

## RECONFIGURABLE REINFORCEMENT LEARNING NETWORKS

In humans, the process of learning is not only driven by the environment and structure of the brain. The development of the brain-structure itself defines what learning may take place; thus the conditions and patterns which direct brain formation are primary and total for the success of learning. As artificial intelligence research continually generates and publishes on novel structures discovered by humans, this work is centered around how the novel structures can be discovered using RLNs. This subject is frequently included in the subject of general intelligence, and is famous for both its philosophical and computational complexity, as well as its difficulty in finding funding. There have been previous works on this subject, such as [Consciousness as a State of Matter] [] [].

Specifically, this work presents a unification method for online learning (via Reinforcement Learning) and offline learning (via Backpropagation). In the most general sense, this work demonstrates an approach to the self-structuring of parametric models. First, it is reviewed that Concurrent Markov Decision Processes (CMDPs) can model parametric structure and facilitate optimal behaviour even when subject to large state spaces and generous state uncertainty. Second, it is shown that a variation of CMDPs called Reconfigurable Learning Networks (RLNs) can learn parametric decision networks. RLNs in structure and behaviour turn out to be equivalent to the structure and behaviour of feed-forward neural networks. Lastly, a few empirical examples are demonstrated, beginning with the MINST dataset. Two main contributions are made: First, RLNs can be trained online and offline, using Reinforcement Learning and then Backpropagation; online learning stimulates network growth and adaptation immediately, whereas backpropagation seems to be an ideal phase for network pruning. Second, an RLN can achieve empirical success even when the reward function for the system is changed dynamically. Thus both a degree of empirical success and general learning have been achieved.

In order for a generally intelligent system to operate, solutions to several open problems need to be solved analytically and/or heuristically. In this work, we present the related problem categories in the Introduction (Section 1), and include background on each area. Second, most of this work is focused around the reconfiguration of existing MDP models, so Section 2, Mapping, includes work on transfer learning and analytical analysis. Third, we express how convergence of behaviour policies can be preserved despite online RLN restructuring (Section 3). The tradeoff between network structure and computation time in learning is expressed analytically (Section 4). Lastly, it is shown that RLNs are actually just feed-forward Neural Networks, which adds the ability to use back propagation and other techniques on discovered models (Section 5).

In this work due to the difficulty of the subject matter initially, models are assumed noiseless and stochastically stable. It is expected that later work will broaden this work by considering state uncertainty, and non-stationary problems.

## NOTATION

In general, most online optimization problems can be expressed as fully observable Markov Decision Processes (MDPs) as  $\langle S, A, T, R, \pi \rangle$  tuples:

- $S \subseteq R^n$ : A discrete collection of states.

- $A \subseteq R^n$ : A discrete collection of actions.
- $T(s'|a, s) \in R$ : A stably stochastic transition function, where  $\sum_{s \in S} T(s|a, s') = 1$
- $R(s'|a, s) \in R$ : A stable stochastic reward function
- $\pi : S \times A \rightarrow R$ : A non-negative behaviour policy with the general property,  $\sum_{a \in A} \pi(s, a) = 1$

In general, we can express behaviour in this domain as a policy  $\pi : S \rightarrow R$  [**? looked like this, but would  $\pi : S \rightarrow A$  make more sense?**]. Particular attention is given to the optimal strategy.

In prior work the issue of tractability and subsequent decomposition have been articulated. In this work the subject of learning and generalizing this decomposition work into a General framework is discussed.

Ⓐ Theory	{	Section 1: Background & Introduction	Background of relevant research & RLN introduction
		Section 2: Mapping	a generalized set of mapping & deconstruction operations (parent, child, reward optimization, complexity
		Section 3: Convergence (16):	parent, child, reward optimization, complexity
		Section 4: Worst Case Performance (23):	system behaviour with malformed problems
Ⓑ NN paper	{	Section 5: Neural Networks (24):	RRLN are just feed forward Neural Networks
		$\hookrightarrow (N.)$	

Special topics:

Temporal Difference (A1): how to discover & change time basis/scale

Transitional Learning (A2): how to re-use and generalize transitional models

Financial Systems (A3): how to use with financial systems

Origins (E1-E4): original examples and sketches

Transitional Encoding (E5-E6): Continuous Gaussian mixture models & applications

# 1 BACKGROUND AND INTRODUCTION

## INTRODUCTION: A RECONFIGURABLE REINFORCEMENT LEARNING METHOD

### 1.1 Paper outline

@ General RL-mapping prior knowledge ( $\tilde{R}$  and  $\tilde{T} \rightarrow \tilde{Q}$ )

#### 1.1.1 Introduction

One of the largest issues facing Q-Learning and Reinforcement Learning is the coupling of reward and state behaviours. In humans, this behaviour is called anchoring. A person can learn to associate reward with stimulus via Pavlov's effect. Importantly, humans possess the ability to disassociate reward from the prediction of system dynamics. For robots to be truly flexible learners, they must be able to separately encode the reward in a system from the transition dynamics in a system. This allows one robot to maximize different reward functions without retraining.

#### 1.1.2 Q-Learning Extension

Specifically, given a perfect encoding  $Q^*(s, a)$ , transition and reward are

$$\arg \max_a Q_t^*(s, a) \sim \pi^*(s) = \arg \max_a \sum_{s'} R(s'|a, s)T(s'|a, s) + \gamma V(s') \quad (1)$$

where

$$Q_t^*(s, a) = Q_{t-1}(s, a) + \alpha (Q_{t-1}(s, a) - R(s'|s, a)) + \gamma Q_{t-1}(s^*, a^*) \quad (2)$$

$$Q_t^*(s, a) = Q_{t-1}(s, a) + \alpha (Q_{t-1}(s, a) - R(s'|s, a)) + \gamma \arg \max_{a^*} Q(s', a^*) \quad (3)$$

Unfortunately, as time  $t$  increases in value, the values  $R(s|s)$  and  $T(s|a, s)$  are encoded such that  $f : R \cdot T \rightarrow Q_0$  it is not possible to compute  $f^{-1}(Q)$  due to the loss of information.

To avoid information loss, a dynamics equation  $\tilde{T}$  and observed reward policy  $\tilde{R}$  can be tracked, and used to compute at any time:

$$\tilde{Q}_t^*(s, a) = \sum_{s'} \tilde{T}_{t-1}^*(s'|s, a) \tilde{R}_{t-1}(s', s, a) + \gamma \pi^*(s') \quad (4)$$

where  $\tilde{\pi}^*(s, a) \sim Q^*(s, a)$ .

It is clear that (4) must iterate over all states  $s'$ , which is intractable. However, if optimal policy locality for a function  $L$  is assumed, the resulting expression ?? is optimal.

**Lemma 1:** Policy optimality for the purposes of regression to optimal Q values ( $\tilde{Q}^*, Q^*$ ), consideration of "local" states and "local" acting are required,

$$L : S, A \rightarrow S, A|_{\text{onto}} \quad (5)$$

$$\sum_{s' \in S} \tilde{T}_t(s'|s, a) \tilde{R}_t(s'|s, a) + \gamma \tilde{\pi}(s') = \sum_{s' \in L(s)} \tilde{T}_t(s'|s, a) \tilde{R}_t(s'|s, a) + \tilde{\gamma}(s') \quad (6)$$

$$\arg \max_{a \in A} \tilde{Q}^*(s, a) = \arg \max_{a \in L(a)} \tilde{Q}^*(s, a) \quad (7)$$

$$\tilde{\pi}^*(s|L) = \tilde{\pi}^*(s) \quad (8)$$

The local policy is equivalent to the complete state  $\times$  action space version.

### 1.1.3 Regression

Instead of regressing to Q values directly ( $Q_t \leftarrow f(Q_{t-1})$ ) we instead regress to system dynamics ( $\tilde{T}_t \leftarrow f(T_{t-1})$ ). Specifically:

$$\tilde{T}_{t+1}(s'|s, a) = \frac{\text{freqn}(s'|s, a)}{\text{freqn}(s, a)} \quad \text{or} \quad \tilde{T}_{t+1}(s'|s, a) = P(s'|s, a) \quad (9)$$

$$\underbrace{\tilde{R}_{t+1}(s'|s, a) = R_t(s'|s, a)}_{\text{simple}} \quad \text{or} \quad \underbrace{\tilde{R}_{t+1}(s'|s, a) = \tilde{R}(s'|s, a) + (\tilde{R}(s'|s, a) - R(s', a, s))}_{\text{verbose}} \quad (10)$$

\*—  $\tilde{T} \approx T$  when

Thus it is possible, using either a simple regression model or a verbose training model, to regress to an optimal policy  $\tilde{\pi}^*$ , online, and exploit reward switching.

### 1.1.4 Behaviour Switching

After an agent is trained ( $\tilde{\pi} \sim \tilde{\pi}^*$ ), then it is possible to redefine the reward function  $\tilde{R}$ . Given  $\tilde{R}_t = R_A$ , then  $\tilde{R}_{t+1} \leftarrow R_B$  is possible. At this point the transition dynamics ( $\tilde{T}$ ) and locality principle ( $L$ ) maybe be used to infer  $\tilde{\pi}_{t+1}$ . Just as a human may redefine an objective during learning, agents may also change the definition of success through reassignment of  $\tilde{R}$ .

1.1.4.1 Optimality of Switching: Given a set of reward functions  $\mathbf{R} = \{R_1, R_2, \dots\}$  and a set of time indexes  $\{1, 2, 3, \dots, t, \dots\}$  as time approaches infinity ( $t \rightarrow \infty$ ) the global policy  $\tilde{\pi}_t^*(s)$  will converge on the optimal policy  $\tilde{\pi}^*$ ,  $\tilde{\pi}^* \in \Pi(R_i)$ , where  $\Pi(R_i)$  is the set of all optimal policies for reward function  $R_i$ .

Set of reward functions characterizing a set of goals

$$\mathbf{R} = \{R_1, R_2, \dots, R_i, \dots\} \quad R_i : S \times A \times S = \mathbb{R}^+, \quad R_i \text{ known} \leftarrow \mathbb{R}^+ \quad (11)$$

An “expected quality” of a state-action pair given locality & states.

$$\tilde{Q}_t(s, a|R_i, L) = \sum_{s' \in L(s)} \tilde{T}_t(s'|s, a) R(s'|s, a) + \gamma \arg \max_{a^*} \tilde{Q}_t(s', a^*|R_i, L) \quad (12)$$

Set of optimal policies

$$\Pi, \quad \pi^* \in \Pi \text{ iff } : \pi^*(s) = \arg \max_a Q^*(s', a^*|L, R_i) \quad (13)$$

optimal Q-value

$$Q^*(s, a|L, R_i) = \sum_{s' \in L(s)} T(s'|s, a) R_i(s'|s, a) + \gamma \arg \max_{a^*} Q^* \quad (14)$$

1.

$$Q^*(s, a|L, R_i) = \tilde{Q}_t(s, a|L, R_i), \quad t \rightarrow \infty \quad (15)$$

$$Q^*(s, a|L, R_i) = E \left[ \sum_{s' \in L(s)} T(s'|s, a) R_i(s'|s, a) + \gamma \arg \max_{a^*} Q^*(s', a^*|R_i, L) \right] \quad (16)$$

$$= \sum_{s' \in L(s)} E[T(s'|s, a)] E[R_i(s'|s, a)] + \gamma \arg \max_{a^*} E[Q^*(s', a^*|R_i, L)] \quad (17)$$

$$\text{obvious: } E[Q^*] \equiv Q^*$$

$$\text{obvious: } R_i(s'|s, a) \equiv E[R_i(s'|s, a)]$$

$$= \sum_{s' \in L(s)} E[T(s'|s, a)] R_i(s'|s, a) + \gamma \arg \max_{a^*} Q^*(s', a^*|R_i, L) \quad (18)$$

$$\begin{aligned} * \text{ Prove } \longrightarrow & \quad \text{assert } \lim_{t \rightarrow \infty} \tilde{T}_t(s'|s, a) = E[T(s'|s, a)] \\ & = \lim_{t \rightarrow \infty} \sum_{s' \in L(s)} \tilde{T}_t(s'|s, a) R_i(s'|s, a) + \gamma \arg \max_{a^*} Q^*(s', a^*|R_i, L) \end{aligned} \quad (19)$$

$$Q^*(s, a|L, R_i) = \tilde{Q}_t(s, a|L, R_i), \quad t \rightarrow \infty \quad (20)$$

2.

Trivially:

$$\text{if } \tilde{Q}_t(s, a|L, R_i) \sim Q^*(s, a|L, R_i) \text{ iff } \tilde{T}_t \sim T \quad (21)$$

$$\forall R_j \in \Pi : \tilde{Q}_t(s, a|L, R_j) \sim Q^*(s, a, L, R_j) \quad (22)$$

### 1.1.5 Validation

To the value the efficiency of this model we consider a staged experiment with three reward functions.

Thus, in this experiment we consider training a model, first, on rewarding trading actions that essentially reward low exposure and no loss ( $R_1$ ). During this period  $\tilde{T}$  will mature, and many dynamics of the system are discovered. After gains reach 10%,  $R_2$  becomes active — given short term success in a security, we increase buy in. After this, given success, we reward an agent that tries to go “all in” to grow profits to 50%. Lastly, at 50%, we “clamp” the agent from being rewarded for taking excessive risks.

This reward profile, across industry, is common. Humans design reward profits and procedures that emphasize the result of behaviour. For example, on the trading floor, we are interested in policies that manage the risk-profit profile.

To experiment, we compose two processes:

- ① a Q-learning mechanism
- ② a  $\tilde{T}$ -learning mechanism

For both models, the following definitions are used.

$S$  — exposure  $\times$  profit  $\times$  price, where exposure  $\in [0, 100]$ , profit  $\in \mathbb{R}$ , price  $\{x \in \mathbb{R}^+\}_{x=0}^{100}$ , representing exposure %, profit %, and price history

$A$  — buy 1%, sell 1%, no action

$T = T(s'|s, a)$  — probability of making an action and ending up in a later state

$R$  —

$$R(s'|s, a) = \begin{cases} R_1 \text{ iif profit} < 10\% & R_1 = \frac{\Delta \text{profit}}{\text{exposure}} - \text{exposure} \\ R_2 \text{ iif profit } 10\% - 20\% & R_2 = \Delta \text{profit} \\ R_3 \text{ iif profit } 20\% - 50\% & R_3 = \Delta \text{profit}(\text{exposure}) \\ R_4 \text{ iif profit } 50\%+ & R_4 = \frac{\text{profit}}{\text{exposure}} \end{cases} \quad (23)$$

$\tilde{T}$  — predicted transitional dynamics — encoded using random forrest classification

Approach ① – Q-Learning	Approach ③ – ensemble	Approach ② – T-learning
$R(s' s, a)$ will exhibit a ‘steady stochastic’ input pattern and the average reward $E[R(s' s, a)]$ will be encoded into the $Q(s, a)$	We train separate Q-learning agents, each which are active given the use of $R_1, R_2, R_3, R_4$ , such that $R_i \rightarrow Q^i$ .	$\tilde{Q}(s, a)$ will never encode more than one, singular, reward function $R_t \approx \{R_1, R_2, R_3, R_4\}$ . Thus it is expected that T-learning will vastly out-perform Q-learning.

### 1.1.6 Experiment

We ran standard basic testing on two independent stocks. One Google Inc., and the other, Blockbuster. Below we have plotted each models: Reward over time (Figure ??), Profit over time (Figure ??), Risk over time (Figure ??). All of the models are too juvenile in approach to be profitable. T-learning exhibits a xx% increase in reward acquisition, a xx% increase in profit and a xx% increase in Reward. The ensemble approach failed to outperform the T-learning approach. The highest performance approach is T-learning.

### 1.1.7 Summary

In this paper a method of reusing experience over several reward functions was demonstrated. T-learning allows agents to recycle transtiion information discovered through online learning. The results of using this experiment in the financial domain was a reduction in model complexity, and increased performance. Using a singular Q-learning agent leads to a loss of understanding of the reward function. Using the ensemble approach requires the transition dynamics of the experiment to be encoded four times. Thus, T-learning had both the best performance, and can scale to an unlimited set of reward functions, without retraining. It is the opinion of the author that the encoding of transition and reward information separately should be seen as the fundamental first step when considering model reuse and generalization.

## 1.2 RLN MDP structure

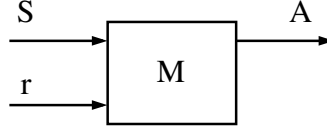
The general approach that is taken to form an RLN is “split” one single MDP into parent and child processes. Doing so assumes the two process are partly independent [CMDP].

1. Largely, the models may be independent.

2. The child and parent may include elements of each other's MDP definition in their own definition.

This section focuses on framing a model. In Sections 2 and later, subjects related to behavior optimality, convergence, and computational complexity are considered.

In order to consider the formation of a reconfigurable reinforcement network, it is required to analytically group all aspects of a process and a behavior policy into one tuples.



To do this, assume a Markov decision process  $M$  which can be internally modeled as a tuple  $M = \langle S, A, T, R, \pi, \tilde{T}, \tilde{R} \rangle$

$S$  – a set of states  $s \in S$  which may be experienced by  $M$

$A$  – a set of actions  $a \in A$  that may be executed

$T$  – a true transitional probability,  $T(s'|a, s)$  expressing the probability of executing an action  $a$  in state  $s$  before ending up in later state  $s'$ .

$R$  – is a reward function which quantifies how desirable a transition  $R(s'|a, s)$  is.  $R : S \times A \times S \rightarrow \mathbb{R}_{\geq 0}$

$\tilde{T}$  – is the current model of  $T$ . The goal of  $\tilde{T}$  is thus  $\tilde{T} \sim T$

$\tilde{R}$  – is the predicted reward of the system, constructed from observation of  $R$ , s.t.  $\tilde{R} \rightarrow R$ .

$\pi$  – is an action selection policy, ideally chosen to maximize expected reward, an optimal policy is denoted  $\pi^*$ .

Ideally

$$\pi^*(s) = \arg \max_a \sum_{s'} \underbrace{R(s'|a, s)T(s'|a, s) + \gamma V(s')}_{\text{expected reward}}$$

and

$$\tilde{\pi}(s) = \arg \max_a \sum_{s'} \tilde{R}(s'|a, s)\tilde{T}(s'|a, s) + \gamma \tilde{V}(s')$$

[Note that  $\tilde{V}$  hasn't been defined in the previous equation.]

## ENCODING

$$\pi^*(s) = \arg \max_a \sum_{s'} R(s'|s, a)T(s'|s, a) + \gamma V(s')$$

where

$$V(s) = \sum_{s'} R(s'|s, a)T(s'|s, a) + V(s').$$

A bellman backup can be used [Bellman backup]. In online applications stochastic gradient descent can be applied to regress to locally optimal solutions. This allows estimation of optimal policy

$$\pi^*(s) = \arg \max_a Q(s, a).$$

To encode the expected reward over all states, typically  $Q$ -values are kept:  $Q(s, a) \sim \sum R(s'|a, s)T(s'|a, s) + \gamma V(s')$  and  $Q_{t+1}(s, a) \leftarrow Q_t(s, a) + \alpha (R(s'|a, s) - Q_t(s, a) + \gamma \arg \max_{a'} Q(s', a'))$ .

To render the process  $M$  separable, it is necessary to decouple the transitional values  $T$  from the reward values  $R$ . Thus, to directly encode  $Q(s, a)$  using  $\pi$  is prohibitive.

knowing:

$$Q : S \times A \rightarrow \mathbb{R}_{\geq 0}$$

indirect encoding:

$$\pi : \{S \times A \times S \times \mathbb{N} | R\} \rightarrow Q$$

If the definitions of  $S$  or  $A$  change, then  $Q$  must be reinitialized. Alternatively,  $\tilde{T}$  and  $\tilde{R}$  are defined as intermediate encoding functions. Thus we define

$$\pi : \{S \times A \times S \times \mathbb{N}\} \rightarrow \tilde{T}, \tilde{R}, Q$$

simple transition

$$\tilde{T}_{t+1}(s'|s, a) = \frac{\text{freq}(s'|s, a)}{\text{freq}(s, a)}$$

simple reward

$$\tilde{R}_{t+1}(s'|s, a) = \tilde{R}(s'|s, a) + \alpha_R \left( \tilde{R}(s'|s, a) - R(s'|s, a) \right)$$

$$f_Q : \tilde{T}_t \times \tilde{R}_t \rightarrow Q_t$$

In this paper we rely on a method of extracting dynamic  $Q$ -values from an encoded transition and reward function  $(\tilde{T}, \tilde{R})$ . The motivation for this encoding is that it allows mapping the transition function into multiple spaces, and allows the reward function to be altered. The significance of this finding is covered in ??? Price wash ???.

### 1.3 Representation and mapping of $M$

Reconfiguring a process  $M$  allows some intractable MDPs to be rendered tractable. As an example, take an MDP  $M$  modeling a 3-dimensional foraging experiment with three thousand positions on the  $x$ ,  $y$ , and  $z$  axes respectively. This process will consume over three billion memory locations and may be impossible to explore. If this system is broken into three sub problems, each targeting a special axis, then only nine thousand memory locations need be consumed. This decreases memory requirements by an exponential factor.

This section presents a method of decomposition that introduces no degeneration of the regressed policy  $\pi(s, a)$ . The trade-off for saving in space is exponentially-increased computation. As in all problems, finding the balance between space and computational complexity is required.

#### 1.3.1 Introduction

The general approach is related to factor analysis, or clustering, eigenvalue decomposition, or projected component analysis. An MDP is split such that a subspace  $S_i \subseteq S$ , and action space  $A_i \subseteq A$  seem independent of subspaces  $S_k, A_k$ .

Specifically, these candidate subproblems can be considered as a replacement for the centralized MDP under some conditions. Two candidate subproblems, although independent, are “explorable” and “concurrently



maximized". Explorability means that for a candidate split  $C_x = \langle S_i, A_i, S_k, A_k \rangle$ , the transition function  $T(s_i|a_i, s_i) = P(s_i|a_i, s_i)$  exhibits full readability. Second, the problem needs to be "concurrently maximized". In this Section it will be shown that a completely explorable candidate,  $C_x$ , can also be concurrently maximized. First, definitions of explorability are articulated. Second, concurrent maximization is discussed.

**1.3.1.1 Candidates:** A candidate split is a tuple  $C_x = \langle S_i, A_i, S_k, A_k \rangle$ , where  $S_i \cup S_k = S$ , and  $A_i \cup A_k = A$ . A candidate also comes with a dynamic region  $G_i : S \rightarrow P(S_i)$ ,  $G(S) \subseteq S_i$ . Given these notations explorability can be defined.

**1.3.1.2 Explorability:** *i*-explorable: a candidate split is *i*-explorable in  $C_x$  if all states are "reachable" and independent of  $S_k, A_k$ :

$$\forall s'_i, \exists s_i, a_i : T(s'_i|a_i, s_i) > 0$$

where  $s'_i \in G(s_i)$ ,  $s_i \in G(s_i)$ ,  $a_i \in A_i$ ,  $a_k = \emptyset$  (inaction).

Completely *i*-explorable: Given  $G(s_i) = S_i$ , then a candidate  $C_x$  can be said to be completely explorable. This is generally intractable to compute, but may be inferred or deduced in other manners.

Explorable: A candidate split  $C_x$  is explorable if it is *i*-explorable and *k*-explorable for  $G_i(s_i)$  and  $G_k(s_k)$ .

**1.3.1.3 Concurrent maximization:** Given a completely explorable candidate  $C_x$ , it is possible to regress to two independent and optimal behaviour policies simultaneously,  $\pi_i^*, \pi_k^*$  if it is also the case that the optimal policy  $\pi^*(s) = \pi^*(s_i) \cup \pi^*(s_k)$ . (See concurrent maximization proof.)

Concurrent maximization can occur when (a) the policies are learned in a synchronous manner, where  $a_i$  and  $a_k$  are executed in the same instant, or asynchronously, where  $a_i$  and  $a_k$  are executed in independent instants. As explained in the concurrent maximization proof, both the synchronous and asynchronous methods bear separate conditions.

Synchronous conditions:

- i.  $\bar{S}_k$  is a (i) stably stochastic random variable, or (ii)  $\bar{S}_k$  is chosen such that  $E \left[ R_t \left( \begin{array}{c|cc} s'_i & a_i & s_i \\ \bar{s}'_k & a_k & \bar{s}_k \end{array} \right) \right] > E \left[ R_{t+1} \left( \begin{array}{c|cc} s_i & a_i & s'_i \\ \bar{s}'_k & a_k & \bar{s}_k \end{array} \right) \right]$ .
- ii.  $\bar{R}_t(s_i|a_i, s_i) = E \left[ \sum_{s'_k} \sum_{a_k} \pi(a_k|s_k) T(s'_k|a_k, s_k) R \left( \begin{array}{c|cc} s'_i & a_i & s_i \\ s'_k & a_k & \bar{s}_k \end{array} \right) + V(s'_i, s'_k) \right]$
- iii. Candidate  $C_x$  is completely explorable.
- iv.  $\pi_k^*(s_k)$  is either (a) stably stochastic, or (b) reward maximization.

When a synchronous method is used,  $a_i$  and  $a_k$  can be executed in tandem according to  $\pi_k()$  and  $\pi_i()$ . However, continual evaluation of  $\bar{R}_t(\bar{s}'_i|a_i, s_i)$  is required according to ii, and in the general case requires a computational complexity of  $O(|A_k| |s_k|)$  to compute per time step. This may be intractable.

Asynchronous conditions:

If  $a_i$  and  $a_k$  are executed in different instances, between reward observance by the regression approach, then only conditions i, iv are required.

v.  $\bar{s}'_k$  is chosen

Lastly, during execution of a policy  $\pi(s)$ , it may be required to evaluate many candidate splits,  $C_1, C_2, C_3, \dots, C_x, \dots$  for the purposes of deciding on an optimal split. For each split  $C_x$  a theorized explorability can be measured, after choice.

- i. the size of  $G(s_i)$  is measured as  $|G(s_i)|$
  - ii. the explored space is tracked  $\text{ex}(C_x) \leq G(s)$
  - iii. the purposeful transitions are tracked,  $\text{count}_p(a \in A_i)$
  - iv. the meddlesome transitions are tracked,  $\text{count}_m(a \in A_k)$
- if  $\frac{|\text{ex}(C_x)|}{|G(s)|} > \mathfrak{Z}_{\text{ex}}$  and  $\frac{\text{count}_p}{\text{count}_m} > \mathfrak{Z}_m$  then  $C_x$  is taken as a global candidate.

(Is that symbol  $\xi$  or  $\mathfrak{Z}$ ?)

#### Concurrent Maximization Proof:

Synchronous: a synchronous case is where  $a_i \in A_i$  and  $a_k \in A_k$  can be executed at the same time.

$$\begin{aligned}
 \pi_i^*(s_i) &\sim \arg \max_{a_i} \sum_{s'_i \in G(s_i)} \left[ T(s'_i|a_i, s_i) \sum_{s'_k \in G(s_i)} \sum_{a_k} \left( \pi(a_k|s_k) T(s'_k|a_k, s_k) R \left( \begin{array}{c|cc} s'_i & a_i & s_i \\ s'_k & a_k & s_k \end{array} \right) + V(s'_i, s'_k) \right) \right] \\
 &= \arg \max_{a_i} \sum_{s'_i \in G(s_i)} P(s'_i|a_i, s_i) \sum_{s'_k \in G(s_i)} \sum_{a_k} P(a_k|s_k) P(s'_k|a_k, s_k) R \left( \begin{array}{c|cc} s'_i & a_i & s_i \\ s'_k & a_k & s_k \end{array} \right) + V(s'_i, s'_k) \\
 &\quad \downarrow \quad \text{chain rule} \\
 &= \arg \max_{a_i} \sum_{s'_i \in G(s_i)} P(s'_i|a_i, s_i) \sum_{s'_k \in G(s_i)} \sum_{a_k} P(s'_k, a_k|s_k) R \left( \begin{array}{c|cc} s'_i & a_i & s_i \\ s'_k & a_k & s_k \end{array} \right) + V(s'_i, s'_k) \\
 &\quad \downarrow \quad \left( \text{using } \tilde{R}_k \left( \begin{array}{ccc} s'_i & a_i & s_i \\ & & s_k \end{array} \right) = E \left[ P(S_k, A_k|S_k) R \left( \begin{array}{ccc} S'_i & A_i & S_i \\ S'_k & A_k & S_k \end{array} \right) \right] + E[V(S_i, S_k)] \right)
 \end{aligned}$$

**Comment: in handwritten notes, this last equation had ambiguous bracketing. I couldn't resolve it to my satisfaction.**

$$\begin{aligned}
 &= \arg \max_{a_i} \sum_{s'_i \in G(s_i)} P(s'_i|a_i, s_i) \tilde{R}_k(s'_i, a_i, s_i, s_k) \\
 &\quad \downarrow \quad \left( \text{using } \tilde{R}(s'_i, a_i, s_i) \approx \tilde{R}(s'_i, a_i, s_i, s_k \sim S_k) \text{ where } s_k \sim S_k \text{ is a random sample from the state } S'_k \right) \\
 &= \arg \max_{a_i} \sum_{s'_i \in G(s_i)} P(s'_i|a_i, s_i) \tilde{R}_k(s'_i, a_i, s_i) \\
 &= \pi_i^*(s_i)
 \end{aligned}$$

$\therefore$  with computation of  $\tilde{R}_k(s_i, a_i, s_i, s_k)$  the estimation of  $\pi_i^*(s_i)$  is possible. The complexity to compute  $\tilde{R}_k(s_i|a_i, s_i)$  with perfect information is  $O(|S_k||A_k|)$ , which may be intractable to compute each time instant  $t$ .

Conditions:

1. stable stochastic  $\pi_k^*(s)$  or one that increases  $(R_k) \ E[R_{t+1}()] > E[R_t()]$
2. independent  $T(s'_k|a_k, s_k)$
3. a stably stochastic  $s_k \sim S_k$

asynchronous: assume  $a_k = \emptyset$ , always.

$$\begin{aligned}
\pi_i^*(s_i) &\sim \arg \max_{a_i} \sum_{s'_i \in G(s_i)} \left( T(s'_i|a_i, s_i) \sum_{s'_k} \sum_{a_k} \left( \pi(a_k|s_k) T(s'_k|a_k, s_k) \left( R \begin{pmatrix} s'_i & a_i & s_i \\ s'_k & a_k & s_k \end{pmatrix} + V(s'_i, s'_k) \right) \right) \right) \\
&\downarrow \\
&a_k = \emptyset \\
&\downarrow \\
&\sim \arg \max_{a_i} \sum_{s'_i \in G(s_i)} \left( T(s'_i|a_i, s_i) \sum_{s'_k} T(s'_k|s_k) \left( R \begin{pmatrix} s'_i & a_i & s_i \\ s'_k & \emptyset & s_k \end{pmatrix} + V(s'_i, s'_k) \right) \right) \\
&\downarrow \text{ since } a_k = \emptyset, a_k \neq \emptyset \quad s'_k, s_k \text{ decided based on } a_i, \text{ and stably stochastic} \\
&\sim \arg \max_{a_i} \sum_{s'_i \in G(s_i)} T(s'_i|a_i, s_i) \left( R \begin{pmatrix} s'_i & a_i & s_i \\ \bar{s}'_k & \emptyset & \bar{s}_k \end{pmatrix} + V(s'_i, \bar{s}'_k) \right), \\
&\quad \text{where } \left. \begin{array}{l} \bar{s}_i \sim S_i \\ \bar{s}'_k \sim S_k \end{array} \right\} \begin{array}{l} \text{randomly drawn} \\ \text{stably stochastic} \end{array} \\
&\bar{R}(s'_i, a_i, s_i) \leftarrow R \begin{pmatrix} s'_i & a_i & s_i \\ \bar{s}'_k & a_i & s_i \end{pmatrix} \\
&\quad \text{a) } \bar{s}_k, \bar{s}_i \text{ chosen as stably stochastic} \\
&\bar{R}(s'_i, a_i, s_i) \sim R(s'_i, a_i, s_i) \quad \text{when:} \quad \text{or} \\
&\quad \text{b) } \bar{s}_k \text{ chosen s.t. } E[R_{k,t+1}()] > E[R_{k,t}()]
\end{aligned}$$

Considering synchronous and asynchronous solutions, when  $\pi_i^*$  and  $\pi_k^*$  run independently, then both  $\bar{s}_k$  and  $\bar{s}_i$  will be chosen according to the “maximization choice”.

synchronous:  $\pi^*(s_i \cup s_k) \sim \pi_i^*(s_i) \cup \pi_k^*(s_k)$  iff  $\exists \tilde{R}_k(s'_i|a_i, s_i) \sim R_k(s'_i|a_i, s_k)$  which is computed each timestep with complexity  $O(|A_k||S_k|)$  and maybe intractable

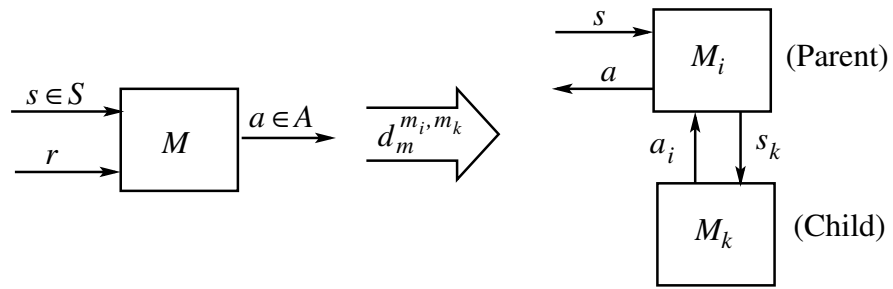
asynchronous:  $\pi^*(s_i \cup s_k) \sim \pi_i^*(s_i) \cup \pi_k^*(s_k)$  iff  $\bar{s}_k, \bar{s}_i$  are chosen according to the “maximum choice” which they trivially are by definition.

Splitting is done using a decomposition function

$$d_M^{M_i, M_k} = M \longrightarrow \left\{ M_i, M_k \left| \begin{array}{l} S_i \times (S_k / \{s_i\}) = S, S_k \times (S_i / \{s_k\}) = S \\ A = A_i \cup A_k \\ \tilde{T} \sim d^{-1}(d(\tilde{T})), d(\tilde{T}) = \tilde{T}_i, \tilde{T}_k \\ \tilde{R} \sim d^{-1}(d(\tilde{R})), d(\tilde{R}) = \tilde{R}_i, \tilde{R}_k \end{array} \right. \right\}$$

where  $d$  represents belief mapping functions that decompose and recompose  $\tilde{T}$  and  $\tilde{R}$ . This allows  $M$  to be mapped as new spaces and observations are encountered. The decomposition process breaks one MDP into a

parent and child:



Although many decomposition strategies are possible, this work presents a specific approach related to parent-child decomposition. This specific decomposition allows for assured policy convergence (see Section ??)

### Definitions

$$\underline{M} \quad S = (S_i/S_k) \times (S_k/S_i)$$

$$A = A_i \cup A_k$$

$$T = P(S \times A \times S)$$

$$R = \text{real, positive, convergent stochastic as } t \rightarrow \infty$$

$$R(s', a, s) = R \begin{pmatrix} s'_i & a_i & s_i \\ s'_k & a_k & s_k \end{pmatrix}$$

The system can be broken into the following MDP definitions

### $M_i$ – Parent

$S_i$  – a collection of states,  $s_i \in S_i$

$A_i$  – a collection of actions,  $a_i \in A_i$

$\tilde{T}(s'_i|s_i, a_i)$  – the observed probability of executing action  $a_i$  in state  $s_i$  and ending up in state  $s'_i$

$$\left. \begin{matrix} \tilde{T} \\ \tilde{R} \end{matrix} \right\} \text{ Covered Pages on BII p12-14}$$

$P(s'_i|s_i, a_i)$  is observed directly and is  $\tilde{T}(s'_i|s_i, a_i)$

$$\tilde{R}_t \left( \begin{matrix} s'_i \\ a'_k \end{matrix} \middle| \begin{matrix} a_i & s_i \\ a_k \end{matrix} \right) = \tilde{R}_t \left( \begin{matrix} s'_i \\ s'_k \end{matrix} \middle| \begin{matrix} a_i & s_i \\ a_k & s_k \end{matrix} \right) = \tilde{R}(s'|a, s)$$

$S, T_i, S_k, S'_k$  are not directly observable by process  $M_i$ , by design. Importantly, some facts are known about  $a_k$  and  $a'_k$  which will be exploited in Section ??.

ii)  $a_k = \pi_k(s_k)$

iii)  $a'_k = \pi_k(s_k)$

iiii)  $(s_k, s'_k)$  chosen indirectly by  $\pi_k(\cdot)$  in a manner that assuming monotonic increase in reward,

$$\text{as } t \rightarrow \infty \quad E[R_{t+1}(\cdot)] \geq E[R_t(\cdot)].$$

Local convergence of this process on a behavior policy  $\pi_i^*(\cdot)$  is assured.

### $M_k$ – child

$S'_i$  – all child states,  $s_k \in S_k$

$a_k \in A_k$

$$\tilde{T}(s'_k | s_k, a_k)$$

$$R_t(s'_k, a_k, s_k) = R_t \left( \begin{smallmatrix} s_i & a_i & s_i \\ s'_k & a_k & s_k \end{smallmatrix} \right) \text{ s.t. } s_i, s'_i \text{ are chosen by another process, and}$$

$$\boxed{A^*} \longrightarrow E[R_{t+1}(\cdot)] \geq E[R_t(\cdot)]$$

It is direct to note that both processes  $M_i$  and  $M_k$  are guaranteed to converge on locally optimal policies  $\pi_i(\cdot)$ ,  $\pi_k(\cdot)$  as the limit of time approaches infinity.

In following sections next steps are discussed:

- $\pi_i^* \times \pi_k^* \rightarrow \pi^*$  utilizing the degenerate-optimal policies (Section ??)
- $d_M^{M_i, M_k}$  – how to map values while preserving information, avoiding degeneracies (Section 2)

## 2 MAPPING FUNCTION

After defining Parent and Child MDPs, it is important to note how to map information from a centralized MDP onto a parent-child pair: (equation  $d_{M_i}^{M_j, M_k} d_{M_i}^{M_j, M_k}$ ). To expedite convergence and computational proofs that follow, a mapping function  $d_{M_i}^{M_j, M_k}$  is assume to have Basic Requirements (Sec. 2.1). Afterward the process is explained (Sec. 2.2).

### 2.1 Basic mapping requirements

$$S_i, S_j, S_k: S_j \times (S_k/S_j) \supseteq S_i, S_k \times (S_j/S_k) \supseteq S_i$$

Sets  $S_k$  and  $S_j$  must be able to be combined to make the original set.

$$A_i, A_j, A_k: A_i \subseteq A_j \cup A_k$$

Also, action sets must combine to recover the original action set.

$$T_i, R_i \sim \text{unknown/unknowable, stable decomposition}$$

The true transition and reward functions may remain unknown.

assumed
---------

more → \* important to select so that  $\tilde{T}_i$  &  $\tilde{T}_k$  seem independent

### 2.2 Mapping process

Transition mapping:

$$\begin{aligned} \exists f_1 : \tilde{T}_i &\rightarrow \tilde{T}_j, \tilde{T}_k, \text{invertible}; \tilde{T}_i = f_1 \left( f^{-1} \left( \tilde{T}_i \right) \right) \\ \exists f_2 : \tilde{R}_i &\rightarrow \tilde{R}_j, \tilde{R}_k, \text{invertible}; \tilde{R}_i = f_2 \left( f_2^{-1} \left( \tilde{R}_i \right) \right) \end{aligned}$$

Transition function mapping: knowing  $s_x \in S_x, s_y \in S_y$  (1 way)

$$\text{Goal } \exists f : P(S_x|A, S_x) \leftarrow P(S|A, S)$$

$$\text{knowing } P(S_x \times S_y|A_x \cup A_y, S_x \times S_y) = P(S|A, S)$$

Clearly:

$$P(s'_x|s_x, a_x) = \sum_{s'_y} \sum_{a_y} \sum_{s_y} P(s'_x, s'_y|s_x, s_y, a_x, a_y) P(a_y|S_y) P(s_y)$$

$$\therefore \tilde{T}(s'_x|s_x, a_x) = \sum_{s'_y} \sum_{a_y} \sum_{s_y} \tilde{T}(s'|a, s) \underbrace{\tilde{\pi}(a_y|s_y)}_{\text{require policy mapping}} P(s_y)$$

Action mapping:

Action mapping (1 way)

Next, we can consider an action mapping where actions from  $A$  can be randomly assigned to  $A_x, A_y$ :  $A_x \leftarrow \{a \in A' | A' \subseteq A\}$ ,  $A_x, A_y \subseteq A$ ,  $A_x \cup A_y = A$ ,  $A_x \neq \{\}$ ,  $A_y \neq \{\}$ .

General approach: High reward for ???ve states

Given  $\tau \in \mathbb{R}$ ,  $\pi(a|s)$ ,  $\tilde{Q}(a, s)$  then

$$A_x \leftarrow A_x \cup \left\{ a \left| \underbrace{\pi(a|s)\tilde{Q}(a, s)}_{\text{condition}} > \tau \right. \right\}$$

or, more usefully/generally

$$A_x \leftarrow A_x \cup \left\{ a \left| \underbrace{\left( \pi(a|S_s^{R'}) \right) \tilde{Q}(a, S^{*'})}_{\text{condition}} > \tau \right. \right\}$$

where

$$S_s^{*'} = \{S' | S/s_s^* \neq S'\} \quad (\text{see p. ???})$$

Condition options:

$$\begin{array}{l} \text{*reformulation over set } S \text{ vs. } s \in S \end{array} \quad \left\{ \begin{array}{ll} \text{a) } \pi(a|s)\tilde{Q}(a, s) > \tau & \dots \text{ High reward} \\ \text{b) } \pi(a|s)\tilde{Q}(a, s) > \tau, \quad \pi(a|s) > 0 & \dots \text{ small reward} \end{array} \right.$$

(Network approach)



### 3 MAPPING AS A PROCESS

In this section it is described how, and when, the function  $d_{M_i}^{M_j, M_k}$  is applied. In fact, instead of structing the application of  $d_{M_i}^{M_j, M_k}$  as human chosen conditional statements, it is possible to consider when to “split” and “merge” and MDP  $M_i$  into  $M_j, M_k$  as a learning problem. This problem is expressed as the mapping problem, and modelled as a learning problem MDP-map ( $M_{\text{map}}$ ):

$$M_{\text{map}}^i = \langle S_{\text{map}}^i, A_{\text{map}}^i, T_{\text{map}}^i, R_{\text{map}}^i \rangle$$

$$S_{\text{map}} = f(d(M_i, \mathcal{P}(M_j), P(M_k))) , \quad \text{s.t. } f : d_{M_i}^{M_j, M_k} \rightarrow \mathbb{R}$$

where  $d(M_i, M_j, M_k) = d_{M_i}^{M_j, M_k}$ .

The state encodes all possible configurations for the MDP split:

[FIGURE]

$A_{\text{map}}^i$  = a set of actions formed from the the decomposition and recomposition functions

$$A_{\text{decomposition}} \cup \{a_r\}$$

where  $a_r$  is a special recomposition function.

$$T_{\text{map}}^i = \begin{cases} 1, & (S = S_1 \text{ and } a \neq a_r) \text{ or } (S \neq S_1 \text{ and } a = a_r) \\ 0, & \text{otherwise} \end{cases}$$

$$R_{\text{map}}^i = \sum_{e=\text{epoch}} R^i(e)$$

Given  $(S_{\text{map}}, A_{\text{map}}, T_{\text{map}}, R_{\text{map}}^i, \pi_{\text{map}})$ , applied to  $M = \langle S_x, A_x, \tilde{T}_x, R_x, \tilde{\pi}_x \rangle$ , we may trivially define  $M_y = \langle S_y, A_y, \tilde{T}_y, R_y, \pi_y \rangle$  in a method consistent with Bush, p. 74, with  $M_x$  being the parent process and  $M_y$  being the child.

- a) for  $R_{\text{map}}^i$ , there are five versions  $i \in \{1, \dots, 5\}$
- b) Given  $M_x, M_y$ , a merge is also possible, so recovering  $M$
- c) we can perform temporal Sink actions on an MDP (Book I, p. 88)
  - $\hookrightarrow$  reduce resolution
  - $\hookrightarrow$  re-increase resolution

#### Actions

$\therefore$  Seven “actions” can be performed on an MDP: ( $M_R$ )

State

$$\left\{R_{\text{map}}^i\right\}_{i=0}^5 \cup \{\text{merge}\} \times \left\{\text{scale up}^i\right\}_{i=0}^5 \cup \{\text{normal}\} \cup \{\emptyset\}$$

Reward

$$R(s', a, s) = \sum_{l \in e} R(l) \quad \text{reward during a trajectory}$$

$e$  = epoch

Transition

-easy to explain in MS Word

$$T = \begin{cases} 1 - \text{allow ???} \\ 0 - \text{otherwise} \end{cases}$$

Reward function mapping (5 way)

knowing  $R(\{s_x, s_y\} | a, \{s'_x, s'_y\}) = R(s | a, s)$

Mapping Policy (1 way)

finding:  $f\tilde{\pi}(s) \rightarrow f\tilde{\pi}(s_y)$

$$\tilde{\pi}(a_y | s_y) \leftarrow \sum_{s_x} \sum_{a_x} \tilde{\pi}(\{a_x, a_y\} | \{s_x, s_y\}) P(s_x)$$

### 3.1

#### 3.1.1 Proof of Decomposition Correctness (Process 1)

The following section demonstrates

#### MDP Policy Decomposition

First, it is true that the optimal policy  $\pi^*(s)$  for a centralized MDP will arise from maximizing the expected reward.

$$\pi^*(s) = \arg \max_a \sum_{s'} R(s, a, s') P(s'|s, a) + \gamma V(s')$$

It is now shown that a set of MDPs  $(M_i, M_k)$  can be observed to create an optimal policy  $\exists f.s.t. f(\pi_i^*, \pi_k^*) \sim \pi^*$ .

Given  $\pi_i^*(S_i, A_k), \pi_k^*(S_k)$

$$1. \quad \pi^*(s_i, s_k) = \arg \max_{a_i, a_k} \sum_{s'_i} \sum_{s'_k} R((s_i, s_k), (a_i, a_k), (s'_i, s'_k)) P((s'_i, s'_k)|(s_i, s_k), (a_i, a_k))$$

\* Lemma 1 – augmentation with  $a_k$  where  $a'_k = \pi^*(s'_k)$

$$2. \quad \pi^*(s_i, s_k) = \arg \max_{a_i, a_k} \sum_{s'_i} \sum_{s'_k} R((s_i, s_k, a_k), (a_i, a_k), (s'_i, s'_k, a'_k)) P((s'_i, s'_k, a'_k)|(s_i, s_k), (a_i, a_k))$$

\* Lemma 2 – Simplification

$$\overbrace{\arg \max_{a_i, a_k} \equiv \arg \max_{a_i} \arg \max_{a_k}}^{\text{separability}}$$

$$3. \quad \pi^*(s_i, s_k) = \arg \max_{a_i, a_k} \sum_{s'_i} R((s_i, a_k), (a_i, a_k), (s'_i, a'_k)) P(s'_i, a'_k|a_i, a_k, (s_i, s_k))$$

#### Lemma 3

\* separation of  $a_k$ , and  $a_k \leftarrow \pi_k^*(s_k)$

$$4. \quad \pi^*(s_i, s_k) = \left( \arg \max_{a_i} \sum_{s'_i} R((s_i, a_k), a_i, (s'_i, a'_k)) P(s'_i, a'_k|a_i, (s_i, s_k)) \right)$$

\* Lemma 4

$$a_i = \pi_i^*(s_i) \longrightarrow \cup \left( \arg \max_{a_k} \sum_{s'_k} R(s_k, a_k, s'_k, a'_k) P(s'_k|a_k, s_k) \right)$$

$$5. \quad \pi^*(s_i, s_k) = \pi_i^*(s_i, a_k) \cup \pi^*(s_k)$$

Thus, given  $\pi_i^*$  and  $\pi_{k'}^i$ , it is possible to infer  $\pi^*(s_i, s_k)$ .

We now prove lemmas involved.

Lemma 1:

Lemma 2:

Lemma 3:

Lemma 4:



### 3.1.2 Proof of Decomposition Correctness (Process #2)

#### MDP Policy Decomposition

$$\pi^*(s) = \arg \max_a \sum_{s'} R(s, a, s') P(s'|s, a) + \gamma V(s)$$

Given  $\pi_i^*(S_i, A_k), \pi^*(S_k)$

$$1. \quad \pi^*(s_i, s_k) = \arg \max_{a_i, a_k} \sum_{s'_i} \sum_{s'_k} R((s_i, s_k), (a_i, a_k), (s'_i, s'_k)) P((s'_i, s'_k)|(s_i, s_k), (a_i, a_k))$$

Note:  $R((s_i, s_k, a_k), (a_i, a_k), (s'_i, s'_k)) \leftarrow R((s_i, s_k), (a_i, a_k), (s'_i, s'_k))$

$R((s_i, s_k, a_k), (a_i, a_k), (s'_i, s'_k, a'_k)) \leftarrow R((s_i, s_k), (a_i, a_k), (s'_i, s'_k))$

\*assume separability  $\longrightarrow$

$$2. \quad \pi^*(s_i, s_k) = \arg \max_{a_i} \arg \max_{a_k} \sum_{s'_i} \sum_{s'_k} R((s_i, s_k, a_k), (a_i, a_k), (s'_i, s'_k, a'_k)) P((s'_i, s'_k, a'_k)|\cdot)$$

$$= \arg \max_{a_i} \sum_{s'_i} \sum_{s'_k} R \left( (s_i, s_k, a_k), \begin{array}{c|c} a_i & s'_i \\ s'_k & \\ \hline a'_k & \end{array} \right) P \left( \begin{array}{c|c} s'_i & s_i \\ s'_k & a_i, s_k \\ \hline a'_k & a_k \end{array} \right) \\ \cup \arg \max_{a_k} \sum_{s'_k} R \left( \begin{array}{c|c} s_i & a_i \\ s_k & a_k \\ \hline a_k & \end{array}, \begin{array}{c|c} s'_i & \\ s'_k & \\ \hline a'_k & \end{array} \right) P \left( \begin{array}{c|c} s'_i & s_i \\ s'_k & a_i, s_k \\ \hline a'_k & a_k \end{array} \right)$$

$$* \quad \text{let } a_k^* = \arg \max_{a_k} \sum_{s'_k} R \left( \begin{array}{c|c} s_i & a_i \\ s_k & a_k \\ \hline a_k & \end{array}, \begin{array}{c|c} s'_i & \\ s'_k & \\ \hline a'_k & \end{array} \right) P \left( \begin{array}{c|c} s'_i & s_i \\ s_k & a_i, s_k \\ \hline a_k & a_k \end{array} \right)$$

$$\text{s. t. } s_i, s'_i \leftarrow \pi_i^*($$

$$a_i^* = \pi_i^*(s)$$

$$3. \quad = \arg \max_{a_i} \sum_{s'_i} R \left( \begin{array}{c|c|c} s_i & a_k & s'_i \\ a_k & a_k^* & a'_k \end{array} \right) P \left( \begin{array}{c|c} s'_i & a_i, s_i \\ a'_k & a_k^*, a_k \end{array} \right) \cup \pi_k^*(s_k) = a_k$$

$$= \pi_i^*(s_i|\pi_k^*) \cup \pi_k^*(s_k)$$

## 4 CONVERGENCE

Policy convergence is considered for the parent and child learning process. Since both policies utilize information from external learning processes the proofs are non-trivial.

### 4.1 For the Parent MDP $M_k$

1. Assume the child's expected return has a monotonic expected value.

$$E[R_t(s'_k|a_k, s_k)] \geq E[R_{t+1}(s'_k|a_k, s_k)]$$

2. Define reward as

$$R_t(s'_i|a'_i, s'_i) = R \left( \begin{array}{c|cc} s'_i & a_i & s_i \\ s'_k & a_k & s_k \end{array} \right)$$

$$\text{i. } a_k = \pi_k(s_k)$$

We de-

$$(s'_k, s_k) \text{ result from } a_i \text{ s.t.}$$

$$\text{ii. } s'_k \sim T(S_k|\pi_i(s_k), s_k)$$

$$\text{iii. } s_k \sim T(S_k|\pi_i(s_k^*), s_k^*)$$

fine the reward of the parent to be the reward received by the child process during observation ( $M_i$  observing  $M_k$ ).

- A) Parent is effective

$$0 > \min_{a_i} E \left[ R \left( \begin{array}{c|cc} s'_i & a_i & s_i \\ s'_k & a_k^* & s_k \end{array} \right) \right] - E \left[ R \left( \begin{array}{c|cc} s'_i & A_i & s_i \\ s'_k & A_k & s_k \end{array} \right) \right]$$

Child behaviour policy  $\pi_k$  is effective if the expected reward increases. In this case we consider  $a_k$  to be optimal.

- B) Child is convergent

$$E \left[ R_{t+1} \left( \begin{array}{c|cc} \cdot & a_i & \cdot \\ \cdot & \pi_{k,t+1} & \cdot \end{array} \right) \right] \geq E \left[ R_{t+1} \left( \begin{array}{c|cc} \cdot & a_i & \cdot \\ \cdot & \pi_{k,t} & \cdot \end{array} \right) \right]$$

It is convergent if the centralized MDP's reward increases when observed.

assume some policy  $\pi_k(\cdot)$  is both effective and convergent, then:

$$\textcircled{A} + \textcircled{B} \rightarrow \textcircled{C}$$

- C) The system must be convergent

- The parent chooses actions to maximize  $R_t(A)$
- The child chooses actions to maximize  $R_t(B)$

$$\boxed{*} \quad \lim_{t \rightarrow \infty} R_t(s'_i|a_i, s_i) \sim R_{t+1}(s'_i|a_i, s_i)$$



## 4.2 For the child

\* trivial

1. Assume the child's expected return has a monotonic expected value.

$$E[R_t(s'_k|a_k, s_k)] \geq E[R_{t+1}(s'_k|a_k, s_k)]$$

2. Define reward as

$$R_t(s'_i|a'_i, s'_i) = R \left( \begin{array}{c|cc} s'_i & a_i & s_i \\ s'_k & a_k & s_k \end{array} \right)$$

$$\text{i. } a_k = \pi_k(s_k)$$

We de-

$(s'_k, s_k)$  result from  $a_i$  s.t.

$$\text{ii. } s'_k \sim T(S_k|\pi_i(s_k), s_k)$$

$$\text{iii. } s_k \sim T(S_k|\pi_i(s_k^*), s_k^*)$$

fine the reward of the parent to be the reward received by the child process during observation ( $M_i$  observing  $M_k$ ).

- A) Parent is effective

$$0 > \min_{a_i} E \left[ R \left( \begin{array}{c|cc} s'_i & a_i & s_i \\ s'_k & a_k^* & s_k \end{array} \right) \right] - E \left[ R \left( \begin{array}{c|cc} s'_i & A_i & s_i \\ s'_k & A_k & s_k \end{array} \right) \right]$$

Child behaviour policy  $\pi_k$  is effective if the expected reward increases. In this case we consider  $a_k$  to be optimal.

- B) Child is convergent

$$E \left[ R_{t+1} \left( \begin{array}{c|cc} \cdot & a_i & \cdot \\ \cdot & \pi_{k_{t+1}} & \cdot \end{array} \right) \right] \geq E \left[ R_{t+1} \left( \begin{array}{c|cc} \cdot & a_i & \cdot \\ \cdot & \pi_{k_{t+1}} & \cdot \end{array} \right) \right]$$

It is convergent if the centralized MDP's reward increases when observed.

assume some policy  $\pi_k(\cdot)$  is both effective and convergent, then:

$$\textcircled{A} + \textcircled{B} \rightarrow \textcircled{C}$$

- C) The system must be convergent

- The parent chooses actions to maximize  $R_t(A)$
- The child chooses actions to maximize  $R_t(B)$

$$* \quad \lim_{t \rightarrow \infty} R_t(s'_i|a_i, s_i) \sim R_{t+1}(s'_i|a_i, s_i)$$

## 5 WORST CASE PERFORMANCE

### 5.1 Computational Complexity

Total Mapping

$$* A_R = \left\{ \begin{array}{l} \text{State 1, State 2, Action 1, Action 2} \\ \text{merge, time up, time down} \end{array} \right\}$$

Given  $S \in \mathbb{R}^n$ , define dimensions  $\{i_s\}_{i_s=1}^n$

$A \in \mathbb{N}^m$ , define dimensions  $\{i_r\}_{i_r=1}^m$

Then, with an initial MDP  $M = \langle S, A, T, R, \pi, M_R \rangle$ , all possible “sub mdps”  $M_1, M_2, M_3, \dots$  represent the family of MDPs which can be created from  $M$ ,  $\mathcal{P}(M) = \{M_x | S_x \subseteq S, A_x \subseteq A, R, \pi \text{ from MDPs } ???\}$  and each member  $M_x$  is characterized by a language  $I_{sx} \subseteq \{i_s\}_{i_s=1}^n$  or  $I_{sy} \subseteq \{i_r\}_{i_r=1}^m$  where  $I_{sx} \times I_{sy}$  defines a space  $S_R$ , for the reconfiguration MDP to explore, with actions from  $A_R$ .

$$I \left| \begin{array}{l} \text{Reward is defined as average expected reward over an epoch } e. \\ \text{in terms of transition} \end{array} \right.$$

$$* S_R = \mathcal{P}\left(\{i_s\}_{i_s=1}^n\right) \times \mathcal{P}\left(\{i_r\}_{i_r=1}^m\right) \quad \longleftarrow \text{exponential increase in space (stupid!)}$$

Problems

- 1) exponential space consumption
- 2) how to handle chaining/nesting
- 3) how to structure action choice policy

## 6 NN

① Summarize Regression ( $y = mx + b$ )

points  $\{1, 2, 3, \dots\}$

- easiest starting point is choosing “the best line”
- need a metric to define this idea of “best” mathematically
- describe the line:  $f(x_1) = m(x_1) + b$

$$m = \frac{y_2 - y_1}{x_2 - x_1}$$

idea — choose line closest to points by taking sum of differences,

$$\text{minimize } \sum_i e_i, \quad e_i = y_i - f(x_i) = y_i - m_i x_i + b m_0 x_0, \quad x_0 = 1 \quad (\text{by convention})$$

- apply square to remove negative/positive problem

$$e_i = (y_i - m_i x_i + b)^2$$

$$e_i = \left( \sum_{d=0}^D y_{id} - m_d x_{id} \right)^2$$

look up linear algebra proof

## Logistic Regression Summary

 $\mathbb{A}_i$ 

## Logistic sigmoid

- not vulnerable to being ??? fit
- (other solutions)
- good for classification
- idea  $f : \mathbb{R} \rightarrow (0, 1)$
- good for probability then!

$$\sigma(a) = \frac{1}{1 + e^{-a}}$$

## Probability review

- Probability density
- CDF

$$\text{density} = P(x_1 < x < x_2)$$

$$\text{CDF} = \int_{x_1}^{x_2} P(X = x) dx$$

$$\text{CDF} \approx \sum_{x=0}^{x_2} P(X = x)$$

$$P(X = x|Y = y) = \frac{P(X=x, Y=y)}{P(Y=y)}$$

$$P(Y = y) = \sum_{x \in X} P(X = x, Y = y)$$

define  $P(x|c_1), P(x, c_2) P(c_1), P(c_2)$

Likelihood class

a fun experiment

$$\begin{aligned}
 &= \frac{P(x|c_1)P(c_1)}{P(x|c_1)P(c_1) + P(x|c_2)P(c_2)} \\
 &= \frac{P(x|c_1)P(c_1)}{P(x|c_2)P(c_2) + P(x|c_1)P(c_1)} \\
 &= \frac{B}{A+B} \\
 &= \frac{\frac{B}{B}}{\frac{A+B}{B}} \\
 &= \frac{1}{\frac{A}{B} + 1} \\
 &= \frac{1}{1 + \exp\left(\ln\left(\frac{A}{B}\right)\right)}
 \end{aligned}$$

$$\begin{aligned}
&= \frac{1}{1 + \exp(-\ln(\frac{B}{A}))} \\
&= \sigma(a) \\
&= \frac{1}{e^{-a}}
\end{aligned}$$

log ratio:

$$\begin{aligned}
a &= \ln \frac{P(x|c_1)P(c_1)}{P(x|c_2)P(c_2)} \\
a &= -\ln \frac{B}{A}
\end{aligned}$$

Given  $x$ , how to decide between  $c_1$  and  $c_2$ ?

① Be racist (prejudice)

$$\frac{P(c_1)}{P(c_2)} \quad \text{— odds of observing } c_1 \text{ or } c_2 \text{ in nature}$$

② racism light:

$$P(c_1|x) = \frac{P(x|c_1)}{P(x|c_1) + P(x|c_2)}$$

— odds of observing  $x$  in  $c_1$  situation

— odds of observing  $x$  in  $c_2$  situation

$$P(c_1|x) = \frac{P(x|c_1)}{\sum_{c \in \{c_1, c_2\}} P(x|c)}$$

③ No prejudice:

— Latent variable explanation

forms for  $P(x|c_k) \dots$  First, count  $P(x|c_k) = \frac{\sum_{c_k \in T} 1}{\sum_{c_k} 1}$   $T = \text{true set}$

impractical as  $\lim_{|T| \rightarrow \infty} \Omega(P(x|c_k)) = \infty$

—  $P(x|c_k)$  — parametric, try Gaussian (scales to  $|T| \rightarrow \infty$ , and realistic in ???)

$= P(x|c_k) = \mathcal{N}(x|c_k, M_k, \Sigma)$ , s.t.  $u_k \sim c_k$

$$P(x|c_k) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} \exp \left\{ -\frac{1}{2} (x - M_k)^T \Sigma^{-1} (x - M_k) \right\}$$

$D$  – dimension

Some derivation work

-----

$$\begin{aligned} a &= \ln \left( \frac{P(x|c_k)}{P(x|c_2)} \right) + \ln \left( \frac{P(c_k)}{P(c_2)} \right) \\ a &= \ln \left( \frac{\exp \left\{ -\frac{1}{2} (x - M_k)^T \Sigma^{-1} (x - M_k) \right\}}{\exp \left\{ -\frac{1}{2} (x - M_2)^T \Sigma^{-1} (x - M_2) \right\}} \right) + \ln \left( \frac{P(c_k)}{P(c_2)} \right) \end{aligned}$$

↓ magic (quadratic terms cancel due to ???  $\Sigma$  matrix)

$$\left( \Sigma^{-1} (u_1 - u_2)^T \right) x - \frac{1}{2} u_k^T \Sigma u_k + \frac{1}{2} u_k^T \Sigma^{-1} u_{k2} + \ln \left( \frac{P(c_k)}{P(c_2)} \right)$$

or  $a_k = \ln(P(x|c_k)P(c_k))$  is an independent Gaussian:

## B) Multiclass logistic regression

$$P(c_k|\phi) = \frac{e^{a_k}}{\sum_j e^{a_j}}$$

discriminate

$$a_k = w_k^T \phi$$

Generative version

$$a_k = \ln \frac{P(x|c_k)P(c_k)}{\sum_j P(x|c_j)P(c_j)}$$

a lot of computation ("fuck it")

for multiclass check out

$$a_k \approx w_k^T \phi$$

$$\frac{\partial P(c_k|\phi)}{\partial a_j} = P(c_k|\phi) (1 - P(c_j|\phi))$$

$$\frac{\partial P(c_k|\phi)}{\partial a_j} = P(c_k|\phi) \left( 1 - \frac{e^{a_j}}{\sum_l e^{a_l}} \right)$$

$$= \frac{e^{a_k}}{\sum_l e^{a_l}} \left( 1 - \frac{e^{a_j}}{\sum_l e^{a_l}} \right)$$

$$= \frac{e^{a_k}}{\sum_l e^{a_l}} - \frac{e^{a_k} e^{a_j}}{\sum_l (e^{a_l})^2}$$

$$= P(c_k|\phi) - P(c_k|c_j, \phi)$$

$$= P(c_k, c_j|\phi) \quad \leftarrow c_k \text{ without } c_j \text{ included (just } c_k)$$

$$\frac{\partial P(c_k|\phi)}{\partial a_j} = P(c_k|\phi) (I_{kj} - P(c_j|\phi))$$

# Feed Forward Neural Netowrks

$$P(c_k|\phi) = \sigma(a_k)$$

change of terms

$$P(c_k|\phi) \rightarrow y_k$$

$$c_k \sim w$$

$$y(x, w) = f \left( \sum_{j=1}^M w_j \phi_j(x) \right)$$

$$(\cdot) \sim \text{layer} \quad y(x, w) = \sum_{j=1}^M w_j \sigma(x) \quad y_k = \sigma(a_k)$$

$$\left( \sum_j \right) Y = \begin{bmatrix} P(c_1|\phi) \\ \vdots \\ P(c_m|\phi) \end{bmatrix}$$

$$a_j^{(x)} = w_k^T X = \sum_{i=0}^D w_{ji}^{(1)} + w_{j0}^{(1)} \quad (\text{intercept term from gaussian})$$

$$a_j = \sum_{i=0}^D w_{ji}^{(1)} x_i + w_{j0}^{(1)}$$

$$P(c_j|x) = \sigma(a_j)$$

$$y_k = \sigma(a_j)$$

$$y_j(x, w) = \sigma \left( \sum_{j=1}^M w_j^{(1)} x_i + w_{j0}^{(1)} \right)$$

↓

$$y_j(x, w) = \sigma \left( \sum_{j=1}^M w_j \phi_j(x) + w_{j0} \right), \quad \text{s.t. } \phi_j(x) = x_i$$

$$y_j(x, w) = \sigma \left( \sum_{j=1}^M w_j \phi_j^{(1)}(x) + w_{j0} \right) \quad \text{s.t. } \begin{cases} \phi_j^{(i)}(x) = \sum_{j=1}^M w_j \phi_j^{(i)}(x) + w_{j0} \\ \phi_j^{i=M}(x) = x_j \end{cases}$$



## Training

### Error function:

$$E(w) = \frac{1}{2} \sum_{n=1}^N \|y(x_n, w) - t_n\|^2$$

view as probability:  $P(t|x, w) = \mathcal{N}(t|y(x, w), \beta^{-1})$

Do

$$P(t|X, W, \beta) = \prod_{n=1}^N P(t_n|x_n, w, \beta)$$

↓

$$-\ln(P(t|x, w, \beta)) = \frac{\beta}{2} \sum \{y(x_n, w) - t_n\}^2 - \frac{N}{2} \ln \beta + \frac{N}{2} \ln(2\pi)$$

$$\therefore E(w) = \frac{1}{2} \sum_{n=1}^N \{y(x_n, w) - t_n\}^2$$

Probabilistic treatment  $\approx$  Sum of squared errors

1) find  $w_{mL}$  (gradient descent)

then  $\frac{1}{\beta_{mL}} = \frac{1}{N} \sum_{n=1}^N \{y(x_n, w_{mL})\}^2$

for output:

$$\frac{\partial E}{\partial a_k} = y_k - t_k \quad (a_k \approx y_k \text{ in output units})$$

(only  $k$  varies)

$$\begin{aligned} \left( \frac{\partial E_w}{\partial a_k} \right) &= \frac{1}{2} \frac{\partial \{y(x_k, w) - t_k\}^2}{\partial a_k} \\ &= \frac{2}{2} \{y(x_k, w) - t_k\} \left( \frac{\partial y_k(x_k, w)}{\partial a_k} - 0 \right) \\ &= y(x_k, w) - t_k \\ &= y_k - t_k \end{aligned}$$

2) Back Propagation

Ⓐ weight derivative (to do with  $dE$ )

Ⓑ weight adjusted (application of  $dE$ )

Childishly:

$$E(w) = \sum_{n=1}^N E_n(w)$$

$$y_k = w_{ki} x_i$$

$$E_n = \frac{1}{2} \sum_k (y_{nk} - t_{nk})^2 \quad \text{---(input pattern } n\text{)}$$

$$y_{nk} = y_k(x_n, w)$$

$$\frac{\partial E}{\partial w_{ji}} = (y_{nj} - t_{nj}) x_{ni}$$

(callback to Logistic Regression)

$$a_j = \sum_i w_{ji} z_i, \quad z_j = h(a_j), \quad y(\cdot)$$

Chain rule:

$$\frac{\partial E_n}{\partial w_{ji}} = \frac{\partial E_n}{\partial a_j} \frac{\partial a_j}{\partial w_{ji}}$$

or

$$\frac{\partial E_n}{\partial w_{ji}} = \delta_j \frac{\partial a_j}{\partial w_{ji}}$$

$$\frac{\partial E_n}{\partial w_{ji}} = \delta_j z_i$$

for outputs

$$\frac{\partial E_n}{\partial w_{ji}} = (y_k - t_k) z_i$$

for hidden

$$\frac{\partial E_n}{\partial a_j} = \sum_k \frac{\partial E_n}{\partial a_k} \frac{\partial a_k}{\partial a_j}$$

$$\delta_j = h'(a_j) \sum_k w_{kj} \delta_k$$

Temporal Differences

(This is a type of reconfiguration)

(see Book 2 p 16 for rewrite)

Paper: On Static and Temporal Difference for MDPs (extends Reconfigurable ???, Page 88)

want to mix a trajectory  $\{l\}_{i=0}^{\infty}$  with ???  $l_{i=0}^{\infty}$ , s.t.  $l_i = \{s_i, a_i, s'_i\}$ ,  $l_i \in L$

Given a trajectory,  $l_i^j$ , we want a compressed representation.

$$F_c(r_c) = \left\{ f_c : L^n \times j \rightarrow L \mid \text{rel} \left[ r_c \left( R(l_{i_1}^{i_2}) \right), r_c \left( R(l_{i_2}^{i_3}) \right) \right] = \text{rel} \left[ R \left( f_c(l_{i_1}^{i_2}) \right), R \left( f_c(l_{i_2}^{i_3}) \right) \right] \right\}$$

s.t.  $f_c(l_{i_1}^{i_2}) = l_{c_0}^{c_1} = l_{c_0}$

and a reward compressor

$$R_c = \{r_c : R(L^n) \times j \rightarrow R(L)\}^{\infty}, \quad A_c = \{f_{a_c} : L \rightarrow \{\pi_c\}^{\infty}\}$$

$$T_c(l_{c_0}, a_c \in A_c, l_{c_1})$$

$$\begin{aligned} T(s_1, \pi_c, s_2) &= \sum_{a \in \pi_c(s_1)} \pi_c(a|s_1) T(s_1, a, s_2) \\ &= \sum_{a \in \pi_c(s_1)} P(a|s_1) P(s_2|s_1, a) \\ &= \sum_{a \in \pi(s_1)} P(a, s_2|s_1) \\ &= P(s_2|s_1, \pi_c) \\ &= T(s_1, \pi_c, s_1) \end{aligned}$$

$$T(s_1, \pi_c, s_3) = \sum_{s_i \in \tilde{S}} T(s_1, \pi_c, s_i) T(s_i, \pi_c, s_2) \quad \leftarrow \text{[should this be } s_3?]$$

$$T(s_1, \pi_c, s_3) = T_{\pi_c}(s_1, \tilde{s}_2) \cdot T_{\pi_c}(\tilde{s}_2, s_3)$$

$$T(s_1, \pi_c, s_n) = T_{\pi_c}(s_1, \tilde{s}_2) \left( \prod_{i=2}^{n-2} T_{\pi_c}(\tilde{s}_i, \tilde{s}_{i+1}) \right) T_{\pi_c}(\tilde{s}_{n-1}, s_n)$$

$$T(s, f_{a_c}(l), s_n) = \prod_{i=1}^{n-1} T_{\pi_c}(\tilde{s}_i, \tilde{s}_{i+1}), \quad \tilde{s}_1 = \{s_1\}, \tilde{s}_n = \{s_n\}$$

(might be intractable to learn but can guess by  
(from page 85:) → link in and rewrite Both.

reconfiguration ideas:

- ☐  $g(\cdot)$  should be a learned output function (regressed)
- ☐ the states of nodes should be able to be remapped to outputs of other MDPs
- ☐ actions can be large, multi-node, set-outputs as opposed to single, discrete or continuous actions

### Transitional Learning (Idea)

→ The generalization of Q-Learning to a theory relates to learning transitional models

assume a process  $\langle S, A, T, R \rangle$

knowing

$$\textcircled{1}: Q_t(s, a) = Q_{t-1}(s, a) + \alpha [R(s, a) + \gamma \max_{a^*} Q_{t-1}(s', a^*)]$$

$$\textcircled{2}: E[R(s, a)] = \gamma \sum_{s' \in S} T(s, a, s') \max_{a^*} E[R(s', a)]$$

$$\textcircled{3}: \lim_{t \rightarrow \infty} Q_t(s, a) = E[R(s, a)]$$

which is the general method  $\textcircled{1}$  used to converge on expected reward given in  $\textcircled{2}$ . We believe  $\textcircled{3}$ .

Turns out (I think) we can do better, replacing Q-Learning with transitional learning by making three substitutions.

$\textcircled{4}$  use  $P_t(s'|s, a)$  as a replacement for  $T(s, a, s')$

$\textcircled{5}$  limit recursion depth using a threshold  $A_{\text{th}}$

$\textcircled{6}$  model an approximation for  $S$  called  $L(s)$ , s.t.  $L(s) \subset S$ , with a basis characterized by  $\textcircled{4}$

Thus we can define an approximation for  $\textcircled{2}$ :

$$\bar{R}_\gamma(s, a|t, \theta) = \begin{cases} \sum_{s' \in L(s)} P_t(s'|s, a) \max_{a^*} \bar{R}_\gamma(s, a|t, \theta\gamma), & \theta\gamma > A_{\text{th}} \\ 0, & \text{otherwise} \end{cases}$$

also:

$$\lim_{A_{\text{th}} \rightarrow 0} \lim_{t \rightarrow \gamma} \bar{R}_\gamma(s, a|t, \theta) = E[R(s, a)] \quad \text{if } T(s, a, s') = P_t(s'|s, a)$$

Lastly,  $L(s)$  and  $P_t(s'|s, a)$  can be considered.  $\{s, s', w\} \in S \times S \times \mathbb{R}$  characterizes

$$L(s) = \{s' \in \{s, s', w\} \mid w > \text{th } B\}$$

## Financial Domain

$$S = \text{Security} = \{V_S, P_S\} \quad V_s = \{V_t\}_{t=0}^N, P_s = \{P_t\}_{t=0}^N$$

Goal:  $V_{N+1}^*, P_{N+1}^*$

and, a set of Securities =  $\{S_i\}_{i=0}^M$

approaches

- predict future prices
- predict direction

## Build Bayesian model

feature vector space  $d = \{1, 0, -1\}$  (direction)

- 1) choose stocks with
  - a) high beta
  - b) high volume
- 2) find:  $P(d | \{d\}_{N-\text{frame}}^N)$  ?  $\dots$  assume independence  
 frame = 4  $\rightarrow$  81 states  
 1-week of data:  
 2400 states for training  
 4800 states ( $\sim$  60 in each state)
- 3) build  $f_r(d | \{d\}_{t=N-\text{frame}}^N)$   $\mathcal{O}(4800)$
- 4) run  $P(d|D) = \frac{f_r(d|D)}{\sum_{d' \neq d} f_r(d'|D)}$

## Testing

- a) Dump  $f_r$  for all 81 states
- b) Dump High/Low  $P(d|D)$

Naive

### Probabilistic Trading

- 1)  $\forall S_i \in M$ : calculate  $f_r^i(x_t | \overbrace{x_{t-1}, x_{t-2}, x_{t-3}, x_{t-4}}^{h_{t-1}})$   
 $f_r^i(x | h_{t-1}) = \sum_{t=\text{start\_history}}^T I \left\{ \sum_{t'=t}^{\text{frame}+t} I \{d_{t'} = x_{t'}\} = 4 \right\}$
- 2)  $P^i(x | h_t - 1) \leftarrow \frac{f_r^i(x | h_{t-1})}{\sum_{x'} f_r^i(x' | h_{t-1})}$   
 $z_t^i = I \left\{ \sum_{t'=t}^{\text{fs}} I \{d_{t'} = x_{t'}\} = \text{fs} \right\}$   
 $z_t^i \in \{0, 1\}$

### Probabilistic Market ???

Given a family  $F = \{S^i, S^j, \dots, S^k\}, |F| \in \mathbb{N}$

$$f_r^F(x | h_{t-1}) = \sum_{S^i \in F} \sum_{t=\text{start\_present}} I \left\{ \sum_{t'=t} \{d_{t'}^i = x_{t'}^i\} = \text{frame} \right\}$$

$\text{start\_history} \ll \text{start\_present}$

$$P^F(x | h_{t-1}) \leftarrow \frac{f_r^F(x | h_{t-1})}{\sum_{x'} f_r^F(x' | h_{t-1})}$$

### Probabilistic correlation ??? func:

$$f_r^{i,F} \left( x^i \middle| h_{t-1}^i, h_{t-1}^j \right) = \sum_{t=\text{start\_history}}^T I \left\{ d_t^i = x^i \right\} \cdot z_t^i \cdot z_t^j$$

have access to  $P^i, P^F, P^{i,j}$

## 7 IMPLEMENTATION

In this section implementation of the RRLN system is considered. Thus, to make this possible Section 7.1 considers heuristics used to render the theoretical system implementable. Section ?? considers experimental results.

### 7.1 Heuristics

#### 7.1.1 Tractability

It is not possible to map infinite state spaces in practice, so it is advantageous to set up  $f_Q$  on a subspace. The core goal during exploration of an MDP  $M_i$  is to learn Q-values and use Q-values as a decision-making aid.

1. Recall Section ??, encoding:  $Q(s, a) \leftarrow \sum_{s' \in S} \tilde{R}(s'|a, s) \tilde{T}(s'|a, s) + \gamma Q(s', a)$
2.  $\pi_t(s, a) = \frac{e^{Q(s, a)}}{\sum_{s' \in S} 1 + e^{Q(s', a)}}$ ,  $\pi_t^*(s) = \arg \max_a \pi(s, a)$

Clearly exploring the space  $S$  each iteration is impossible, thus ???  $\mathfrak{Z}$  as a local subset. Thus using  $\mathfrak{Z}_{(S)}$  it is possible to use functions  $f_Q$  and  $f_m$  to regress to  $Q(S \times A)$  images.

$$\mathfrak{Z}_{(S)} \subseteq S \times A \times S$$

$\mathfrak{Z}_{(S)}$  is tractable size! and local!

$$f_m : Q_{t+1}(S \times A) \times Q_t(\mathfrak{Z}) \rightarrow Q_t(S \times A)$$

#### 7.1.2 Learning

We also need to keep  $\tilde{T}$  and  $\tilde{R}_t$  updated and can employ SGD regression instances to do this. As an MDP  $M_i$  is exploring the environment the process in Section 7.1.1 will always supply the latest and most up-to-date Q-values for decision-making. Thus, the values  $\tilde{T}$  and  $\tilde{R}$  must be updated online to facilitate this. To do so both  $\tilde{T}$  and  $\tilde{R}$  can be updated with stochastic gradient descent:

$$f_T : \tilde{T}_{t-1}(\mathfrak{Z}) \times \{s, a, s'\} \rightarrow \tilde{T}_t(\mathfrak{Z})$$

$$f_R : \tilde{R}_{t-1}(\mathfrak{Z}) \times \{R(s'|a, s)\} \rightarrow \tilde{R}_t(\mathfrak{Z})$$

④ Implementation: I use instances of stochastic gradient descent to regress  $\boxed{f_T}$  and  $\boxed{f_R}$

- $\tilde{T}_t(s'|a, s) \leftarrow f_T(s, a, s', T_{t-1}) = \tilde{T}_{t-1}(s'|a, s) + \alpha_T \left( \frac{f_r(s, a, s')}{f_r(s, a)} - T_{t-1}(s'|a, s) \right)$
- ...  $\tilde{R}_t \leftarrow$  user defined / observed from sensors

where  $f_r(s, a, s')$  and  $f_r(s, a)$  reflect visitation frequencies.

$f_Q$  is more difficult,

$$f_m : Q_{t+1}(s, a) \leftarrow f_Q(\mathfrak{Z}, \tilde{T}_t, \tilde{R}_t) = \sum_{s' \in \mathfrak{Z}} \tilde{T}(s'|s, a) \tilde{R}(s'|s, a) + \gamma V(s)$$

where  $V(s) \approx \arg \max_a Q_{t-1}(s, a)$ .



## PAGE 17

On The Generalization and reuse of transitional knowledge #2

---

- ① Setting the state, taking a general FOMDP given usual expectations (stably stochastic etc.)  $m = \langle S, A, T, R \rangle$  want to find  $\pi^* : \{S \times A\} \cup Q(S, A) \rightarrow A$  s. t. for some value function  $V(s)$ ,  $\pi^*(s) = \arg \max_a \sum_{s'} T(s'|a, s) R(s'|a, s) + \gamma V(s')$

Traditionally, convergence can be found directly, using stochastic gradient descent

$$Q_{t+1}(s, a) = Q_t(s, a) + \alpha \left( R(s'|s, a) - Q_t(s, a) + \gamma \arg \max_a (Q_t(s', a^*)) \right)$$

which is limited because as  $Q(S, A)$  converges, it becomes difficult to adjust to changes in  $R(S, A, S)$ .

---

- ② Optimization objectives change, meaning the basis of  $Q(S, A)$  is typically malleable in real-life scenarios. In this paper we present a method for separating transitional models and reward models. We hold reward and transitional functioning separate as  $\tilde{T}$  and  $\tilde{R}$ ; and attempt to regress to true values s. t.  $\tilde{T} \approx T$  and  $\tilde{R} \approx R$ . We then develop a  $Q_{\text{map}}$  function  $f_Q$  to ???  $Q(S, A)$  space as needed:

$$f_Q : \tilde{T}(S \times A \times S) \times \tilde{R}(S \times A \times S) \rightarrow Q_t(S, A)$$

**PAGE 18**

MDP: Linearization of Reward/Optimal Policy

 $P_a - P_a(i, j)$  represents  $T(s_i, a, s_j)$        $S$  – all states $\gamma$  – decay factor       $(0, 1)$  $\pi$  – policy $V^\pi(s)$  – typical value function $\mathbf{V}^\pi$  – vector of all values  $\{V^\pi(s_1), \dots, V^\pi(s_n)\}$  $\prec$  and  $\preceq$  denote strict and non-strict vectoral inequality. $\mathbf{R}$  – vector of reward (like  $\mathbf{V}^\pi(s)$ )

for optimal reward:

$$(P_{a_i} - P_a)(I - \gamma P_{a_i})^{-1} \mathbf{R} \succeq 0 \quad \Leftrightarrow$$

Proof (cool as fuck):

$$a_1 \equiv \pi(s) \in \arg \max_{a \in A} \sum_{s'} P_{s_a}(s') V^\pi(s') \quad \forall s \in S$$

$$\sum_{s'} P_{s_{a_1}} \geq \sum_{s'} P_{s_a}(s') V^\pi(s') \quad \forall s \in S, a \in A$$

$$\vdots \quad \begin{array}{c} \nwarrow \\ \nearrow \end{array} a_1 \text{ is Pareto efficient (!)}$$

$$P_{a_1} \mathbf{V}^\pi \succeq P_a \mathbf{V}^\pi \quad \forall a \in A \setminus a_1 \quad (\text{non-strict improvement})$$

$$\vdots$$

$$P_{a_1}(I - \gamma P_{a_1})^{-1} \mathbf{R} \succeq P_a(I - \gamma P_{a_1})^{-1} \mathbf{R} \quad \forall a \in A \setminus a_1$$

The hard part to verify:  $\mathbf{V}^\pi = (I - \gamma P_{a_1}) \mathbf{R}$

**PAGE 19**Transitional Learning Continued

- $\{s, s', w\}$  can be controlled to both represent the state space and accurately represent  $P + (s'|s', a)$
- $A + h, B + h$  and  $\gamma$  can be controlled to speed the algorithm  $R_\gamma$
- Q-learning can still be used, if calculation of  $\bar{R}_\gamma$  is too “slow”.
- the reward function  $R(s, a, s')$  can be redefined at an instant to allow immediate re-calculation of a policy  $\bar{R}_\gamma(s, a)$ .

Possible experiments:

- show speed of convergence is greater, due to the “storing” of the transitional model across all actions
- show that the generalized learning allows for redefinition of the reward function.

**PAGE 20**

Policy

Convergence of “Bad MDP”

Question, given  $\pi_i^*$ , is it possible to find

$$f : \pi_j^*, \pi_k^* \rightarrow \pi_i^*$$

$$f(\pi_j(s))$$

a) suppose  $f(\pi_j(s, \pi_k(s))) = \pi_j(s, \pi_k(s)) \cup \pi_k(s)$

b)  $S_k, S_j$  – assume a subspace that is independent of effect by  $A$ ,  $S_j \in S_i$

$A_k, A_j$  – assume a subset of  $A_j$

$T_k, T_j$  – assume  $T(s_i, a, s'_j) = 0 \quad \forall (s_i, a_j, s_j) \in S_j \times A_j \times S_j$

$R_k, R_j$  – assume  $R(s_i, a, s_j)$

In this case,  $A_j$  must effect  $T(s_k, a_k, s'_k)$  and  $A_k$  must effect  $T(s_j, a_j, s'_j)$

effect how:

$$\exists a_j, a_k \sum_{s'_j \neq s_j} T \left( s'_j \left| \begin{array}{c} a_j \\ a_k \end{array} \right. , s_j \right) > 0$$

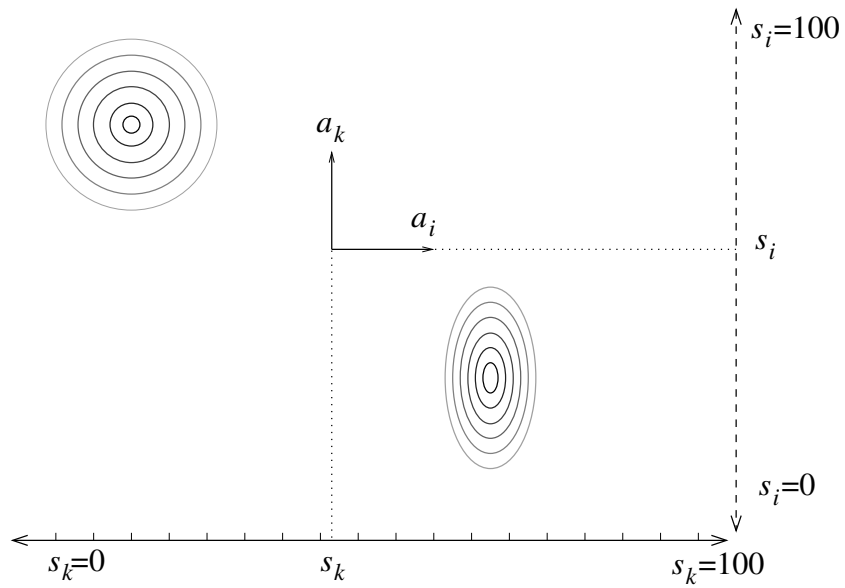
\*(Derive this conclusion from  $T(s_i, a_i, s_i)$ )\*

c)  $M_j, M_k$  execute concurrently, meaning at each time  $t \exists (a_j, a_k) \in A_i$ , chosen by  $a_j \sim \pi_j(s_j, a_k) \quad a_k \sim \pi_k(s_k)$

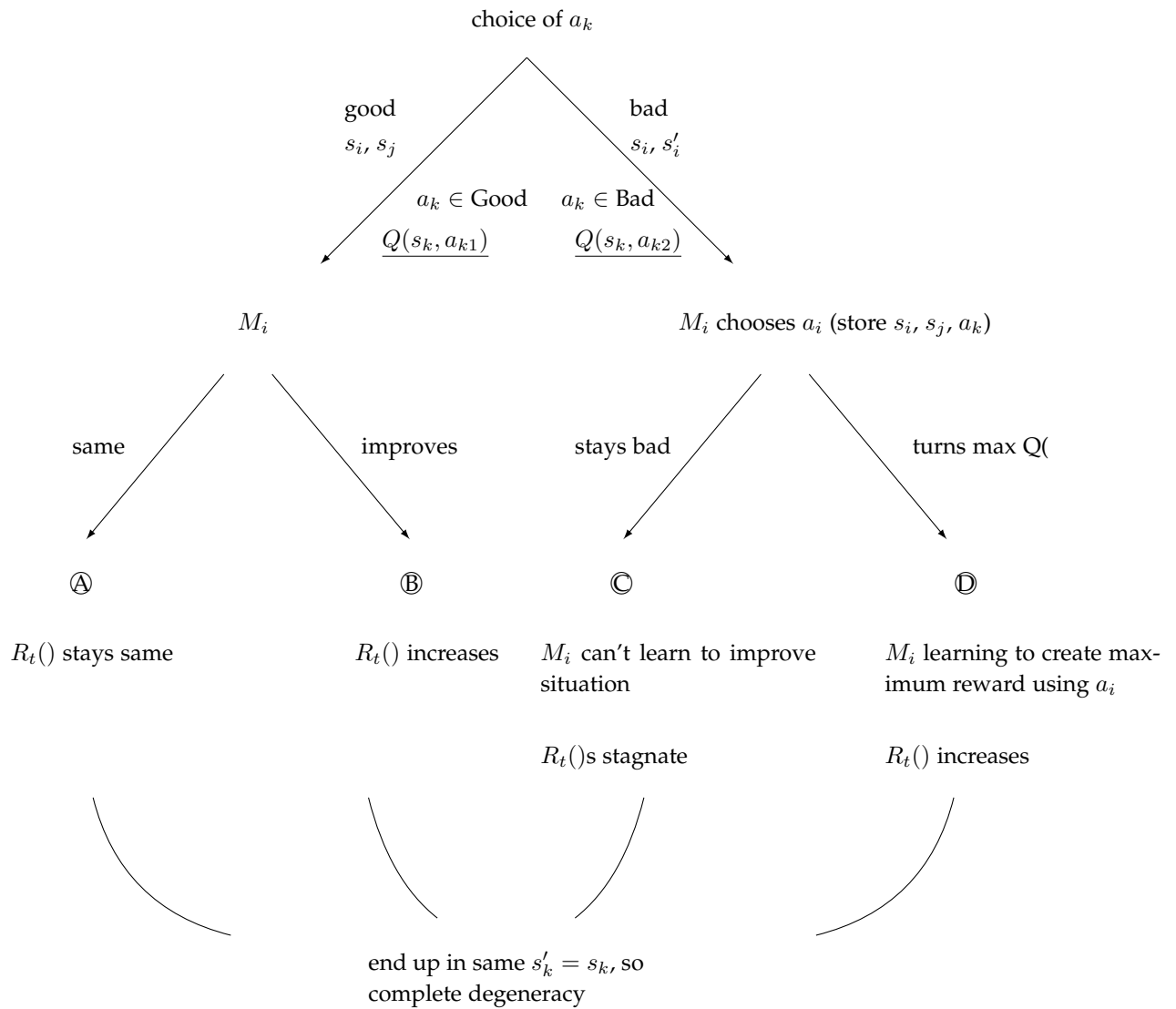
**PAGE 21**

child convergence assume during concurrent learning

- $R_t$  must be time-monotonic and convergent stochastic.
- Evaluate selection of  $s'_i$  and  $s_i$ , assuming worst case:  $a_k$  has no impact on  $s'_k, s_k$  and only an impact on  $s_i$  and  $s'_i$



consider policy learning: execution of  $\pi_k(s_k)$  yielding  $s'_k$



As  $t \rightarrow \infty$ , i) Ⓐ & Ⓒ will never be selected by  $M_k()$

ii) if Ⓑ > Ⓓ  $Q(a_k \in \text{Good}, s_k) > Q(a_k \in \text{Bad})$  and then  $a_k \in \text{Good}$  will be chosen

**PAGE 22**

Child convergence rewrite, using time index

$$m_1, \dots, m_\alpha$$

$$m_1: (s_i, s_k) \quad \text{pre-existing } s_i = 1, s_k = 1$$

$$m_2: a_k = \pi_k(s_k) = (x) \text{ — worst case, effects only } s_i, s_j,$$

→  
supposes

$$\rightarrow m_3: a_i = \pi_i(s_i, a_k)$$

$$s'_i \leftarrow 2$$

$$s'_k \leftarrow 2$$

$$m_4: Q(s_k, a_k) \leftarrow \text{update: } R(s_k, a_k, s'_k)$$

$$m_5: Q(s_i, a_k, s_i) \leftarrow \text{update: } R(s_i, a_k, a_i, s'_k)$$

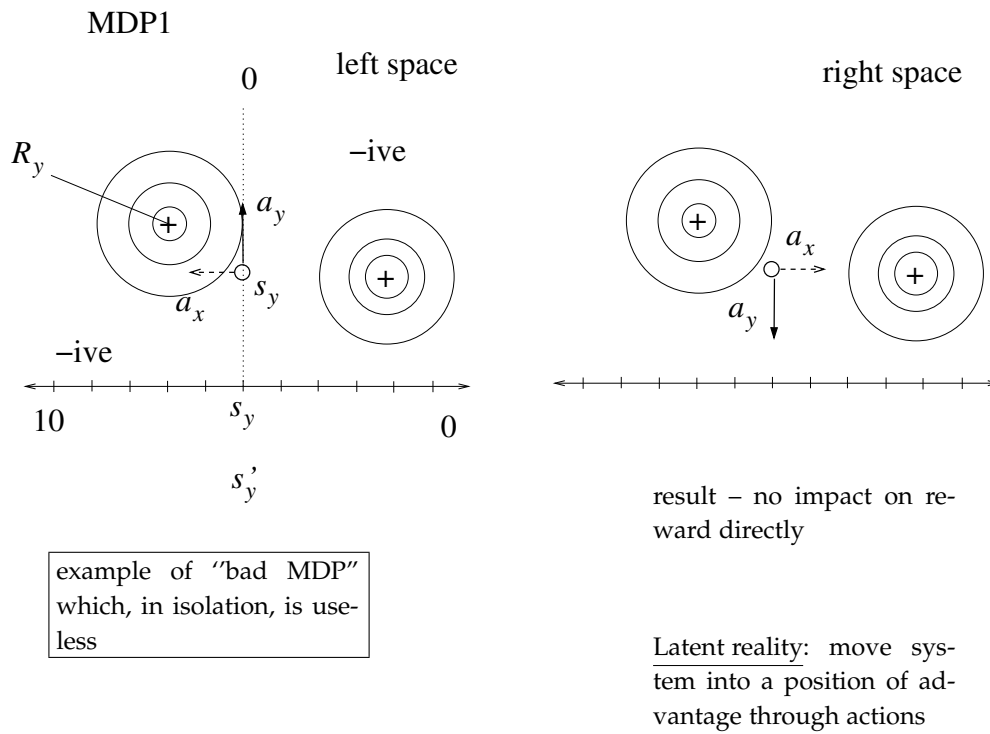
$$\hookrightarrow \text{update } \pi_i(s_i, a_k)$$

so:  $m_2$  assumed  $\pi(s_i, a_k)(m_3)$  to be convergent stochastic and monotonic in reward ???, which is assured by updating  $Q(s_i, a_k, a_i)$  at  $m_5$

$$\pi_k(s_k) \rightarrow Q(s_k, a_k) \quad \pi_i(s_i, a_k) \rightarrow Q(s_k, a_k, s_k)$$

## 8 PAGE 23

### Bad MDP Decomposition (worst possible case)



ways to understand:

- 1) analyze actions of subsystem
  - 2) analyze effectiveness of subsystem
- ↔ Do until ave coordination

only ability: coordinate acting with subsystem, s.t., an understanding of the relationship of your action & subsystem action arises