^*(\phi(h))$  for P. The size of the MDP is bounded (uniformly for any P) by

$$|\mathcal{S}| \le \left(\frac{3}{\varepsilon'(1-\gamma)^3}\right)^{|\mathcal{A}|}$$
 and if (25) holds even by  $|\mathcal{S}| = O\left(\frac{|\mathcal{A}|}{\varepsilon'(1-\gamma)^{1+?}}\right)$ 

**Proof.** For S defined in (26) we have

$$|\mathcal{S}| = \left( \left\lfloor \frac{1}{\varepsilon(1-\gamma)} \right\rfloor + 1 \right)^{|\mathcal{A}|} = \left( \left\lfloor \frac{2}{\varepsilon'(1-\gamma)^3} \right\rfloor + 1 \right)^{|\mathcal{A}|} \le \left( \frac{3}{\varepsilon'(1-\gamma)^3} \right)^{|\mathcal{A}|}$$

where in the last inequality we have assumed  $\varepsilon' \leq \frac{1}{1-\gamma}$ . (For  $\varepsilon' > \frac{1}{1-\gamma}$  the theorem is trivial, since any policy is  $\varepsilon'$ -optimal). For  $\mathcal{S}$  defined in (27) the derivation is similar. The theorem now follows from the considerations in the paragraphs before the theorem.

**Discussion.** A valid question is of course whether Theorem 11 is just an interesting theoretical insight/curiosity or of any practical use. After all,  $\phi$  depends on  $Q^*$  (or  $V^*$  and  $\Pi^*$ ), but if we knew  $Q^*$ ,  $\Pi^*$  would readily be available and the detour through p and  $\pi^*$  pointless.

Theorem 11 reaches relevance by the following observation: If we start with a sufficiently rich class of maps  $\Phi$  that contains at least one  $\phi$  approximately representing  $Q^*(h,\cdot)$ , and have a learning algorithm that favors such  $\phi$ , then Theorems 5–9 tell us that we do not need to worry about whether  $P_{\phi}$  is MDP or not; we "simply" use/learn MDP p instead. Theorem 11 tells us that this allows for extreme aggregation far beyond MDPs.

This program is in parts worked out in the next two sections, but more research is needed for its completion. Learning p from (real) P-samples is considered in Section 7 and learning  $\phi$  in Section 8.

## 7 Reinforcement Learning

In RL, P and therefore p are unknown. We now show how to learn p from samples from P. For this we have to link B to the distribution over histories induced by P and to the behavior policy  $\Pi_B$  the agent follows. We still assume  $\phi$  is given.

Behavior policy  $\Pi_B$ . Let  $\Pi_B: \mathcal{H} \leadsto \mathcal{A}$  be the behavior policy of our RL agent, which in general is non-stationary due to learning, often stochastic to ensure exploration, and (usually) different from any policy considered so far  $(\Pi^*, \pi^*, \tilde{\Pi}, \pi^0, \Pi, \pi)$ . Note that a sequence of policies  $\Pi_1, \Pi_2, ...$  where each  $\Pi_t$  is learnt from  $h_t$  and used at time t (or for some number of steps) is nothing but a single non-stationary policy  $\Pi_B(h_t) = \Pi_t(h_t) \forall t, h_t$ , so  $\Pi_B$  indeed includes the case of policy learning.

Choice of B. The interaction of agent  $\Pi_B$  with environment P stochastically generates some history  $h_t$  followed by action  $a_t$  with joint probability, say  $P_B(h_t a_t)$ . We use subscripts B and/or  $\phi$  to indicate dependence on  $\Pi_B$  and/or  $\phi$ . A natural choice for B(h|sa) in (14) would be to condition of  $P_B$  on  $s_t a_t$ . We now show that this does not work and how to fix the problem. We can get  $P_{\phi B}(h_t|s_t a_t)$  from P and  $\Pi_B$  and several other useful distributions as follows:

$$P_{B}(h_{t+1}|h_{t}) = P(o_{t+1}r_{t+1}|h_{t}a_{t})\Pi_{B}(a_{t}|h_{t}) \quad [h_{t+1} = h_{t}a_{t}o_{t+1}r_{t+1}]$$

$$P_{B}(h_{n}) = \prod_{t=0}^{n-1} P_{B}(h_{t+1}|h_{t}), \quad P_{B}(h_{t}a_{t}) = \Pi_{B}(a_{t}|h_{t})P_{B}(h_{t})$$

$$P_{\phi B}(s_{t}a_{t}) = \sum_{h_{t}:\phi(h_{t})=s_{t}} P_{B}(h_{t}a_{t}), \quad P_{\phi B}(h_{t}|s_{t}a_{t}) = \frac{P_{B}(h_{t}a_{t})}{P_{\phi B}(s_{t}a_{t})} \llbracket \phi(h_{t}) = s_{t} \rrbracket$$

$$P_{\phi B}(s_{t+1}r_{t+1}|s_{t}a_{t}) = \sum_{h_{t}:\phi(h_{t})=s_{t}} P_{\phi}(s_{t+1}r_{t+1}|h_{t}a_{t})P_{\phi B}(h_{t}|s_{t}a_{t}) \quad [\text{see (5) for def. of } P_{\phi}] (28)$$

$$P_{\phi B}(s_{t}a_{t}s_{t+1}r_{t+1}) = P_{\phi B}(s_{t+1}r_{t+1}|s_{t}a_{t})P_{\phi B}(s_{t}a_{t})$$

 $P_{\phi B}(h_t|s_t a_t)$  has the following properties:

$$P_{\phi B}(h_t|s_t a_t) \ge 0$$
 and  $\sum_{h_t \in \mathcal{H}_t} P_{\phi B}(h_t|s_t a_t) = \sum_{h_t: \phi(h_t) = s_t} P_{\phi B}(h_t|s_t a_t) = 1 \quad \forall t, s_t, a_t \quad (29)$ 

This is close to the required condition (14) for B but crucially different. The sum in (14) is over histories of all lengths while in (29) the sum is limited to histories of length t. It is easy to miss this difference due to the compact notation. Technically  $P_B$  is a probability measure on infinite sequences  $\mathcal{H}_{\infty}$  and  $P_B(h_t)$  is short for  $P_B(\Gamma_{h_t})$  where  $\Gamma_{h_t}$  is the set of infinite histories starting with  $h_t$ , i.e.  $P_B(h_t)$  is the probability that the infinite history starts with  $h_t$  ( $\sum_{h_t \in \mathcal{H}_t} P_B(h_t) = 1 \forall t$ ). On the other hand, B(h) is a probability distribution over finite histories of mixed length ( $\sum_{h \in \mathcal{H}} B(h) = 1$ ); similarly for  $P_B$  and B conditioned on / parameterized by s and a.

We can fix this mismatch by introducing weights  $w_t: \mathcal{S} \times \mathcal{A} \leadsto [0;1]$  and define

$$B(h_t|sa) := w_t(sa)P_{\phi B}(h_t|s_t = s, a_t = a) \ \forall t, \quad \text{where} \quad \sum_{t=1}^{\infty} w_t(sa) = 1 \ \forall s, a \ (30)$$

which now satisfies (14) (due to  $\sum_{h\in\mathcal{H}} = \sum_{t=1}^{\infty} \sum_{h_t\in\mathcal{H}_t}$ ). MDP p can now be repre-

sented as

$$p(s'r'|sa) = \sum_{t=1}^{\infty} w_t(sa) \sum_{h_t \in \mathcal{H}_t} P_{\phi}(s_{t+1} = s', r_{t+1} = r'|h_t, a_t = a) P_{\phi B}(h_t|s_t = s, a_t = a)$$

$$= \sum_{t=1}^{\infty} w_t(sa) P_{\phi B}^t(s'r'|sa)$$
(31)

That is, p is the w-weighted time-average of  $P_{\phi B}^t$ . The first equality follows from (15) and (30); the second one from (28). We also introduced the shorthand  $P_{\phi B}^t(s'r'|sa) := P_{\phi B}(s_{t+1} = s', r_{t+1} = r'|s_t = s, a_t = a)$ .

Choice of  $w_t$ . If  $P_{\phi B}^t$  in (31) is stationary, i.e. independent of t, then  $p(s'r'|sa) = P_{\phi B}^t(s'r'|sa)$  for all t, since the weights sum to one, and estimation is easy. Note that in general we cannot estimate non-stationary  $P_{\phi B}^t$ , since for each t we have only one sample available, but we will see that estimation of p is still possible. Assume we have observed  $h_n$ , and choose

$$w_t(sa) := \frac{P_{\phi B}^t(sa)}{\sum_{t=1}^n P_{\phi B}^t(sa)} \quad \text{for } t \le n \quad \text{and} \quad 0 \quad \text{for } t > n$$
 (32)

Inserting this into (31) and using (28) gives

$$p(s'r'|sa) = \frac{\frac{1}{n} \sum_{t=1}^{n} P_{\phi B}^{t}(sas'r')}{\frac{1}{n} \sum_{t=1}^{n} P_{\phi B}^{t}(sa)}$$
(33)

We estimate numerator and denominator separately.

Law of large numbers. For t=1,2,3,... let  $X_t \in \{0,1\}$  be binary random variables with expectation  $\mathbb{E}[X_t]$ . Define  $n_1 = \sum_{t=1}^n X_t$  be the number of sampled 1s. The strong law of large numbers says that

$$\frac{n_1}{n} - \frac{1}{n} \sum_{t=1}^{n} \mathbb{E}[X_t] \stackrel{n \to \infty}{\longrightarrow} 0 \quad \text{almost surely} \quad \text{under weak conditions}$$
 (34)

Note that the law holds far beyond i.i.d. random variables under a variety of conditions [Faz06, VGS05] which we collectively call 'weak conditions'. It is not even necessary for  $n_1/n$  to converge.

Estimation of p. Now fix some (s,a), and let  $X_t := [s_t = s, a_t = a]$ . (Here we assume that variables in  $h_t$  are random variables and sas'r' are realizations.) Then

$$n(sa) := n_1 = \sum_{t=1}^{n} X_t = \#\{t \le n : s_t = s, a_t = a\}$$

is the number of times action a is taken in state s, and  $\mathbb{E}[X_t] = P(X_t = 1) = P_{\phi B}^t(sa)$ , hence (34) implies

$$\frac{n(sa)}{n} - \frac{1}{n} \sum_{t=1}^{n} P_{\phi B}^{t}(sa) \stackrel{n \to \infty}{\longrightarrow} 0 \quad \text{a.s. under weak conditions}$$
 (35)

Similarly for  $Y_t := [s_t a_t s_{t+1} r_{t+1} = sas'r']$  and  $n(sas'r') := \sum_{t=1}^n Y_t$  we have

$$\frac{n(sas'r')}{n} - \frac{1}{n} \sum_{t=1}^{n} P_{\phi B}^{t}(sas'r') \stackrel{n \to \infty}{\longrightarrow} 0 \quad \text{with } P\text{-probability } 1$$
 (36)

under weak conditions. (35) and (36) via (33) are nearly sufficient to imply

$$\frac{n(sas'r')}{n(sa)} - p(s'r'|sa) \xrightarrow{n \to \infty} 0 \text{ almost surely}$$
 (37)

A sufficient but by far not necessary condition is

$$\lim_{n \to \infty} \inf \frac{n(sa)}{n} > 0 \quad \text{almost surely}$$
(38)

**Theorem 12** (p-estimation) For B defined in (30) and (32) we have: If (36) and (38) hold, then (37) holds. For example, if  $Y_t$  are stationary ergodic processes, then (36) and (38) hence (37) hold for all state-action pairs that matter (i.e. for those occurring with non-zero probability).

**Proof.** We introduce the following (*n*-dependent) shorthands:

$$\bar{X} := \frac{n(sa)}{n}, \qquad \bar{x} := \frac{1}{n} \sum_{t=1}^{n} P_{\phi B}^{t}(sa), \quad \alpha := \liminf_{n \to \infty} \frac{n(sa)}{n},$$

$$\bar{Y} := \frac{n(sas'r')}{n}, \qquad \bar{y} := \frac{1}{n} \sum_{t=1}^{n} P_{\phi B}^{t}(sas'r')$$

With these abbreviations, assumption (36) implies (35), i.e.

$$\bar{Y} - \bar{y} \to 0 \text{ implies } \bar{X} - \bar{x} = \sum_{s'r'} \bar{Y} - \sum_{s'r'} \bar{y} = \sum_{s'r'} [\bar{Y} - \bar{y}] \to 0$$
 (39)

since  $\mathcal S$  and  $\mathcal R$  have been assumed finite. Now

$$\begin{split} \left|\frac{n(sas'r')}{n(sa)} - p(s'r'|sa)\right| &= \left|\frac{\bar{Y}}{\bar{X}} - \frac{\bar{y}}{\bar{x}}\right| \leq \left|\frac{\bar{Y}}{\bar{X}} - \frac{\bar{y}}{\bar{X}}\right| + \left|\frac{\bar{y}}{\bar{X}} - \frac{\bar{y}}{\bar{x}}\right| \\ &= \frac{1}{\bar{X}}|\bar{Y} - \bar{y}| + \frac{\bar{y}}{\bar{X}\bar{x}}|\bar{x} - \bar{X}| \leq \frac{1}{\bar{X}}\left(|\bar{Y} - \bar{y}| + |\bar{x} - \bar{X}|\right) \stackrel{n \to \infty}{\longrightarrow} 0 \quad \text{a.s.} \end{split}$$

The first inequality is just the triangle inequality. The second inequality follows from  $\bar{y} \leq \bar{x}$ . The limit is zero, since almost surely  $\limsup_{n\to\infty} [1/\bar{X}] = 1/\alpha < \infty$  and  $\bar{Y} - \bar{y} \to 0$  and  $\bar{X} - \bar{x} \to 0$ . Hence (37) holds. Finally, for stationary ergodic  $Y_t$ , we have  $\bar{y} = \frac{1}{n} \sum_{t=1}^{n} \mathbb{E}[Y_t] = \mathbb{E}[Y_1] = constant$ , and hence  $\bar{x} = \sum_{s'r'} \bar{y} = constant$ . Therefore

(36) holds by 
$$\bar{Y} = \frac{1}{n} \sum_{t=1}^{n} Y_t \xrightarrow{\text{ergodicity}} \mathbb{E}[Y_1] \stackrel{\text{stationarity}}{=} \frac{1}{n} \sum_{t=1}^{n} \mathbb{E}[Y_t] = \bar{y},$$

(38) holds by 
$$\liminf_{n\to\infty} \bar{X} \stackrel{(39)}{=} \bar{x} \stackrel{\text{stationarity}}{=} \mathbb{E}[X_1] = P_{\phi B}^1(sa) \stackrel{\text{assumption}}{>} 0$$

**Discussion.** Limit (37) shows that standard frequency estimation for p will converge to the true p under weak conditions. If  $P_{\phi}$  is MDP, samples are conditionally i.i.d. and the 'weak conditions' are satisfied. But the law of large numbers and hence (37) holds far beyond the i.i.d. case [FK01, VGS05], e.g. for stationary ergodic processes. Condition (38) that every state-action pair be visited with non-vanishing relative frequency can be significantly relaxed. Stationarity is also not necessary, and indeed often does not hold due to a non-stationary environment P or a non-stationary behavior policy  $\Pi_B$  (or both).

Other choices for  $w_t$  are possible, e.g. we could multiply numerator and denominator of (32) by some arbitrary positive function  $u_t(as)$ , which leads to a weighted average estimator.

We estimate p in order to estimate  $q^*$  and ultimately  $\pi^*$ . This is model-based RL. We can also learn  $\pi^*$  model-free. For instance, condition (37) should be sufficient for Q-learning to converge to  $Q^*$ .

Q-learning and other RL algorithms designed for MDPs have been observed to often (but not always) perform well even if applied to non-MDP domains. Our results appear to explain why, but this calls for further investigations.

# 8 Feature Reinforcement Learning

The idea of FRL is to  $learn \phi$  [Hut09c]. FRL starts with a class of maps  $\Phi$ , compares different  $\phi \in \Phi$ , and selects the most appropriate one given the experience  $h_t$  so far. Several criteria based on how well  $\phi$  reduces P to an MDP have been devised [Hut09b, Hut09a] and theoretically [SH10] and experimentally [NSH11] investigated [Ngu13]. Theorems 5–9 show that demanding  $P_{\phi}$  to be approximately MDP is overly restrictive. Theorem 11 suggests that if we relax this condition, much more substantial aggregation is possible, provided  $\Phi$  is rich enough.

(F)RL deals with the case of unknown P. We first discuss learning  $\phi$  for the unrealistic case of exact aggregation ( $\varepsilon = 0$ ) and infinite sample size ( $n = \infty$ ). This serves as a useful guide to work out its generalization to the realistic but significantly more complex case of approximate aggregation based on finite sample size. Finally we discuss a family of recent algorithms (BLB and extensions [Ngu13]) that appear to nearly have the right properties for our purpose. This section is more a collection of ideas and outlook towards algorithms exploiting and motivating the usefulness of the new insights obtained in the previous sections.

Search for exact  $\phi$  based on infinite sample size. Since we are now concerned with comparing different  $\phi \in \Phi$ , we subscribe quantities with  $\phi$  when necessary. Consider the unrealistic case of infinite sample size  $(n = \infty)$  and a search for exact reductions  $\phi$ . We call a reduction  $\phi: \mathcal{H} \to \mathcal{S}_{\phi}$  exact iff  $Q^*(h,a) = q_{\phi}^*(s,a)$  and  $\Pi^*(h) = \pi_{\phi}^*(s)$  for all  $s = \phi(h)$  and a.

Even for  $n=\infty$ , P hence  $Q^*$  needed for  $\Pi^*$  is (usually) not estimable (from  $h_\infty$ ). On the other hand, for each  $\phi \in \Phi$ ,  $p=p_\phi$  can be determined (exactly) by (37) (under weak conditions). From  $p_\phi$  we can determine  $q_\phi^*$  and  $\pi_\phi^*$  via (8) and (9). The solution always satisfies the reduced Bellman equations exactly, even for very bad reductions, e.g. single state  $\phi(h) \equiv 0 \,\forall h$ . So the reduced problem is not sufficient to judge the quality of  $\phi$ . An alternative to assuming  $n=\infty$  is to assume that P is known, which also allows to determine  $p_\phi$ , etc. So what follows applies to stochastic planning as well

Coarsening and refining reductions  $\phi$ : Let us now coarsen  $\phi$ , i.e. further merge some partitions  $\phi^{-1}(s)$ . In the simplest case we just merge two states into one. In general, consider coarsening  $\chi: \mathcal{S}_{\phi} \to \mathcal{S}_{\psi}$  and coarser reduction  $\psi: \mathcal{H} \to \mathcal{S}_{\psi}$  such that  $\psi(h) = \chi(\phi(h))$ . We also call  $\mathcal{S}_{\phi}$  a refinement of  $\mathcal{S}_{\psi}$ . For example, U-trees [McC96, UV98] and Kd-trees have been used in RL [EGW05], where expanding a leaf corresponds to splitting a state. Or in  $\phi$ DBN,  $S_{\phi} = \{0,1\}^d$  is a binary feature vector, where removing one component corresponds to pairwise combining  $2^d$  states to  $2^{d-1}$  states [Hut09a].

Ordering reductions in  $\Phi$ : We can partially order  $\Phi$  as follows:

$$\psi \prec \phi : \Leftrightarrow q_{\phi}^* \text{ and } \pi_{\phi}^* \text{ are constant on all } s_{\phi} \in \chi^{-1}(s_{\psi}) \text{ for all } s_{\psi} \text{ and } a$$
  
 $\Leftrightarrow q_{\phi}^*(s_{\phi}, a) = q_{\psi}^*(s_{\psi}, a) \text{ and } \pi_{\phi}^*(s_{\phi}) = \pi_{\psi}^*(s_{\psi}) \text{ for all } s_{\psi} = \chi(s_{\phi}) \text{ and } a.$ 

 $\psi \prec \phi$  means  $\psi$  is a better reduction than  $\phi$  since it leads to the same optimal q-value and policy as  $\psi$  does, but is more parsimonious (coarser) than  $\phi$ . If  $q_{\phi}^*$  or  $\pi_{\phi}^*$  is not constant on  $\psi$ -partitions, coarsening  $\phi$  to  $\psi$  and using  $\psi$  (potentially) leads to suboptimal solutions.

Enriching the order  $\prec$ :  $\prec$  is a transitive but 'very' partial order. Two maps are incomparable if neither is a refinement of the other. We can enrich order  $\prec$  as follows: For any two maps  $\psi$  and  $\psi'$ , the map  $\phi(h) := (\psi(h), \psi'(h)) \in \mathcal{S}_{\phi} = \mathcal{S}_{\psi} \times \mathcal{S}_{\psi'}$  refines both. Define  $\psi \prec_{\times} \psi'$  iff  $\psi \prec \phi \prec \psi'$ . Extended order  $\prec_{\times}$  is still not total. The remaining incomparable cases are: Case  $\psi \prec \phi \succ \psi'$ : This is only possible if  $q^*$  and  $\pi^*$  of  $\psi$  and  $\psi'$  (and  $\phi$ ) coincide. A secondary criterion based on the relative complexity of  $\psi$  and  $\psi'$  could decide the case, e.g.  $\psi \prec_{\times} \psi'$  iff  $|\mathcal{S}_{\psi}| < |\mathcal{S}_{\psi'}|$ . Case  $\psi \succ \phi \prec \psi'$ : Both  $\psi$  and  $\psi'$  are inferior to  $\phi$ . If class  $\Phi$  is closed under cartesian product,  $\phi$  should be favored over  $\psi$  and  $\psi'$  so their relative order is not or less important.

Search for  $\phi$ : Assume  $\Phi$  contains at least one exact reduction. Then the  $\prec_{\times}$ -minimal elements in  $\Phi$  are exactly the maximally coarse exact  $\phi \in \Phi$ . If  $\Phi$  is closed under arbitrary coarsening, then there is a unique minimizer (modulo isomorphism). If  $\Phi$  is also closed under cartesian product, the same holds for  $\prec$ . This implies that any exhaustive search for a  $\prec_{\times}$ -minimum in  $\Phi$  will give an exact  $\phi$  with minimal number of states, say  $\phi_0$ . Now Theorem 7 tells us that  $q_{\phi_0}^*$  and  $\pi_{\phi_0}^*$  are the optimal value and policy also of the original process P, irrespective of whether  $P_{\phi_0}$  is Markov or not. So while the conditions of Theorem 7 cannot be verified in practice, the theorem justifies a search procedure based on  $(q_{\phi}^*, \pi_{\phi}^*)$  that ignores the (non-)Markov structure of  $P_{\phi}$ .

Search for approximate  $\phi$  based on finite sample size. The principle approach in the previous paragraph is sound, but needs to be generalized in various ways before it can be used: Real sample size is finite, which means we only have access to approximations  $\hat{q}_{\phi}^*$  and  $\hat{\pi}_{\phi}^*$  via estimation  $\hat{p}_{\phi}$  of  $p_{\phi}$ . The criterion for exact equality  $q_{\phi}^* = q_{\psi}^*$  in  $\prec$  needs to be replaced by a suitable  $\hat{q}_{\phi}^* \approx \hat{q}_{\psi}^*$ , which should be done anyway, since real-word problems seldom allow for exact reductions.  $\approx$  should be chosen so as to come with statistical guarantees; e.g. Kolmogorov-Smirnov tests have been used in [McC96]. A suitable  $\hat{\pi}_{\phi}^* \approx \hat{\pi}_{\psi}^*$  requires more effort (see outlook). For large  $\Phi$  this also requires appropriate regularization, i.e. penalizing complex  $\phi$  [Hut09c]. To ensure  $\hat{q}^* \to q^*$  for  $n \to \infty$ , we need proper exploration strategies [SLL09]. Finally, we want an efficient search procedure in  $\Phi$ , rather than exhaustive search. This will be heuristic or will require strong assumptions on  $\Phi$  [Ngu13]. All but the last point raised above have or should have general solutions (see next paragraph).

Utilizing existing algorithms. The BLB algorithm [MMR11] and its extensions IBLB [NMRO13] and improvements OMS [NOR13] solve most of the problems above and can (nearly) readily be used for our purpose.

The BLB family uses the same basic FRL setup from [Hut09c] used also here. The authors consider a countable class  $\Phi$  assumed to contain at least one  $\phi$  such that  $P_{\phi}$  is an MDP (6). They consider average reward, rather then discounting, and analyze regret, which (in general) requires some assumption on the mixing rate or 'diameter' of the MDP. They prove that the total regret grows with  $\tilde{O}(n^{1/2...2/3})$ , depending on the algorithm.

Their algorithms and analyses rely on UCRL2 [JOA10], an exploration algorithm for finite-state MDPs. Going through the BLB proofs, it appears that the condition that  $P_{\phi}$  is an MDP can be removed if p (15) is used instead, modulo the analysis of UCRL2 itself. The proofs for the bounds for UCRL2 exploit that s',r' conditioned on s,a are i.i.d., which is true if  $P_{\phi}$  is Markov but not in general. Asymptotic versions should remain valid under the 'weak conditions' alluded to in (37). With some stronger assumptions that guarantee good convergence rates, the regret analysis of UCRL2 should remain valid too. Formally, the use of Hoeffding's inequality for i.i.d. need to be replaced by comparable bounds with weaker conditions, e.g. Azuma's inequality for martingales.

There is one serious gap in the argument above. BLB uses average reward while our theorems are for discounted reward. It is often possible to adapt algorithms and proofs which come with regret bounds for average reward to PAC bounds for discounted reward or vice versa. This would have to be done first: either a PAC version of BLB by combining MERL [LHS13] with UCRL $\gamma$  [LH12], or average reward versions of the bounds derived in this paper.

## 9 Miscellaneous

Action permutation instead of policy condition. We can rename actions without changing the underlying problem: Let  $A: \mathcal{A} \to \tilde{\mathcal{A}}$  be a bijection, and define  $\tilde{P}(o'r'|h\tilde{a}) := P(o'r'|ha)$ , where  $\tilde{a} := A(a)$ . Clearly, all results for P also hold for  $\tilde{P}$  if a is replaced by  $\tilde{a}$  everywhere, in particular  $\tilde{\Pi}(h) := A(\Pi(h))$ . In general, this is of little use. Things become more interesting if we allow the bijection A to be history-dependent, which we can do since our results hold for any, even non-stationarity,  $\tilde{P}$ . This allows us to devise an  $A: \mathcal{A} \times \mathcal{H} \to \tilde{\mathcal{A}}$  such that  $A(\Pi(h);h) = constant$  for the policy  $\Pi$  of interest. For example, for  $\tilde{\mathcal{A}} := \mathcal{A}$ , this is achieved by a permutation that swaps action  $a = \Pi(h)$  with some arbitrary but fixed action  $a^1 \in \mathcal{A}$ , and leaves all other actions unchanged:

$$A(a;h) := \begin{cases} a^1 & \text{if} & \Pi(h) = a \\ \Pi(h) & \text{if} & \Pi(h) \neq a = a^1 \\ a & \text{else} \end{cases}$$

Since  $\tilde{\Pi}(h) \equiv A(\Pi(h);h) \equiv a^1$  is constant, the  $\phi$ -uniformity condition for  $\tilde{\Pi}$  in Theorems 5, 6 and 9 becomes vacuous. While this transformation is theoretical interest, it only becomes practically useful if we can somehow learn the function A without knowledge of  $\Pi$ , and in particular for  $\Pi^*$ . We could also allow non-bijective A that merge actions that have (approximately) the same (optimal) Q-value.

## 10 Discussion

**Summary.** Our results show that RL algorithms for finite-state MDPs can be utilized even for problems P that have arbitrary history dependence and history-to-state reductions/aggregations  $\phi$  that induce  $P_{\phi}$  that are also neither stationary nor MDP. The only condition to be placed on the reduction is that the quantities of interest, (Q-)Values and (optimal) Policies, can approximately be represented. This considerably generalizes previous work on feature reinforcement learning and MDP state aggregation and allows for extreme state aggregations beyond MDPs. The obtained results may also explain why RL algorithms designed for MDPs sometimes perform well beyond MDPs.

**Outlook.** As usual, lots remains to be done. A list of the more interesting remaining tasks and open questions follows:

- While the approximate  $\phi$ -uniformity condition on  $Q^*$  in Theorem 7 is very weak compared to bisimilarity, uniformity of  $V^*$  in Theorem 9 is even weaker (Theorem 11 shows how much of a difference this can make). It is an Open Problem 10 whether an analogue of Theorem 7ii also holds for Theorem 9 beyond  $\varepsilon = 0$ .
- An algorithm learning  $\phi$  beyond MDPs that comes with regret or PAC guarantees has yet to be developed. This could be done by generalizing the partial order

- $\prec_{\times}$  to  $n < \infty$ , or by adapting the class and proofs of BLB algorithms, or by integrating MERL with UCRL $\gamma$ . All bounds contain  $\frac{1}{1-\gamma}$  to some power. Can the exponents be improved? For which environments/examples are the bounds tight?
- The trick to use a-dependent  $Q^*$  as a-independent map  $\phi$  in Section 6 was to vectorize  $Q^*$  in a. Unfortunately this leads to a state-space size exponential in  $\mathcal{A}$ . Solution  $\phi$  based on  $(V^*,\Pi^*)$  pair is only linear in  $\mathcal{A}$ , but rests on Open Problem 10. Are there other/better ways of dealing with actions? Other extreme aggregations  $\phi$ , or are a-dependent  $\phi$  possible?
- Are average-reward total-regret versions of our discounted reward results possible, under suitable mixing rate conditions?
- For small discrete action spaces typical for many board games, the exact conditions on  $\Pi$  are met. For continuous action spaces as in robotics, we can simply discretize the action space, introducing another  $\varepsilon$ -error, but action-continuous versions of our results would be nicer. Except for Theorem 7, any interesting generalization should replace the exact by approximate  $\phi$ -uniformity conditions on  $\Pi$ .
- Our theorems and/or proof ideas should allow to extend existing convergence theorems for RL algorithms such as Q-learning and others from MDPs to beyond MDPs.
- The bisimulation conditions of classical state aggregation results are for reward and transition probabilities. It would be interesting to derive explicit weaker conditions for them that still imply our conditions on (Q-)Values.

## References

- [EGW05] D. Ernst, P. Geurts, and L. Wehenkel. Tree-based batch mode reinforcement learning. *Journal of Machine Learning Research*, 6:503–556, 2005.
- [Faz06] I. Fazekas. On a general approach to the strong laws of large numbers. Technical report, Faculty of Informatics, University of Debrecen, Hungary, 2006.
- [FK01] I. Fazekas and O. Klesov. A general approach to the strong law of large numbers. Theory of Probability & Its Applications, 45(3):436–449, 2001.
- [FPP04] N. Ferns, P. Panangaden, and D. Precup. Metrics for finite Markov decision processes. In *Proc. 20th conf. on Uncertainty in Artificial Intelligence (UAI'04)*, pages 162–169, 2004.
- [GDG03] R. Givan, T. Dean, and M. Greig. Equivalence notions and model minimization in Markov decision processes. *Artificial Intelligence*, 147(1–2):163–223, 2003.
- [Hut05] M. Hutter. Universal Artificial Intelligence: Sequential Decisions based on Algorithmic Probability. Springer, Berlin, 2005.
- [Hut09a] M. Hutter. Feature dynamic Bayesian networks. In *Proc. 2nd Conf. on Artificial General Intelligence (AGI'09)*, volume 8, pages 67–73. Atlantis Press, 2009.
- [Hut09b] M. Hutter. Feature Markov decision processes. In *Proc. 2nd Conf. on Artificial General Intelligence (AGI'09)*, volume 8, pages 61–66. Atlantis Press, 2009.
- [Hut09c] M. Hutter. Feature reinforcement learning: Part I: Unstructured MDPs. *Journal of Artificial General Intelligence*, 1:3–24, 2009.

- [JOA10] T. Jaksch, R. Ortner, and P. Auer. Near-optimal regret bounds for reinforcement learning. *Journal of Machine Learning Research*, 11:1563–1600, 2010.
- [LH12] T. Lattimore and M. Hutter. PAC bounds for discounted MDPs. In *Proc. 23rd International Conf. on Algorithmic Learning Theory (ALT'12)*, volume 7568 of *LNAI*, pages 320–334, Lyon, France, 2012. Springer.
- [LH14] T. Lattimote and M. Hutter. General time consistent discounting. *Theoretical Computer Science*, 519:140–154, 2014.
- [LHS13] T. Lattimore, M. Hutter, and P. Sunehag. The sample-complexity of general reinforcement learning. *Journal of Machine Learning Research*, W&CP: ICML, 28(3):28–36, 2013.
- [McC96] A. K. McCallum. Reinforcement Learning with Selective Perception and Hidden State. PhD thesis, Department of Computer Science, University of Rochester, 1996.
- [MMR11] O.-A. Maillard, R. Munos, and D. Ryabko. Selecting the state-representation in reinforcement learning. In *Advances in Neural Information Processing Systems* (NIPS'11), volume 24, pages 2627–2635, 2011.
- [Ngu13] P. Nguyen. Feature Reinforcement Learning Agents. PhD thesis, Research School of Computer Science, Australian National University, 2013.
- [NMRO13] P. Nguyen, O. Maillard, D. Ryabko, and R. Ortner. Competing with an infinite set of models in reinforcement learning. *JMLR WS&CP AISTATS*, 31:463–471, 2013.
- [NOR13] O.-A. Maillard P. Nguyen, R. Ortner, and D. Ryabko. Optimal regret bounds for selecting the state representation in reinforcement learning. *JMLR W&CP ICML*, 28(1):543–551, 2013.
- [NSH11] P. Nguyen, P. Sunehag, and M. Hutter. Feature reinforcement learning in practice. In *Proc. 9th European Workshop on Reinforcement Learning (EWRL-9)*, volume 7188 of *LNAI*, pages 66–77. Springer, September 2011.
- [Ort07] R. Ortner. Pseudometrics for state aggregation in average reward Markov decision processes. In *Proc. 18th International Conf. on Algorithmic Learning Theory* (ALT'07), volume 4754 of LNAI, pages 373–387, Sendai, Japan, 2007.
- [Put94] M. L. Puterman. Markov Decision Processes Discrete Stochastic Dynamic Programming. Wiley, New York, NY, 1994.
- [RN10] S. J. Russell and P. Norvig. Artificial Intelligence. A Modern Approach. Prentice-Hall, Englewood Cliffs, NJ, 3rd edition, 2010.
- [SB98] R. S. Sutton and A. G. Barto. Reinforcement Learning: An Introduction. MIT Press, Cambridge, MA, 1998.
- [SH10] P. Sunehag and M. Hutter. Consistency of feature Markov processes. In *Proc.* 21st International Conf. on Algorithmic Learning Theory (ALT'10), volume 6331 of LNAI, pages 360–374, Canberra, Australia, 2010. Springer.
- [SLL09] A. L. Strehl, L. Li, and Michael L. Littman. Reinforcement learning in finite MDPs: PAC analysis. Journal of Machine Learning Research, 10:2413–2444, 2009.
- [UV98] W. T. B. Uther and M. M. Veloso. Tree based discretization for continuous state space reinforcement learning. In AAAI, pages 769–774, 1998.
- [VGS05] V. Vovk, A. Gammerman, and G. Shafer. Algorithmic Learning in a Random World. Springer, New York, 2005.

## A List of Notation

### General notation

[R] = 1 if R=true and =0 if R=false (Iverson bracket)

 $\#\mathcal{X}$  size of set  $\mathcal{X}$ 

 $\varepsilon,\delta$  small non-negative real numbers

|z| largest integer  $\leq z$ 

### Original history-based process

 $\mathcal{O}, \mathcal{R}, \mathcal{A}$  = finite observation, reward, action spaces.

 $o_t r_t a_t \in \mathcal{O} \times \mathcal{R} \times \mathcal{A} = \text{observation, reward, action at time } t$ 

 $t \le n \in \mathbb{N}$  = any time  $\le$  sample size

 $P,Q,V,\Pi$  = Probability, (Q-)Value, Policy of original history-based Process

 $\Pi^*, \tilde{\Pi}, \Pi_B$  = optimal, approximately optimal, behavior Policy

 $h \in \mathcal{H}$  =  $(\mathcal{O} \times \mathcal{R} \times \mathcal{A})^* \times \mathcal{O} \times \mathcal{R}$  = possible histories of any length

 $h' = hao'r' = \text{successor history of } h \in \mathcal{H}$ 

 $h_t = o_1 r_1 a_1 \dots o_t r_t = \text{history up to time } t$ 

 $\mathcal{H}_t$  =  $(\mathcal{O} \times \mathcal{R} \times \mathcal{A})^{t-1} \times \mathcal{O} \times \mathcal{R}$  = history of length t

P(o'r'|ha) = probability of next observation&reward given history&action

### Reduction/aggregation from history to states

 $S_{\phi}$  = finite state space induced by  $\phi$  (range of  $\phi$ )

 $\phi: \mathcal{H} \to \mathcal{S}_{\phi} = \text{reduction/map/aggregation from histories to states}$ 

 $s_t = \phi(h_t) \in \mathcal{S} = \text{state at time } t$ 

 $P_{\phi}(s'r'|ha) = \text{marginalized } P\text{-probability over state\&reward given history\&action}$ 

B(h|sa) = dispersion probability. Stochastic "inverse" of  $\phi$ 

 $\langle Q(h,a)\rangle_B = B$ -average over  $\{\hat{h}: \phi(\hat{h}) = \phi(h)\}$ 

 $w_t(sa)$  = non-negative weight function  $\sum_{t=1}^{\infty} w_t(sa) = 1 \ \forall sa$ 

 $P_B(h)$  = probability of h from P interacting with  $\Pi_B$ 

 $P_{\phi B}()$  = (partially)  $\phi$ -reduced, marginalized, conditionalized  $P_B$ 

 $\prec, \prec_{\times}$  = (extended) ordering of  $\phi$  w.r.t. quality ( $n = \infty$  so far only)

### Finite state Markov decision process (MDP)

S = finite state space

 $p,q,v,\pi$  = probability, (q-)value, policy of MDP

s,a,s',r' = stat, action, successor state, reward

n(sas'r') = number of times sas'r' appears in  $h_{n+1}$ 

 $\gamma \in [0;1)$  = discount factor