

Contributed Article in Neural Networks, Vol. 19, No. 1, pp. 62–75, 2005

# The Asymptotic Equipartition Property in Reinforcement Learning and its Relation to Return Maximization

Kazunori Iwata<sup>1</sup>

Kazushi Ikeda<sup>2</sup>

Hideaki Sakai<sup>2</sup>

<sup>1</sup> Faculty of Information Sciences, Hiroshima City University  
3-4-1 Ozuka-Higashi, Asa-Minami-Ku, Hiroshima, 731-3194, JAPAN  
Phone: +81-82-830-1679 Fax: +81-82-830-1679  
Email: [kiwata@im.hiroshima-cu.ac.jp](mailto:kiwata@im.hiroshima-cu.ac.jp)

<sup>2</sup> Department of Systems Science, Graduate School of Informatics,  
Kyoto University  
Yoshida-Honmachi, Sakyo-ku, Kyoto, 606-8501, JAPAN  
Phone: +81-75-753-5501 Fax: +81-75-753-4755  
Email: [kazushi@i.kyoto-u.ac.jp](mailto:kazushi@i.kyoto-u.ac.jp), [hsakai@i.kyoto-u.ac.jp](mailto:hsakai@i.kyoto-u.ac.jp)

### **Abstract**

We discuss an important property called the asymptotic equipartition property on empirical sequences in reinforcement learning. This states that the typical set of empirical sequences has probability nearly one, that all elements in the typical set are nearly equi-probable, and that the number of elements in the typical set is an exponential function of the sum of conditional entropies if the number of time steps is sufficiently large. The sum is referred to as stochastic complexity. Using the property we elucidate the fact that the return maximization depends on two factors, the stochastic complexity and a quantity depending on the parameters of environment. Here, the return maximization means that the best sequences in terms of expected return have probability one. We also examine the sensitivity of stochastic complexity, which is a qualitative guide in tuning the parameters of action-selection strategy, and show a sufficient condition for return maximization in probability.

**Keywords:** reinforcement learning, Markov decision process, information theory, asymptotic equipartition property, stochastic complexity, return maximization

# 1 Introduction

In information theory the weak law of large numbers is known as the asymptotic equipartition property (AEP) which was first stated in (Shannon, 1948) and then developed by the type method in (Csiszár & Körner, 1997; Csiszár, 1998). When a sequence of random variables is drawn many times, independently and according to an identical probability distribution, the AEP states that there exists the typical set of the sequences with probability nearly one, that all elements in the typical set are nearly equi-probable, and that the number of elements in the typical set is given by an exponential function of the entropy of the probability distribution. In addition, the number of elements in the typical set is quite small compared to the number of possible sequences. If the AEP also holds on empirical sequences generated from a Markov decision process (MDP) in reinforcement learning (RL), it facilitates the analysis of the learning process since most of our attention can be focused on the typical set of the empirical sequences. This leads us to the question of whether or not the AEP holds for an empirical sequence. The fact is that a similar AEP holds but it is more complicated than the original AEP. Using the type method, first we introduce an information-theoretic formulation for almost stationary ergodic MDPs in general and then describe the AEP that holds on the empirical sequences. From the AEP, we indicate the existence of an important factor called the stochastic complexity which consists of the sum of conditional entropies and elucidate that the return maximization is characterized by two factors, the stochastic complexity and a quantity which depends on the parameters of environment. Here, the return maximization means that the probability of best sequences that yield the maximal expected return goes to probability one. Also, useful knowledge for tuning the parameters of action-selection strategy is described by examining the sensitivity of the stochastic complexity. Furthermore, we show that the stochastic complexity is derived from the algorithmic complexity which was explored by Chaitin (Chaitin, 1977, 1987).

The organization of this paper is as follows. We introduce some notation and the type of empirical sequence in Section 2. Section 3 shows the main theorems associated with the AEP. Using the AEP we analyze the RL process in Section 4. Finally, we give some conclusions in Section 5. Appendices A and B are the related theorems to the AEP and those proofs, respectively.

## 2 Preliminaries

We concentrate on the discrete-time MDP with discrete states and actions in this paper. Let  $\mathcal{S} \stackrel{\text{def}}{=} \{s_1, s_2, \dots, s_I\}$  be the finite set of states of the environment,  $\mathcal{A} \stackrel{\text{def}}{=} \{a_1, a_2, \dots, a_J\}$  be the finite set of actions, and  $\mathbb{R}_0 \stackrel{\text{def}}{=} \{r_1, r_2, \dots, r_K\} \subset \mathbb{R}$  be the finite set of rewards which are discrete real numbers. Notice that  $|\mathcal{S}| = I$ ,  $|\mathcal{A}| = J$ , and  $|\mathbb{R}_0| = K$ . We assume that elements in these sets are recognized without error by the learner, hereinafter called the agent. We denote a time step by  $t$ . The stochastic variables of state, action, and reward at time step  $t$  ( $t = 1, 2, \dots$ ) are written as  $s(t)$ ,  $a(t)$ , and  $r(t)$ , respectively. The agent improves the policy by observing one-by-one each element of the empirical sequence that is generated by the interactions between the agent and the environment, as shown in Figure 1.

[Figure 1 about here.]

Now let us consider the empirical sequence of  $n$  time steps,

$$s(1), a(1), s(2), r(2), a(2), \dots, s(n), r(n), a(n), r(n+1).$$

Let  $r(n+1) = r(1)$  for notational convenience and let  $\mathbf{x} = \{s(t), a(t), r(t)\}_{t=1}^n$  denote the empirical sequence of  $n$  time steps. The state sequence, action sequence, and reward sequence of the empirical sequence  $\mathbf{x} \in (\mathcal{S} \times \mathcal{A} \times \mathbb{R}_0)^n$  are denoted by  $\mathbf{s} = \{s(t)\}_{t=1}^n$ ,  $\mathbf{a} = \{a(t)\}_{t=1}^n$ , and  $\mathbf{r} = \{r(t)\}_{t=1}^n$ , respectively. We use the term return to express the sum of rewards.

Let  $q_i \stackrel{\text{def}}{=} \Pr(s(1) = s_i)$  be the initial probability distribution and  $\mathbf{q} \stackrel{\text{def}}{=} \{q_1, q_2, \dots, q_I\}$  where  $q_i > 0$  for all  $i$ . The empirical sequence is drawn according to an ergodic MDP specified by the following two conditional probability distribution matrices. Henceforth, the conditional probability distribution matrix is simply called the matrix. The policy matrix which the agent determines is an  $I \times J$  matrix defined by

$$\mathbf{\Gamma}^\pi \stackrel{\text{def}}{=} \begin{pmatrix} p_{11} & p_{12} & \dots & p_{1J} \\ p_{21} & p_{22} & \dots & p_{2J} \\ \vdots & \vdots & \ddots & \vdots \\ p_{I1} & p_{I2} & \dots & p_{IJ} \end{pmatrix} = \begin{pmatrix} \mathbf{P}^{(1)} \\ \mathbf{P}^{(2)} \\ \vdots \\ \mathbf{P}^{(I)} \end{pmatrix}, \quad (1)$$

where  $p_{ij} \stackrel{\text{def}}{=} \Pr(a(t) = a_j | s(t) = s_i)$ . According to this matrix, the agent selects an action in a state at each time step. Note that  $\mathbf{\Gamma}^\pi$  is actually time-varying because the agent improves the policy in the process of RL. However,  $\mathbf{\Gamma}^\pi$  tends to be constant as the policy goes to be optimal by the learning. The state transition matrix of the environment is an  $IJ \times IK$  matrix defined by

$$\mathbf{\Gamma}^\mathbf{T} \stackrel{\text{def}}{=} \begin{pmatrix} p_{1111} & p_{1112} & \dots & p_{11IK} \\ p_{1211} & p_{1212} & \dots & p_{12IK} \\ \vdots & \vdots & \ddots & \vdots \\ p_{IJ11} & p_{IJ12} & \dots & p_{IJIK} \end{pmatrix} = \begin{pmatrix} \mathbf{P}^{(11)} \\ \mathbf{P}^{(12)} \\ \vdots \\ \mathbf{P}^{(IJ)} \end{pmatrix}, \quad (2)$$

where  $p_{ijj'k} \stackrel{\text{def}}{=} \Pr(s(t+1) = s_{i'}, r(t+1) = r_k | s(t) = s_i, a(t) = a_j)$ . The agent does not know the matrix  $\mathbf{\Gamma}^\mathbf{T}$  of the environment but can estimate it by observing the results for an action. We assume that  $\mathbf{\Gamma}^\mathbf{T}$  is constant and that for simplicity of analysis  $\mathbf{\Gamma}^\pi$  is temporarily fixed for  $n$  time steps where  $n$  is sufficiently large. For notational simplicity we define  $\mathbf{\Gamma} \stackrel{\text{def}}{=} (\mathbf{\Gamma}^\pi, \mathbf{\Gamma}^\mathbf{T})$ . Since MDPs are characterized by the finite sets, the initial probability distribution, and the matrices, we denote the MDP by  $\text{M}(\mathcal{S}, \mathcal{A}, \mathbb{R}_0, \mathbf{q}, \mathbf{\Gamma})$ .

## 2.1 Type of Empirical Sequence

Let  $n_i$  ( $n_i \leq n$ ) denote the number of times that a state  $s_i \in \mathcal{S}$  occurs in the empirical sequence of  $n$  time steps,  $\mathbf{x} = (\mathbf{s}, \mathbf{a}, \mathbf{r}) \in (\mathcal{S} \times \mathcal{A} \times \mathbb{R}_0)^n$ . In a similar manner, let  $n_{ij}$  ( $n_{ij} \leq n_i$ ) be the number of occurrences of  $t$  such that  $(s(t), a(t)) = (s_i, a_j) \in \mathcal{S} \times \mathcal{A}$  in the empirical sequence. With an additional ‘‘cyclic’’ convention that  $s(n)$ ,  $a(n)$ , and  $r(n+1) = r(1)$  precede  $s(1)$ ,  $a(1)$ , and  $r(2)$ , let  $n_{ijj'k}$  ( $n_{ijj'k} \leq n_{ij}$ ) denote the number of occurrences of  $t$  such that  $(s(t), a(t), s(t+1), r(t+1)) = (s_i, a_j, s_{i'}, r_k) \in \mathcal{S} \times \mathcal{A} \times \mathcal{S} \times \mathbb{R}_0$  in the empirical sequence. Note that the cyclic convention is for simplicity of development. The discussions in this paper strictly hold even if we do not assume this convention. The relationship among the non-negative numbers  $n$ ,  $n_i$ ,  $n_{ij}$ , and  $n_{ijj'k}$  is expressed as

$$n = \sum_{i=1}^I n_i = \sum_{i=1}^I \sum_{j=1}^J n_{ij} = \sum_{i=1}^I \sum_{j=1}^J \sum_{i'=1}^I \sum_{k=1}^K n_{ijj'k}. \quad (3)$$

Now we define the type of  $s_i \in \mathcal{S}$  by

$$f_i = \frac{n_i}{n}. \quad (4)$$

The type is generally called the empirical distribution (Han & Kobayashi, 2002, p. 42) because we can regard each sequence as a sample from a stochastic process. Also, the joint type of  $(s_i, a_j) \in \mathcal{S} \times \mathcal{A}$  is defined as

$$f_{ij} = \frac{n_{ij}}{n}. \quad (5)$$

Let us denote all the types and the joint types by

$$\mathbf{F}_S \stackrel{\text{def}}{=} (f_1, f_2, \dots, f_I), \quad (6)$$

and

$$\mathbf{F}_{S\mathcal{A}} \stackrel{\text{def}}{=} \begin{pmatrix} f_{11} & f_{12} & \dots & f_{1J} \\ f_{21} & f_{22} & \dots & f_{2J} \\ \vdots & \vdots & \ddots & \vdots \\ f_{I1} & f_{I2} & \dots & f_{IJ} \end{pmatrix}, \quad (7)$$

respectively. In this case we say that the state sequence  $\mathbf{s}$  and the state-action sequence  $(\mathbf{s}, \mathbf{a})$  have the type  $\mathbf{F}_S$  and the joint type  $\mathbf{F}_{S\mathcal{A}}$ , respectively.

**Conditional Type Relative to Policy** If  $n_i > 0$  for all  $i$ , then the conditional type  $g_{ij}$  of  $(s_i, a_j) \in \mathcal{S} \times \mathcal{A}$  given a state sequence  $\mathbf{s} \in \mathcal{S}^n$  is defined as

$$n_{ij} \stackrel{\text{def}}{=} g_{ij} n_i. \quad (8)$$

However, if there exists  $i$  such that  $n_i = 0$ , then we can not uniquely determine the conditional type (see Example 2.1). To avoid such a case, we consider the set of action sequences given any state sequence  $\mathbf{s}$  having the type  $\mathbf{F}_S$  and an  $I \times J$  matrix  $\Phi^\pi : \mathcal{S} \rightarrow \mathcal{A}$  expressed as

$$\Phi^\pi \stackrel{\text{def}}{=} \begin{pmatrix} g_{11} & g_{12} & \dots & g_{1J} \\ g_{21} & g_{22} & \dots & g_{2J} \\ \vdots & \vdots & \ddots & \vdots \\ g_{I1} & g_{I2} & \dots & g_{IJ} \end{pmatrix} = \begin{pmatrix} \mathbf{G}_{(1)} \\ \mathbf{G}_{(2)} \\ \vdots \\ \mathbf{G}_{(I)} \end{pmatrix}. \quad (9)$$

In short,  $n_{ij}$  is decided by  $n_i$  and  $g_{ij}$  for every  $i, j$ . The set of action sequences, which is uniquely determined in this way, is referred to as  $\Phi^\pi$ -shell (Csiszár & Körner, 1997, p. 31) of  $\mathbf{s}$  and denoted by  $\mathcal{C}^n(\Phi^\pi, \mathbf{s})$ . The entire set of possible matrices  $\Phi^\pi$  for any state sequence with the type  $\mathbf{F}_S$  is simply written as  $\Lambda_n^\pi$ .

**Example 2.1** Let  $I = J = 2$ , the state sequence  $\mathbf{s} = (s_1, s_1, s_1, s_1) \in \mathcal{S}^4$ , and the action sequence  $\mathbf{a} = (a_1, a_1, a_2, a_2) \in \mathcal{A}^4$ . Then, from the definition of (8) we obtain  $g_{11} = g_{12} = 1/2$ . Also, because of  $n_2 = 0$ , letting  $g_{21} = \omega$  where  $0 \leq \omega \leq 1$  we have  $g_{21} = \omega$  and  $g_{22} = 1 - \omega$ . Therefore, we can not uniquely determine the conditional type.

**Example 2.2 ( $\Phi^\pi$ -shell)** Let  $I = J = 2$ , again. For the state sequence  $\mathbf{s} = (s_1, s_1, s_1, s_2) \in \mathcal{S}^4$  with the type  $\mathbf{F}_S = (3/4, 1/4)$  and the matrix,

$$\Phi^\pi = \begin{pmatrix} g_{11} & g_{12} \\ g_{21} & g_{22} \end{pmatrix} = \begin{pmatrix} 2/3 & 1/3 \\ 0 & 1 \end{pmatrix}, \quad (10)$$

the  $\Phi^\pi$ -shell of  $\mathbf{s}$  is  $\mathcal{C}^4(\Phi^\pi, \mathbf{s}) = \{(a_1, a_1, a_2, a_2), (a_1, a_2, a_1, a_2), (a_2, a_1, a_1, a_2)\}$ .

**Conditional Markov Type Relative to State Transition** In a slightly different manner we need to deal with the conditional Markov type<sup>1</sup>. We consider the set of state-reward sequences such that the

<sup>1</sup>For Markov type, see (Davisson, Longo, & Sgarro, 1981).

joint type is  $\mathbf{F}_{S,A}$  given any action sequence and an  $IJ \times IK$  matrix  $\Phi^T : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{S} \times \mathbb{R}_0$  designated by

$$\Phi^T \stackrel{\text{def}}{=} \begin{pmatrix} g_{1111} & g_{1112} & \cdots & g_{11IK} \\ g_{1211} & g_{1212} & \cdots & g_{12IK} \\ \vdots & \vdots & \ddots & \vdots \\ g_{IJ11} & g_{IJ12} & \cdots & g_{IJIK} \end{pmatrix} = \begin{pmatrix} \mathbf{G}_{(11)} \\ \mathbf{G}_{(12)} \\ \vdots \\ \mathbf{G}_{(IJ)} \end{pmatrix}. \quad (11)$$

The set of state-reward sequences is referred to as  $\Phi^T$ -shell and denoted by  $\mathcal{C}^n(\Phi^T, \mathbf{F}_{S,A})$ . The entire set of possible matrices  $\Phi^T$  such that the joint type is  $\mathbf{F}_{S,A}$  for any action sequence is simply written as  $\Lambda_n^T$ .

For simplicity, we define  $\Phi \stackrel{\text{def}}{=} (\Phi^\pi, \Phi^T)$  and  $\Lambda_n \stackrel{\text{def}}{=} \Lambda_n^\pi \times \Lambda_n^T$ . The set of empirical sequences that consists of the  $\Phi^\pi$ -shell and  $\Phi^T$ -shell is called the  $\Phi$ -shell and denoted by  $\mathcal{C}^n(\Phi, \mathbf{F}_S, \mathbf{F}_{S,A})$ . The structure of the  $\Phi$ -shell is depicted in Figure 2. When a joint-type  $\mathbf{F}_{S,A}$  and a matrix  $\Phi^T$  are given, the  $\Phi^T$ -shell having the type  $\mathbf{F}_S$  is uniquely determined and then the combination of each element in the  $\Phi^T$ -shell and a matrix  $\Phi^\pi$  produces the  $\Phi^\pi$ -shell. Therefore, the  $\Phi$ -shell is uniquely determined. Notice that

$$|\mathcal{C}^n(\Phi, \mathbf{F}_S, \mathbf{F}_{S,A})| = \sum_{(s', r') \in \mathcal{C}^n(\Phi^T, \mathbf{F}_{S,A})} |\mathcal{C}^n(\Phi^\pi, s')|. \quad (12)$$

In this case we write that the empirical sequence has the conditional type matrix  $\Phi$ .

[Figure 2 about here.]

## 2.2 V-typical and W-typical sequences

In order to prove the AEP on empirical sequences, we have to introduce the **V**-typical sequence with respect to the state sequences and the **W**-typical sequences with respect to the state-action sequences.

**Definition 2.1 (V-typical and W-typical sequences)** *We assume the existence of the following two unique stationary probability distributions,*

$$\mathbf{V} \stackrel{\text{def}}{=} (v_1, v_2, \dots, v_I), \quad (13)$$

$$\mathbf{W} \stackrel{\text{def}}{=} \begin{pmatrix} w_{11} & w_{12} & \cdots & w_{1J} \\ w_{21} & w_{22} & \cdots & w_{2J} \\ \vdots & \vdots & \ddots & \vdots \\ w_{I1} & w_{I2} & \cdots & w_{IJ} \end{pmatrix}, \quad (14)$$

and assume that  $\mathbf{F}_S$  and  $\mathbf{F}_{S,A}$  tend to  $\mathbf{V}$  and  $\mathbf{W}$  as  $n \rightarrow \infty$ , respectively. The stationary probability distributions are uniquely determined by the MDP,  $\mathbb{M}(\mathcal{S}, \mathcal{A}, \mathbb{R}_0, \mathbf{q}, \Gamma)$ . In this case, there exists a sequence of positive  $\kappa_n$  such that  $\kappa_n \rightarrow 0$  as  $n \rightarrow \infty$ , and if the type  $\mathbf{F}_S$  of a state sequence  $\mathbf{s} \in \mathcal{S}^n$  satisfies

$$D(\mathbf{F}_S \| \mathbf{V}) = \sum_{i=1}^I f_i \log \frac{f_i}{v_i} \leq \kappa_n, \quad (15)$$

then we call the state sequence a **V**-typical sequence. The set of **V**-typical sequences is denoted by  $\mathcal{C}_{\kappa_n}^n(\mathbf{V}) \stackrel{\text{def}}{=} \{\mathbf{s} \in \mathcal{S}^n | D(\mathbf{F}_S \| \mathbf{V}) \leq \kappa_n\}$ . In a similar manner, there exists a sequence of positive  $\xi_n$  such that  $\xi_n \rightarrow 0$  as  $n \rightarrow \infty$ , and if

$$D(\mathbf{F}_{S,A} \| \mathbf{W}) = \sum_{i=1}^I \sum_{j=1}^J f_{ij} \log \frac{f_{ij}}{w_{ij}} \leq \xi_n \quad (16)$$

holds, then the state-action sequences  $(\mathbf{s}, \mathbf{a}) \in (\mathcal{S} \times \mathcal{A})^n$  are referred to as  $\mathbf{W}$ -typical sequences. We define the set of  $\mathbf{W}$ -typical sequences as  $\mathcal{C}_{\xi_n}^n(\mathbf{W}) \stackrel{\text{def}}{=} \{(\mathbf{s}, \mathbf{a}) \in (\mathcal{S} \times \mathcal{A})^n | D(\mathbf{F}_{\mathcal{S}\mathcal{A}} \| \mathbf{W}) \leq \xi_n\}$ .

In the rest of this section, we will introduce a few basic conventions in information theory (Cover & Thomas, 1991, Chapter 2). Let us use the convention that  $0 \log 0 = 0$  henceforth. The function  $H$  indicates the entropy. For instance, we write the entropy of  $\mathbf{P}_{(i)}$  in (1) for any  $i$  as

$$H(\mathbf{P}_{(i)}) = - \sum_{j=1}^J p_{ij} \log p_{ij}, \quad (17)$$

and describe its conditional entropy given  $\mathbf{V}$  as

$$H(\mathbf{\Gamma}^\pi | \mathbf{V}) = \sum_{i=1}^I v_i H(\mathbf{P}_{(i)}). \quad (18)$$

Also, as used in (15) and (16), the divergence is designated by the function  $D$ . The divergence between  $\mathbf{\Phi}^\pi$  and  $\mathbf{\Gamma}^\pi$  given  $\mathbf{F}_S$  is denoted as

$$D(\mathbf{\Phi}^\pi \| \mathbf{\Gamma}^\pi | \mathbf{F}_S) = \sum_{i=1}^I f_i D(\mathbf{G}_{(i)} \| \mathbf{P}_{(i)}), \quad (19)$$

where

$$D(\mathbf{G}_{(i)} \| \mathbf{P}_{(i)}) = \sum_{j=1}^J g_{ij} \log \frac{g_{ij}}{p_{ij}}. \quad (20)$$

### 3 Asymptotic Equipartition Property

In this section, it is elucidated that the empirical sequences generated from almost stationary ergodic MDPs have the AEP. Now we are in a position to give the definitions of the typical sequence and the typical set of empirical sequences, which will lead us to show that the AEP holds on empirical sequences.

**Definition 3.1 ( $\mathbf{\Gamma}$ -typical sequence and  $\mathbf{\Gamma}$ -typical set)** *If the matrix  $\mathbf{\Phi} \in \mathbf{\Lambda}_n$  of the conditional types with respect to an empirical sequence  $\mathbf{x} = (\mathbf{s}, \mathbf{a}, \mathbf{r}) \in (\mathcal{S} \times \mathcal{A} \times \mathbb{R}_0)^n$  satisfies*

$$D(\mathbf{\Phi}^\pi \| \mathbf{\Gamma}^\pi | \mathbf{F}_S) + D(\mathbf{\Phi}^\top \| \mathbf{\Gamma}^\top | \mathbf{F}_{\mathcal{S}\mathcal{A}}) \leq \lambda_n, \quad (21)$$

*for any matrix  $\mathbf{\Gamma}$  and positive number  $\lambda_n$ , then the empirical sequence is called a  $\mathbf{\Gamma}$ -typical sequence. The set of such empirical sequences is also called the  $\mathbf{\Gamma}$ -typical set and denoted by  $\mathcal{C}_{\lambda_n}^n(\mathbf{\Gamma})$ . That is,  $\mathcal{C}_{\lambda_n}^n(\mathbf{\Gamma})$  is given by*

$$\mathcal{C}_{\lambda_n}^n(\mathbf{\Gamma}) \stackrel{\text{def}}{=} \bigcup_{\substack{\mathbf{\Phi} \in \mathbf{\Lambda}_n: \\ D(\mathbf{\Phi}^\pi \| \mathbf{\Gamma}^\pi | \mathbf{F}_S) + D(\mathbf{\Phi}^\top \| \mathbf{\Gamma}^\top | \mathbf{F}_{\mathcal{S}\mathcal{A}}) \leq \lambda_n}} \mathcal{C}^n(\mathbf{\Phi}, \mathbf{F}_S, \mathbf{F}_{\mathcal{S}\mathcal{A}}). \quad (22)$$

Figure 3 illustrates the concept of Definition 3.1. The matrix  $\mathbf{\Phi}$  of the  $\mathbf{\Gamma}$ -typical sequence exists in the neighborhood of  $\mathbf{\Gamma}$ , shown by the shaded circle on the manifold spanned by  $\mathbf{\Gamma}$ .

[Figure 3 about here.]

From the theorems, presented in Appendix A, we can derive the following three theorems regarding the AEP on empirical sequences. We begin with the theorem similar to (Wolfowitz, 1978).

**Theorem 3.1 (Probability of the  $\Gamma$ -typical set)** *If  $\lambda_n \rightarrow 0$  as  $n \rightarrow \infty$  and  $\lambda_n$  satisfies*

$$\lambda_n - \frac{(IJ + I^2JK) \log(n+1) + \log I - \log \nu}{n} > 0, \quad (23)$$

where

$$\nu \stackrel{\text{def}}{=} \min_{1 \leq i, i' \leq I, 1 \leq j \leq J, 1 \leq k \leq K: p_{ij i' k} > 0} p_{ij i' k}, \quad (24)$$

there exists a sequence  $\{\varepsilon_n(I, J, K, \lambda_n)\}$  such that  $\varepsilon_n(I, J, K, \lambda_n) \rightarrow 0$  and then

$$\Pr(\mathcal{C}_{\lambda_n}^n(\Gamma)) = 1 - \varepsilon_n(I, J, K, \lambda_n). \quad (25)$$

Note that  $n\lambda_n \rightarrow \infty$  because of (23). The proof is given in Appendix B.4. This theorem implies that the probability of the  $\Gamma$ -typical set asymptotically goes to one independently of the underlying probabilistic structures,  $\Gamma^\pi$  and  $\Gamma^T$ . Next, the following theorem indicates the fact that all elements in the  $\Gamma$ -typical set are nearly equi-probable.

**Theorem 3.2 (Equi-probability of the  $\Gamma$ -typical sequence)** *If  $\mathbf{s} \in \mathcal{C}_{\kappa_n}^n(\mathbf{V})$ ,  $(\mathbf{s}, \mathbf{a}) \in \mathcal{C}_{\xi_n}^n(\mathbf{W})$ ,  $\mathbf{x} \in \mathcal{C}_{\lambda_n}^n(\Gamma)$  such that  $\kappa_n \rightarrow 0$ ,  $\xi_n \rightarrow 0$ ,  $\lambda_n \rightarrow 0$  as  $n \rightarrow \infty$ , then there exists a sequence  $\{\rho_n(I, J, K, \kappa_n, \xi_n, \lambda_n)\}$  such that*

$$\rho_n(I, J, K, \kappa_n, \xi_n, \lambda_n) \rightarrow 0.$$

Then,

$$\frac{\log \nu}{n} - \rho_n \leq -\frac{1}{n} \log \Pr(\mathbf{x}) - \{\mathbb{H}(\Gamma^\pi|\mathbf{V}) + \mathbb{H}(\Gamma^T|\mathbf{W})\} \leq -\frac{\log \mu}{n} + \lambda_n + \rho_n, \quad (26)$$

where  $\nu$  is given in (24) and

$$\mu \stackrel{\text{def}}{=} \min_{1 \leq i \leq I: q_i > 0} q_i. \quad (27)$$

This theorem is proved in Appendix B.5. Finally, we present the theorem which implies that the number of elements in the  $\Gamma$ -typical set is written as an exponential function of the sum of the conditional entropies.

**Theorem 3.3 (Bound of the number of the  $\Gamma$ -typical sequences)** *If  $\mathbf{s} \in \mathcal{C}_{\kappa_n}^n(\mathbf{V})$ ,  $(\mathbf{s}, \mathbf{a}) \in \mathcal{C}_{\xi_n}^n(\mathbf{W})$ ,  $\mathbf{x} \in \mathcal{C}_{\lambda_n}^n(\Gamma)$  such that  $\kappa_n \rightarrow 0$ ,  $\xi_n \rightarrow 0$ ,  $\lambda_n \rightarrow 0$  as  $n \rightarrow \infty$ , then there exist two sequences,  $\{\zeta_n(I, J, K, \kappa_n, \xi_n, \lambda_n)\}$  and  $\{\eta_n(I, J, K, \kappa_n, \xi_n, \lambda_n)\}$ , such that*

$$\zeta_n(I, J, K, \kappa_n, \xi_n, \lambda_n) \rightarrow 0, \quad \eta_n(I, J, K, \kappa_n, \xi_n, \lambda_n) \rightarrow 0,$$

respectively. Then, the number of elements in the  $\Gamma$ -typical set is bounded by

$$\exp[n\{\mathbb{H}(\Gamma^\pi|\mathbf{V}) + \mathbb{H}(\Gamma^T|\mathbf{W}) - \zeta_n\}] \leq |\mathcal{C}_{\lambda_n}^n(\Gamma)| \leq \exp[n\{\mathbb{H}(\Gamma^\pi|\mathbf{V}) + \mathbb{H}(\Gamma^T|\mathbf{W}) + \eta_n\}]. \quad (28)$$

The proof is given in Appendix B.6. The ratio of the number of  $\Gamma$ -typical sequences to that of all empirical sequences  $\mathbf{x} \in (\mathcal{S} \times \mathcal{A} \times \mathbb{R}_0)^n$  of  $n$  time steps is

$$\frac{|\mathcal{C}_{\lambda_n}^n(\Gamma)|}{(IJK)^n} \leq \exp[n\{\mathbb{H}(\Gamma^\pi|\mathbf{V}) + \mathbb{H}(\Gamma^T|\mathbf{W}) + \eta_n - \log I - \log J - \log K\}] \rightarrow 0, \quad (29)$$

as  $n \rightarrow \infty$ , when the probability distributions of  $\Gamma^\pi$  and  $\Gamma^T$  are not uniform distributions, that is,

$$\mathbb{H}(\Gamma^\pi|\mathbf{V}) < \log I, \quad (30)$$

$$\mathbb{H}(\Gamma^T|\mathbf{W}) < \log J + \log K. \quad (31)$$

Hence, we can say that the  $\Gamma$ -typical set is quite small in comparison to the set of all empirical sequences. Nonetheless, their existence is important enough because the total probability is almost one.



**Remark 3.1** *The equation (28) shows*

$$|\mathcal{C}_{\lambda_n}^n(\Gamma)| \doteq \exp [n \{H(\Gamma^\pi|\mathbf{V}) + H(\Gamma^\pi|\mathbf{W})\}], \quad (32)$$

where the notation  $\doteq$  indicates that both sides are equal to the first order in the exponent, namely,

$$\lim_{n \rightarrow \infty} \frac{1}{n} \log |\mathcal{C}_{\lambda_n}^n(\Gamma)| = \lim_{n \rightarrow \infty} \frac{1}{n} \log \exp [n \{H(\Gamma^\pi|\mathbf{V}) + H(\Gamma^\pi|\mathbf{W})\}], \quad (33)$$

(Cover & Thomas, 1991, p. 55).

## 4 The Role of Stochastic Complexity in Reinforcement Learning

The agent learns the optimal policy via return maximization (RM) in RL. A number of studies have been made on the analysis of the process of RM (Jaakkola, Jordan, & Singh, 1994; Kushner & Yin, 1997; Singh, Jaakkola, Littman, & Szepesvári, 2000), but most of the studies focus on concrete stochastic approximation methods such as temporal difference (TD) learning. The aim here is to explore a more general mechanism of RM, how the probability of the subset of best sequences in terms of expected return is maximized, from a viewpoint of Shannon's ideas. In this section, we state the existence of an important factor called the stochastic complexity and show new insights about the role of the stochastic complexity in RM. We also discuss a sensitivity helpful in tuning the parameters of action-selection (AS) strategy and exhibit a sufficient condition for RM. We first give a review of the TD learning and typical AS strategies.

### 4.1 Temporal Difference Learning and Action Selection Strategy

Let  $Q_{ij}$  denote the estimate of an action-value function (Sutton & Barto, 1998, Chapter 3) with respect to a state-action pair  $(s_i, a_j) \in \mathcal{S} \times \mathcal{A}$ . Let  $\mathcal{A}_i$  be the set of indices of actions available in a state  $s_i \in \mathcal{S}$ . The TD learning is an iterative approximation method to directly update the estimate of the action-value function from an observed event, without explicitly treating the matrix  $\Gamma^\pi$  of the environment. We introduce the one-step version of Q-learning (Watkins & Dayan, 1992), here. At each time step  $t$ , for an observed one-step event  $(s(t), a(t), s(t+1), r(t+1)) = (s_i, a_j, s_{i'}, r_k)$ , the estimate  $Q_{ij}$  is updated by

$$Q_{ij} \leftarrow Q_{ij} + \alpha_t \delta Q_{ij i'}, \quad (34)$$

where  $\alpha_t$  is a learning rate at time step  $t$  and

$$\delta Q_{ij i'} = r_k + \gamma \max_{j' \in \mathcal{A}_{i'}} Q_{i' j'} - Q_{ij}, \quad (35)$$

where  $\gamma$  denotes the discount factor that controls the relative importance of an immediate reward and delayed rewards. The learning rate  $\alpha_t$ , where  $0 \leq \alpha_t \leq 1$ , is gradually decreased with respect to  $t$  such that the trajectory of the mean ordinary differential equation of  $Q_{ij}$  has a limit point. Under certain conditions (Dayan, 1992), it was proved that all  $Q_{ij}$  converge to the expected values with probability one. The convergence theorem was extended to more general versions using the stochastic approximation method in (Jaakkola et al., 1994; Tsitsiklis, 1994).

Next, we review the following two AS strategies that have been employed in many cases.

**Softmax Method** The softmax method (Sutton & Barto, 1998, Chapter 2) is the most popular strategy and is also termed the Boltzmann method when the exponential function is used. Recall that  $p_{ij}$  denotes the probability that the agent chooses an action  $a_j$  in a state  $s_i$ . The policy probability is defined as

$$p_{ij} \stackrel{\text{def}}{=} \pi(\beta, Q_{ij}) = \frac{\exp(\beta Q_{ij})}{Z_i(\beta)}, \quad (36)$$

where the partition function is

$$Z_i(\beta) \stackrel{\text{def}}{=} \sum_{j' \in \mathcal{A}_i} \exp(\beta Q_{ij'}). \quad (37)$$

The parameter  $\beta$  is gradually increased as  $n \rightarrow \infty$  to promote the acceptance of actions which may produce a good return. Let us denote the value of  $\beta$  at time step  $n$  by  $\beta_n$ .

**$\epsilon$ -greedy Method** In the  $\epsilon$ -greedy method (Sutton & Barto, 1998, Chapter 2), with probability  $\epsilon$ , the agent randomly chooses an action from the possible ones. On the other hand, the agent chooses the best action with the largest estimated value with probability  $1 - \epsilon$ . That is,  $p_{ij}$  is given by

$$p_{ij} \stackrel{\text{def}}{=} \pi(\epsilon, Q_{ij}) = \frac{\epsilon}{J_i} + (1 - \epsilon)\theta_{ij}, \quad (38)$$

where  $J_i \stackrel{\text{def}}{=} |\mathcal{A}_i|$  and

$$\theta_{ij} \stackrel{\text{def}}{=} \begin{cases} 1 & \text{if } j = \arg \max_{j' \in \mathcal{A}_i} Q_{ij'} \\ 0 & \text{if } j \neq \arg \max_{j' \in \mathcal{A}_i} Q_{ij'} \end{cases}. \quad (39)$$

The parameter  $\epsilon$  is gradually decreased such that  $\epsilon \rightarrow 0$  as  $n \rightarrow \infty$ . We denote the value of  $\epsilon$  at time step  $n$  by  $\epsilon_n$ .

Whether the softmax AS or the  $\epsilon$ -greedy AS is better is unclear and it may depend on the task and on human factors (Sutton & Barto, 1998, p. 31). Added to this, the explicit role of the parameters  $\beta$  and  $\epsilon$  is also unknown. In the rest of this section, we elucidate the mathematical role of the parameters and the difference between the two strategies by studying their effect in RM.

## 4.2 Stochastic Complexity

We assume that the policy is improved sufficiently slowly such that the AEP holds. Figure 4 illustrates an RL process on the manifold spanned by  $\Gamma$ . This manifold is called the information manifold (IM) (Amari & Han, 1989). Fuller explanation about the figure will be described in the following section. We use  $Q_{ij}^*$  to denote the expected value of  $Q_{ij}$  for all  $i, j$ , henceforth. Let  $p_{ij}^* = \pi(\beta, Q_{ij}^*)$  in the softmax method and  $p_{ij}^* = \pi(\epsilon, Q_{ij}^*)$  in the  $\epsilon$ -greedy method. Let  $\Gamma^{\pi^*}$  be the policy matrix whose components are given by  $p_{ij}^*$ . We define  $\Gamma^* \stackrel{\text{def}}{=} (\Gamma^{\pi^*}, \Gamma^T)$  and write the set of  $\Gamma^*$  as  $\Omega \stackrel{\text{def}}{=} \{\Gamma | \Gamma^\pi = \Gamma^{\pi^*}\}$  for notational convenience. The set  $\Omega$  is given by changing the parameter of AS strategy, such as  $\beta$  and  $\epsilon$ . The optimal policy matrix is denoted by  $\Gamma^{\pi^\dagger}$  whose components are

$$p_{ij}^\dagger = \begin{cases} 1 & \text{if } j = \arg \max_{j' \in \mathcal{A}_i} Q_{ij'}^* \\ 0 & \text{if } j \neq \arg \max_{j' \in \mathcal{A}_i} Q_{ij'}^* \end{cases}. \quad (40)$$

For example, in the softmax method we can write it as  $\Gamma^{\pi^\dagger} = \{p_{ij}^\dagger = \pi(\infty, Q_{ij}^*)\}$ , and in the  $\epsilon$ -greedy method we can also write it as  $\Gamma^{\pi^\dagger} = \{p_{ij}^\dagger = \pi(0, Q_{ij}^*)\}$ . Also, we define  $\Gamma^\dagger \stackrel{\text{def}}{=} (\Gamma^{\pi^\dagger}, \Gamma^T)$ . Figure 5 shows  $\Omega$  in each method. As shown in Figures 5(a) and 5(b) the set of matrices  $\Phi$  such that (21) holds, designated by the shaded circle, depends on  $n$  but not on  $\beta$  or  $\epsilon$ . Note that the number of elements

in  $\mathcal{C}_{\lambda_n}^n(\mathbf{\Gamma})$  depends on  $\beta$  or  $\epsilon$  because the parameter affects the value of  $H(\mathbf{\Gamma}^\pi|\mathbf{V})$ , while the set of such matrices does not depend. We assume that the neighborhood of the optimal matrix on the IM is smooth (differentiable) for the parameters of the AS strategy, such as  $\beta$  and  $\epsilon$ .

[Figure 4 about here.]

[Figure 5 about here.]

According to Lemma A.1, the number of possible conditional type matrices on the IM is determined by  $n$ ,  $I$ ,  $J$ , and  $K$ . As  $n$  increases we can create empirical sequences with arbitrarily close conditional type matrices to  $\mathbf{\Gamma}^*$ . If the environment, or specifically, the state transition matrix  $\mathbf{\Gamma}^T$  is constant,  $\mathbf{\Gamma}$  varies only with the changes of  $\mathbf{\Gamma}^\pi$ . Hence the area of possible  $\mathbf{\Phi}$  on the IM is actually restricted. Now we define a stochastic complexity that will play an important role in the later discussion.

**Definition 4.1 (Stochastic complexity)** *The stochastic complexity (SC) is defined by*

$$\psi(\mathbf{\Gamma}) \stackrel{\text{def}}{=} H(\mathbf{\Gamma}^\pi|\mathbf{V}) + H(\mathbf{\Gamma}^T|\mathbf{W}). \quad (41)$$

This is referred to as complexity because the value of  $\psi(\mathbf{\Gamma})$  is closely related to the algorithmic complexity as will be discussed in Section 4.4.

To understand the role of the SC in RL, it is worth to mention that the SC has a relationship to exploration (or exploitation) (Sutton & Barto, 1998, Chapter 2) in some cases. In short, the SC expresses the randomness of the agent’s policy. Exploration is, in general, to search for policies better than the current one, instead of the simple randomness. One efficient way for such exploratory search is to give a randomness to the policy as is done in the softmax and  $\epsilon$ -greedy methods. In this case, a policy for exploration is to enlarge the set of possible empirical sequences, that is, the  $\mathbf{\Gamma}$ -typical set in order to widely explore the environment. This is because the  $\mathbf{\Gamma}$ -typical set has probability almost one according to Theorem 3.1. On the other hand, using estimates of the action-value function the agent has to select the best action with the largest estimate of the action-value function to maximize the future return. Such a policy for exploitation is to make the  $\mathbf{\Gamma}$ -typical set smaller, so that only few empirical sequences which yield high return are allowed to be generated in practice. Thus, when the agent performs randomized exploration in AS strategy, if the value of  $\psi(\mathbf{\Gamma})$  is large, then the policy is exploratory, and analogously if the value is small, then the policy is exploitative. Of course, since the SC does not assess the rewards of empirical sequence, we have to consider both the SC and the rewards when we argue the original sense of exploration in RL. In (Iwata, Ikeda, & Sakai, 2004) the estimated entropy (like the SC) of return with respect to each state-action was formulated and a novel criterion for AS strategy was proposed by combining the estimated entropy with the estimates of the action-value function.

### 4.3 Stochastic Complexity and Return Maximization

We will show the relationship between the SC and RM in RL. We use the term RM to maximize the probability that the best sequences appear, but does not mean a lucky case in which the best sequences appear unexpectedly.

**Definition 4.2 (Return maximization)** *We denote a proper subset of best sequences by*

$$\mathcal{X}_n^\dagger \stackrel{\text{def}}{=} \{\mathbf{x} \in (\mathcal{S} \times \mathcal{A} \times \mathbb{R}_0)^n | \mathbf{\Phi}^\pi = \mathbf{\Gamma}^{\pi^\dagger} \text{ where } \mathbf{\Phi}^\pi \in \mathbf{\Lambda}_n^\pi\}. \quad (42)$$

*Then, RM means that the subset of best sequences asymptotically has probability one, that is,  $\mathcal{C}_{\lambda_n}^n(\mathbf{\Gamma}) = \mathcal{X}_n^\dagger$  as  $n \rightarrow \infty$ .*

The following theorem states that RM can be performed under a proper AS strategy so that the estimates of the value function eventually converge to the expected values.

**Theorem 4.1 (RM near the end of process)** *If the agent's policy tends to be optimal, then  $\mathcal{X}_n^\dagger \subset \mathcal{C}_{\lambda_n}^n(\mathbf{\Gamma})$  holds and the probability  $\Pr(\mathcal{X}_n^\dagger)$  of RM satisfies*

$$\Pr(\mathcal{X}_n^\dagger) \rightarrow \frac{|\mathcal{X}_n^\dagger|}{|\mathcal{C}_{\lambda_n}^n(\mathbf{\Gamma})|}, \quad (43)$$

for sufficient large  $n$ . Then, as  $n \rightarrow \infty$  and  $\mathbf{\Gamma} \rightarrow \mathbf{\Gamma}^\dagger$ ,  $\Pr(\mathcal{X}_n^\dagger) \rightarrow 1$ .

The proof is given in Appendix B.7.

Hence, we consider that  $\mathcal{X}_n^\dagger \subset \mathcal{C}_{\lambda_n}^n(\mathbf{\Gamma})$  and then reduce the  $\mathbf{\Gamma}$ -typical set such that  $\mathcal{C}_{\lambda_n}^n(\mathbf{\Gamma}) \simeq \mathcal{X}_n^\dagger$ . Here the key points are that

- by updating the estimates the agent has to improve the policy matrix  $\mathbf{\Gamma}^\pi$  as quickly as possible such that the  $\mathbf{\Gamma}$ -typical set includes the empirical sequence having the conditional type matrix  $\mathbf{\Gamma}^{\pi^*}$ , that is,

$$D(\mathbf{\Gamma}^{\pi^*} \|\mathbf{\Gamma}^\pi | \mathbf{F}_S) \leq \lambda_n, \quad (44)$$

(see Figure 4(a)), and then

- the agent is required to shut out empirical sequences except the best sequences from the  $\mathbf{\Gamma}$ -typical set in order to assign high probability to the best sequences (see Figure 4(b)).

The algorithm for the former is simply TD learning. Figure 4(a) illustrates that  $\mathbf{\Gamma}$  on the IM is refined by a TD learning such that the  $\mathbf{\Gamma}$ -typical set includes the empirical sequences having the matrix  $\mathbf{\Gamma}^*$  of the conditional types, that is, (44) holds. It is known that the convergence order of TD learning is at most  $1/\sqrt{n}$  (Kushner & Yin, 1997). After satisfying (44) the agent has to allot higher probability to the  $\mathbf{\Gamma}$ -typical set. The goal of the latter is to make the number of elements in the  $\mathbf{\Gamma}$ -typical set small while satisfying (44). This leads to the result that the subset of the best sequences occurs with high probability because according to Theorem 3.2 all the  $\mathbf{\Gamma}$ -typical sequences of  $n$  time steps have the same probability for sufficiently large  $n$ . From Theorem 3.3 we see that the number of elements in the  $\mathbf{\Gamma}$ -typical set is dependent on the SC  $\psi(\mathbf{\Gamma})$  and the quantity  $\lambda_n$ , and that the smaller each value is, the smaller the number of elements. Recall that by tuning the parameters of the AS strategy the agent can control only the SC. This leads us to the question of how sensitive the parameters such as  $\beta$  and  $\epsilon$  are for controlling the SC. The following theorems answer this question.

**Theorem 4.2 (Relationship between  $\beta$  and SC)** *The value of  $\psi(\mathbf{\Gamma})$  decreases as  $\beta$  increases. The derivative of  $\psi(\mathbf{\Gamma})$  with respect to  $\beta$  is*

$$\frac{d\psi(\mathbf{\Gamma})}{d\beta} = \sum_{i=1}^I v_i \left\{ \frac{-\beta}{2(\mathbf{Z}_i(\beta))^2} \sum_{j=1}^J \sum_{j'=1}^J (Q_{ij} - Q_{ij'})^2 \exp(\beta(Q_{ij} + Q_{ij'})) \right\}. \quad (45)$$

In particular, if  $\beta \rightarrow \infty$ , then

$$\psi(\mathbf{\Gamma}) \rightarrow \mathbf{H}(\mathbf{\Gamma}^T | \mathbf{W}). \quad (46)$$

**Theorem 4.3 (Relationship between  $\epsilon$  and SC)** *The value of  $\psi(\mathbf{\Gamma})$  decreases as  $\epsilon \rightarrow 0$ . The derivative of  $\psi(\mathbf{\Gamma})$  with respect to  $\epsilon$  is*

$$-\frac{d\psi(\mathbf{\Gamma})}{d\epsilon} = \sum_{i=1}^I v_i \left\{ \left(1 - \frac{1}{J_i}\right) \left( \log \frac{\epsilon}{J_i} - \log \left( \frac{\epsilon}{J_i} + 1 - \epsilon \right) \right) \right\}. \quad (47)$$

In particular, if  $\epsilon \rightarrow 0$ , then  $\psi(\mathbf{\Gamma})$  coincides with (46).

Theorems 4.2 and 4.3 are proved in Appendices B.8 and B.9, respectively. The equations (45) and (47) denote the sensitivity for the randomness of the policy. The main difference between the two methods is that the estimates of the action-value function affect the derivative of the SC directly in the softmax method but not in the  $\epsilon$ -greedy method. Theorems 4.2 and 4.3 draw an attention to the important dependence, often overlooked in tuning the parameter. In general, it is difficult to tune  $\beta$  and  $\epsilon$  well and the tuning forms depend only on  $n$  in the literatures. For example,  $\beta_n = c^n$  or  $\epsilon_n = c/n$  is adopted and then  $c$  is optimized by trial-and-errors, although the result of the tuning strongly depends on the values of  $v_i$ ,  $Q_{ij}$ , and  $J_i$  for every  $i, j$ , as explicitly shown in Theorems 4.2 and 4.3. In other words, one of the causes of the difficulty is that the sensitivity can not be considered on the above tuning. In fact, since all the values of  $Q_{ij}$  and  $J_i$  are available for the agent and all the value of  $v_i$  can be approximated by the values of the type, the agent can calculate the sensitivity asymptotically. Accordingly, the sensitivity may be a guide for tuning the parameters appropriately. The importance of knowing the sensitivity has been also pointed out in (Dearden, Friedman, & Russell, 1998), first.

**Example 4.1 (A guide of RM)** *The sensitivity is not something like a quantitative criterion to be directly used by itself in practical issues because of its generality. However, it can be used as a qualitative guide in choosing a tuning which depends on each case. Here, we calculate the sensitivity approximately to gain an insight into the RM speed. Let  $c$  be an arbitrary constant value. When we choose  $\epsilon_n = c/n$ , there exists a non-negative value  $c'$  such that*

$$-\frac{d\psi(\mathbf{\Gamma})}{d\epsilon} = -\sum_{i=1}^I v_i \left\{ \left(1 - \frac{1}{J_i}\right) \log \left(\frac{J_i}{c}n + 1 - J_i\right) \right\} \approx -c' \log n, \quad (48)$$

where  $n$  is sufficiently large. The value of  $c'$  can be computed by  $\{f_i, J_i\}_{i=1}^I$ . Due to  $|\mathcal{C}_{\lambda_n}^n(\mathbf{\Gamma})| \approx \exp(n\psi(\mathbf{\Gamma}))$  for sufficient large  $n$ ,

$$-\frac{d}{d\epsilon} \frac{1}{|\mathcal{C}_{\lambda_n}^n(\mathbf{\Gamma})|} = \frac{n}{|\mathcal{C}_{\lambda_n}^n(\mathbf{\Gamma})|} \frac{d\psi(\mathbf{\Gamma})}{d\epsilon} \approx c'' \frac{n \log n}{\exp(n\psi(\mathbf{\Gamma}))}, \quad (49)$$

where  $c''$  is a non-negative value which depends on  $\psi(\mathbf{\Gamma})$  and  $c'$ . This is a qualitative guide for checking the RM speed,  $d\Pr(\mathcal{X}^\dagger)/d\epsilon$ , near the optimal policy because from Theorem 4.1

$$\Pr(\mathcal{X}^\dagger) \rightarrow \frac{|\mathcal{X}^\dagger|}{|\mathcal{C}_{\lambda_n}^n(\mathbf{\Gamma})|}, \quad (50)$$

for sufficient large  $n$ . Thus, we can estimate the RM speed near the end of learning process and can select a tuning referring to it. If we think that the estimated speed is too fast for a given environment, choose more slower tuning such as  $\epsilon_n = c/\log n$ . Of course, this is a rough utility but have an interesting potential by combining it with other criteria. In the case of the softmax method, similarly, when  $\beta = c^n$ , there exists a non-negative value  $c'$  such that

$$\frac{d\psi(\mathbf{\Gamma})}{d\beta} \approx -c' c^n, \quad (51)$$

where  $n$  is sufficiently large. Then, we have

$$\frac{d}{d\beta} \frac{1}{|\mathcal{C}_{\lambda_n}^n(\mathbf{\Gamma})|} \approx c'' \frac{nc^n}{\exp(n\psi(\mathbf{\Gamma}))}, \quad (52)$$

for sufficient large  $n$ .

Next, we consider another important factor  $\lambda_n$  for making the number of elements in the  $\Gamma$ -typical set smaller. Figure 4(b) shows the changes of  $\lambda_n$  with  $n$  where the lower bound of  $\lambda_n$  is given by (23). There may be a tighter bound in various situations such that MDPs have a deterministic rule because the bound was derived under the condition that  $\Gamma$  has no constraint. In other words, the bound means a sufficient condition for RM. Hence, the first order of the bound is tightest and valid only when the agent takes “randomized” AS strategies<sup>2</sup> such that  $p_{ij} > 0$  for every  $i, j$  in the environments where  $p_{ij i'k} > 0$  for every  $i, j, i', k$ . In such cases, the bound suggests that the convergence rate of  $D(\Phi_n \| \Gamma)$  going to zero is at most  $(\log n)/n$  and its coefficient is  $(IJ + I^2JK)$ . The convergence rate indicates how fast the policy reflects on the structure of empirical sequence. The coefficient also implies that in applications a lot of time steps are required for agreement between the current matrix  $\Gamma$  and the matrix  $\Phi$  of the conditional types regarding the empirical sequence when the state, action, and reward sets are large.

#### 4.4 Stochastic Complexity and Kolmogorov Complexity

In this section we show the relationship between the SC and the Kolmogorov complexity (KC) (Cover & Thomas, 1991; Li & Vitányi, 1997). The SC is also reasonable from the point of view of algorithmic complexity. Let  $l(\mathbf{x})$  denote the length of the sequence  $\mathbf{x}$ . Let  $\mathcal{U}(q)$  be the output of a universal computer  $\mathcal{U}$  when presented with a program  $q$ . The KC of a sequence  $\mathbf{x}$  is defined as the minimal description length of  $q$  (Cover & Thomas, 1991, pp. 147–148).

**Definition 4.3 (KC and conditional KC)** *The KC  $K_{\mathcal{U}}(\mathbf{x})$  of a sequence  $\mathbf{x}$  with respect to a universal computer  $\mathcal{U}$  is defined as*

$$K_{\mathcal{U}}(\mathbf{x}) \stackrel{\text{def}}{=} \min_{q: \mathcal{U}(q)=\mathbf{x}} l(q), \quad (53)$$

*the minimum length over all programs that print  $\mathbf{x}$  and halt. If we assume that the computer already knows the length of the sequence, then we can define the conditional KC knowing  $l(\mathbf{x})$  as*

$$K_{\mathcal{U}}(\mathbf{x}|l(\mathbf{x})) \stackrel{\text{def}}{=} \min_{q: \mathcal{U}(q, l(\mathbf{x}))=\mathbf{x}} l(q). \quad (54)$$

*This is the shortest possible description length if the length of  $\mathbf{x}$  is made available to the computer  $\mathcal{U}$ .*

Since the length  $l(\mathbf{x})$  of an empirical sequence  $\mathbf{x} \in (\mathcal{S} \times \mathcal{A} \times \mathbb{R}_0)^n$  is  $3n$ , consider

$$K_{\mathcal{U}}(\mathbf{x}|3n) = \min_{q: \mathcal{U}(q, 3n)=\mathbf{x}} l(q). \quad (55)$$

Note that  $K_{\mathcal{U}}(\mathbf{x}|3n)$  denotes the algorithmic complexity to print  $\mathbf{x}$  and halt. The following theorem shows that the expected value of  $K_{\mathcal{U}}(\mathbf{x}|3n)$  is asymptotically equal to the SC.

**Theorem 4.4 (Relationship between KC and SC)** *If  $\mathbf{s} \in \mathcal{C}_{\kappa_n}^n(\mathbf{V})$ ,  $(\mathbf{s}, \mathbf{a}) \in \mathcal{C}_{\xi_n}^n(\mathbf{W})$ , and  $\mathbf{x} \in \mathcal{C}_{\lambda_n}^n(\Gamma)$ , then there exists a constant value  $c$  such that*

$$\frac{\log \nu}{n} - \rho_n \leq \frac{1}{n} \mathbb{E}_{\Gamma} [K_{\mathcal{U}}(\mathbf{x}|3n)] - \psi(\Gamma) \leq \rho_n + \frac{\log I}{n} + (IJ + I^2JK) \frac{\log(n+1)}{n} + \frac{c}{n}, \quad (56)$$

*for a computer  $\mathcal{U}$  and all  $n$ . In particular, if  $\Gamma \rightarrow \Gamma^*$  as  $n \rightarrow \infty$ , then*

$$\lim_{n \rightarrow \infty} \frac{1}{n} \mathbb{E}_{\Gamma^*} [K_{\mathcal{U}}(\mathbf{x}|3n)] = \psi(\Gamma^*). \quad (57)$$

The proof is given in Appendix B.10. The SC is so called because of this relationship.

<sup>2</sup>For example, the softmax and  $\epsilon$ -greedy methods with the parameters  $\beta < \infty$  and  $\epsilon > 0$ , respectively.

## 5 Conclusions

In this paper, we have formulated almost stationary ergodic MDPs by the type method and shown that the AEP holds on empirical sequences in such processes. Under a proper AS strategy which guarantees the convergence of the estimates, the RM is characterized by the SC  $\psi(\mathbf{\Gamma})$  and the quantity  $\lambda_n$ . We examined the role of these factors on RM and then derived the sensitivity of the SC, which is a qualitative guide in tuning the parameters of AS strategy. Also, we showed the bound of the convergence speed of the empirical sequences tending to the best sequence in the worst cases. Using the results of (Merhav, 1991; Merhav & Neuhoff, 1992) the discussions in this paper can be readily extend to the more general case where the source of empirical sequences is a unifilar source (Han & Kobayashi, 2002, p. 77) in a similar manner.

## Acknowledgments

We thank the anonymous referees for helpful comments particularly in the role of the SC and in the formulation of MDPs. This study was supported in part by a grant for scientific research (No. 14003714) from the Ministry of Education, Culture, Sports, Science, and Technology of Japan.

## A Related Theorems

We will show a number of theorems related to the AEP. Obviously, from (Csiszár & Körner, 1997, Lemma 2.2) we obtain the following lemma which plays a major role in determining the AEP.

**Lemma A.1 (Number of elements in the set of possible  $\Phi$ )** *The size of  $\Lambda_n^\pi$  is upper bounded by*

$$|\Lambda_n^\pi| \leq (n+1)^{IJ}. \quad (58)$$

Analogously,

$$|\Lambda_n^T| \leq (n+1)^{I^2 JK}. \quad (59)$$

Accordingly, the number of elements in the set of possible  $\Phi$  is upper bounded at most by a polynomial order of  $n$ , that is,

$$|\Lambda_n| \leq (n+1)^{IJ+I^2 JK}. \quad (60)$$

The following lemma states the fact that the discrepancy between the empirical entropy and the entropy asymptotically goes to zero.

**Lemma A.2** *Let  $\Phi \in \Lambda_n$  denote the matrix of the conditional types with respect to the empirical sequence which satisfies  $\mathbf{s} \in \mathcal{C}_{\kappa_n}^n(\mathbf{V})$ ,  $(\mathbf{s}, \mathbf{a}) \in \mathcal{C}_{\xi_n}^n(\mathbf{W})$ , and  $\mathbf{x} \in \mathcal{C}_{\lambda_n}^n(\mathbf{\Gamma})$ . Then, if  $\lambda_n \leq 1/8$ , we obtain*

$$|H(\Phi^\pi | \mathbf{F}_S) - H(\mathbf{\Gamma}^\pi | \mathbf{V})| \leq \sqrt{2\kappa_n} \log J - \sqrt{2\lambda_n} \log \frac{\sqrt{2\lambda_n}}{IJ}, \quad (61)$$

$$|H(\Phi^T | \mathbf{F}_{S\mathcal{A}}) - H(\mathbf{\Gamma}^T | \mathbf{W})| \leq \sqrt{2\xi_n} \log I + \sqrt{2\xi_n} \log K - \sqrt{2\lambda_n} \log \frac{\sqrt{2\lambda_n}}{I^2 JK}. \quad (62)$$

For the proof of this lemma, see Appendix B.1.

Now we show that the number of sequences with the same conditional type matrix increases exponentially for  $n$ .

**Theorem A.1 (Bound of  $|\mathcal{C}^n(\Phi, \mathbf{F}_S, \mathbf{F}_{S\mathcal{A}})|$ )** *For every state sequence  $\mathbf{s} \in \mathcal{S}^n$  with the type  $\mathbf{F}_S$  and matrix  $\Phi^\pi : \mathcal{S} \rightarrow \mathcal{A}$  such that  $\mathcal{C}^n(\Phi^\pi, \mathbf{s})$  is not empty,*

$$\frac{\exp\{nH(\Phi^\pi | \mathbf{F}_S)\}}{(n+1)^{IJ}} \leq |\mathcal{C}^n(\Phi^\pi, \mathbf{s})| \leq \exp\{nH(\Phi^\pi | \mathbf{F}_S)\}. \quad (63)$$

Also, for every action sequence  $\mathbf{a} \in \mathcal{A}^n$  and matrix  $\Phi^T : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{S} \times \mathbb{R}_0$  such that  $\mathcal{C}^n(\Phi^T, \mathbf{F}_{\mathcal{S}\mathcal{A}})$  is not empty and the joint-type is  $\mathbf{F}_{\mathcal{S}\mathcal{A}}$ ,

$$\frac{\exp\{n\mathbf{H}(\Phi^T|\mathbf{F}_{\mathcal{S}\mathcal{A}})\}}{n^{IJ}(n+1)^{I^2JK}} \leq |\mathcal{C}^n(\Phi^T, \mathbf{F}_{\mathcal{S}\mathcal{A}})| \leq I \exp\{n\mathbf{H}(\Phi^T|\mathbf{F}_{\mathcal{S}\mathcal{A}})\}. \quad (64)$$

Therefore, for every  $\mathbf{x} \in (\mathcal{S} \times \mathcal{A} \times \mathbb{R}_0)^n$  with the type  $\mathbf{F}_S$  and the joint-type  $\mathbf{F}_{\mathcal{S}\mathcal{A}}$  and for the matrix  $\Phi$ ,

$$\frac{\exp[n\{\mathbf{H}(\Phi^\pi|\mathbf{F}_S) + \mathbf{H}(\Phi^T|\mathbf{F}_{\mathcal{S}\mathcal{A}})\}]}{n^{IJ}(n+1)^{IJ+I^2JK}} \leq |\mathcal{C}^n(\Phi, \mathbf{F}_S, \mathbf{F}_{\mathcal{S}\mathcal{A}})| \leq I \exp[n\{\mathbf{H}(\Phi^\pi|\mathbf{F}_S) + \mathbf{H}(\Phi^T|\mathbf{F}_{\mathcal{S}\mathcal{A}})\}]. \quad (65)$$

The proof is given in Appendix B.2. There also exist the following bounds on the probability of the  $\Phi$ -shell.

**Theorem A.2 (Bound on probability of  $\mathcal{C}^n(\Phi, \mathbf{F}_S, \mathbf{F}_{\mathcal{S}\mathcal{A}})$ )** For every matrix  $\Phi$  such that  $\mathcal{C}^n(\Phi, \mathbf{F}_S, \mathbf{F}_{\mathcal{S}\mathcal{A}})$  is not empty,  $\Pr(\mathcal{C}^n(\Phi, \mathbf{F}_S, \mathbf{F}_{\mathcal{S}\mathcal{A}}))$  is bounded by

$$\begin{aligned} \frac{\mu \exp[-n\{\mathbf{D}(\Phi^\pi|\Gamma^\pi|\mathbf{F}_S) + \mathbf{D}(\Phi^T|\Gamma^T|\mathbf{F}_{\mathcal{S}\mathcal{A}})\}]}{n^{IJ}(n+1)^{IJ+I^2JK}} \\ \leq \Pr(\mathcal{C}^n(\Phi, \mathbf{F}_S, \mathbf{F}_{\mathcal{S}\mathcal{A}})) \leq \\ \nu^{-1} I \exp[-n\{\mathbf{D}(\Phi^\pi|\Gamma^\pi|\mathbf{F}_S) + \mathbf{D}(\Phi^T|\Gamma^T|\mathbf{F}_{\mathcal{S}\mathcal{A}})\}], \end{aligned} \quad (66)$$

where  $\nu$  and  $\mu$  are defined by (24) and (27), respectively.

The proof is given in Appendix B.3. This theorem implies that empirical sequences with conditional type matrix  $\Phi$  far from  $\Gamma$  are not likely to be generated in practice. The term ‘‘far’’ here means that the divergence between  $\Phi$  and  $\Gamma$  is large. We have mentioned the theorems to be used in the proofs of Theorem 3.1–3.3.

## B Proofs

### B.1 Proof of Lemma A.2

From (Kullback, 1967), if  $\mathbf{s} \in \mathcal{C}_{\kappa_n}^n(\mathbf{V})$ , then

$$\sum_{i=1}^I |f_i - v_i| \leq \sqrt{2\mathbf{D}(\mathbf{F}_S|\mathbf{V})}. \quad (67)$$

Hence, by  $\mathbf{H}(\mathbf{P}_{(i)}) \leq \log J$  and (15),

$$\left| \mathbf{H}(\Gamma^\pi|\mathbf{F}_S) - \mathbf{H}(\Gamma^\pi|\mathbf{V}) \right| = \left| \sum_{i=1}^I (f_i - v_i) \mathbf{H}(\mathbf{P}_{(i)}) \right| \leq \sqrt{2\kappa_n} \log J \quad (68)$$

is satisfied. In the same way as (Csiszár & Körner, 1997, Lemma 2.7), if  $\lambda_n \leq 1/8$ , then

$$\left| \mathbf{H}(\Phi^\pi|\mathbf{F}_S) - \mathbf{H}(\Gamma^\pi|\mathbf{F}_S) \right| = \left| \mathbf{H}(\Phi^\pi, \mathbf{F}_S) - \mathbf{H}(\Gamma^\pi, \mathbf{F}_S) \right| \leq -\sqrt{2\lambda_n} \log \frac{\sqrt{2\lambda_n}}{IJ}. \quad (69)$$

From

$$\left| \mathbf{H}(\Phi^\pi|\mathbf{F}_S) - \mathbf{H}(\Gamma^\pi|\mathbf{V}) \right| \leq \left| \mathbf{H}(\Phi^\pi|\mathbf{F}_S) - \mathbf{H}(\Gamma^\pi|\mathbf{F}_S) \right| + \left| \mathbf{H}(\Gamma^\pi|\mathbf{F}_S) - \mathbf{H}(\Gamma^\pi|\mathbf{V}) \right|, \quad (70)$$

we have (61). The equation (62) can be derived similarly.  $\blacksquare$



## B.2 Proof of Theorem A.1

First, we define

$$|\mathcal{C}^{n_i}(\mathbf{G}_{(i)})| \stackrel{\text{def}}{=} \frac{n_i!}{\prod_{j=1}^J n_{ij}!}. \quad (71)$$

Since the actions given by  $s_i$  have the type  $\mathbf{G}_{(i)}$ , from (Dueck & Körner, 1979) we have

$$\frac{\exp\{n_i \mathbf{H}(\mathbf{G}_{(i)})\}}{(n_i + 1)^J} \leq |\mathcal{C}^{n_i}(\mathbf{G}_{(i)})| \leq \exp\{n_i \mathbf{H}(\mathbf{G}_{(i)})\}. \quad (72)$$

By  $|\mathcal{C}^n(\Phi^\pi, \mathbf{s})| = \prod_{i=1}^I |\mathcal{C}^{n_i}(\mathbf{G}_{(i)})|$ ,  $|\mathcal{C}^n(\Phi^\pi, \mathbf{s})|$  is bounded by

$$\frac{\exp\{n \mathbf{H}(\Phi^\pi | \mathbf{F}_S)\}}{\prod_{i=1}^I (n_i + 1)^J} \leq |\mathcal{C}^n(\Phi^\pi, \mathbf{s})| \leq \exp\{n \mathbf{H}(\Phi^\pi | \mathbf{F}_S)\}. \quad (73)$$

Therefore, by  $\prod_{i=1}^I (n_i + 1)^J \leq (n + 1)^{IJ}$  we obtain (63).

The proof of (64) is different from the above proof because  $\mathcal{C}^n(\Phi^T, \mathbf{F}_{S_A})$  is not only dependent on  $\mathbf{F}_{S_A}$  but also on the initial state  $s(1)$  because of Markov property. So for any  $i, j$  we define

$$|\tilde{\mathcal{C}}^{n_{ij}}(\mathbf{G}_{(ij)})| \stackrel{\text{def}}{=} \frac{n_{ij}!}{\prod_{i'=1}^I \prod_{k=1}^K n_{ij i' k}!}. \quad (74)$$

By following along the same lines as the proof of (63) above, we have

$$\frac{\exp\{n \mathbf{H}(\Phi^T | \mathbf{F}_{S_A})\}}{(n + 1)^{I^2 JK}} \leq |\tilde{\mathcal{C}}^n(\Phi^T, \mathbf{F}_{S_A})| \leq \exp\{n \mathbf{H}(\Phi^T | \mathbf{F}_{S_A})\}. \quad (75)$$

The set  $\tilde{\mathcal{C}}^n(\Phi^T, \mathbf{F}_{S_A})$  allows a unique reconstruction of empirical sequence with  $\mathcal{C}^n(\Phi^\pi, \mathbf{s})$  only when the initial state  $s(1)$  is known. Then, from (Davisson et al., 1981), the upper bound of  $\mathcal{C}^n(\Phi^T, \mathbf{F}_{S_A})$  is

$$|\mathcal{C}^n(\Phi^T, \mathbf{F}_{S_A})| \leq I |\tilde{\mathcal{C}}^n(\Phi^T, \mathbf{F}_{S_A})| \leq I \exp\{n \mathbf{H}(\Phi^T | \mathbf{F}_{S_A})\}, \quad (76)$$

because  $s(1)$  is not specified by the  $I$  sequences. Next, in the same manner as (Davisson et al., 1981), we obtain the lower bound,

$$|\mathcal{C}^n(\Phi^T, \mathbf{F}_{S_A})| \geq \prod_{i=1}^I \prod_{j=1}^J \frac{(n_{ij} - 1)!}{\prod_{i'=1}^I \prod_{k=1}^K n_{ij i' k}!} \geq \frac{1}{n^{IJ}} |\tilde{\mathcal{C}}^n(\Phi^T, \mathbf{F}_{S_A})| \geq \frac{\exp\{n \mathbf{H}(\Phi^T | \mathbf{F}_{S_A})\}}{n^{IJ} (n + 1)^{I^2 JK}}. \quad (77)$$

Thus we have proved that (64) holds. Consequently, from (12) we obtain (65). ■

## B.3 Proof of Theorem A.2

The probability of  $\mathbf{x} \in (\mathcal{S} \times \mathcal{A} \times \mathbb{R}_0)^n$  is

$$\Pr(\mathbf{x}) = \frac{\Pr(s(1))}{\Pr(s(1) | s(n), a(n))} \prod_{t=1}^n \{\Pr(a(t) | s(t)) \Pr(s(t+1), r(t+1) | s(t), a(t))\}, \quad (78)$$

$$= \frac{\Pr(s(1))}{\Pr(s(1) | s(n), a(n))} \left( \prod_{i=1}^I \prod_{j=1}^J p_{ij}^{n_{ij}} \right) \left( \prod_{i=1}^I \prod_{j=1}^J \prod_{i'=1}^I \prod_{k=1}^K p_{ij i' k}^{n_{ij i' k}} \right), \quad (79)$$

$$= \frac{\Pr(s(1))}{\Pr(s(1) | s(n), a(n))} \exp \left[ -n \left\{ \mathbf{D}(\Phi^\pi \| \Gamma^\pi | \mathbf{F}_S) + \mathbf{D}(\Phi^T \| \Gamma^T | \mathbf{F}_{S_A}) + \mathbf{H}(\Phi^\pi | \mathbf{F}_S) + \mathbf{H}(\Phi^T | \mathbf{F}_{S_A}) \right\} \right], \quad (80)$$

where  $\Pr(s(n+1)|s(n), a(n)) = \Pr(s(1)|s(n), a(n))$  by the cyclic convention. With the definitions of (24) and (27), the probability of  $\mathbf{x}$  is bounded by

$$\begin{aligned} \mu \exp[-n \{D(\Phi^\pi \| \Gamma^\pi | \mathbf{F}_S) + D(\Phi^T \| \Gamma^T | \mathbf{F}_{S,A}) + H(\Phi^\pi | \mathbf{F}_S) + H(\Phi^T | \mathbf{F}_{S,A})\}] \\ \leq \Pr(\mathbf{x}) \leq \\ \nu^{-1} \exp[-n \{D(\Phi^\pi \| \Gamma^\pi | \mathbf{F}_S) + D(\Phi^T \| \Gamma^T | \mathbf{F}_{S,A}) + H(\Phi^\pi | \mathbf{F}_S) + H(\Phi^T | \mathbf{F}_{S,A})\}]. \end{aligned} \quad (81)$$

Using (65) and

$$\min_{\mathbf{x} \in \mathcal{C}^n(\Phi, \mathbf{F}_S, \mathbf{F}_{S,A})} |\mathcal{C}^n(\Phi, \mathbf{F}_S, \mathbf{F}_{S,A})| \Pr(\mathbf{x}) \leq \Pr(\mathcal{C}^n(\Phi, \mathbf{F}_S, \mathbf{F}_{S,A})) \leq \max_{\mathbf{x} \in \mathcal{C}^n(\Phi, \mathbf{F}_S, \mathbf{F}_{S,A})} |\mathcal{C}^n(\Phi, \mathbf{F}_S, \mathbf{F}_{S,A})| \Pr(\mathbf{x}), \quad (82)$$

we obtain (66).  $\blacksquare$

#### B.4 Proof of Theorem 3.1

Let us define the set of the matrix  $\Phi$  whose empirical sequence does not belong to the set of the  $\Gamma$ -typical sequences as

$$\Lambda'_n = \{\Phi \in \Lambda_n \mid D(\Phi^\pi \| \Gamma^\pi | \mathbf{F}_S) + D(\Phi^T \| \Gamma^T | \mathbf{F}_{S,A}) > \lambda_n\}. \quad (83)$$

Then,

$$\Pr(\mathcal{C}_{\lambda_n}^n(\Gamma)) = 1 - \Pr\left(\bigcup_{\Phi \in \Lambda'_n} \mathcal{C}^n(\Phi, \mathbf{F}_S, \mathbf{F}_{S,A})\right). \quad (84)$$

Following along the same lines as (Csiszár & Körner, 1997, Theorem 2.15) with (66) we have

$$\Pr\left(\bigcup_{\Phi \in \Lambda'_n} \mathcal{C}^n(\Phi, \mathbf{F}_S, \mathbf{F}_{S,A})\right) \leq \nu^{-1} I(n+1)^{IJ+I^2JK} \exp\left[-n \min_{\Phi \in \Lambda'_n} \{D(\Phi^\pi \| \Gamma^\pi | \mathbf{F}_S) + D(\Phi^T \| \Gamma^T | \mathbf{F}_{S,A})\}\right]. \quad (85)$$

Since  $D(\Phi^\pi \| \Gamma^\pi | \mathbf{F}_S) + D(\Phi^T \| \Gamma^T | \mathbf{F}_{S,A}) > \lambda_n$  when  $\Phi \in \Lambda'_n$ , substituting  $\lambda_n$  for the minimum value we obtain

$$\Pr\left(\bigcup_{\Phi \in \Lambda'} \mathcal{C}^n(\Phi, \mathbf{F}_S, \mathbf{F}_{S,A})\right) \leq \nu^{-1} I(n+1)^{IJ+I^2JK} \exp(-n\lambda_n), \quad (86)$$

$$= \exp\left[-n \left\{ \lambda_n - \frac{(IJ + I^2JK) \log(n+1) + \log I - \log \nu}{n} \right\}\right]. \quad (87)$$

We define

$$\varepsilon_n(I, J, K, \lambda_n) \stackrel{\text{def}}{=} \exp\left[-n \left\{ \lambda_n - \frac{(IJ + I^2JK) \log(n+1) + \log I - \log \nu}{n} \right\}\right], \quad (88)$$

and hence  $\varepsilon_n \rightarrow 0$  as  $n \rightarrow \infty$  if (23) is satisfied. Also, by (84) Theorem 3.1 holds.  $\blacksquare$

#### B.5 Proof of Theorem 3.2

First, let us derive the lower bound. We define

$$\begin{aligned} \rho_n(I, J, K, \kappa_n, \xi_n, \lambda_n) \stackrel{\text{def}}{=} \sqrt{2\kappa_n} \log J + \sqrt{2\xi_n} \log K + \sqrt{2\xi_n} \log I \\ - \sqrt{2\lambda_n} \log \frac{\sqrt{2\lambda_n}}{IJ} - \sqrt{2\lambda_n} \log \frac{\sqrt{2\lambda_n}}{I^2JK}. \end{aligned} \quad (89)$$

By (81) we have

$$-\log \Pr(\mathbf{x}) \geq \log \nu + n \{D(\Phi^\pi \| \Gamma^\pi | \mathbf{F}_S) + D(\Phi^T \| \Gamma^T | \mathbf{F}_{S\mathcal{A}}) + H(\Phi^\pi | \mathbf{F}_S) + H(\Phi^T | \mathbf{F}_{S\mathcal{A}})\}, \quad (90)$$

$$\geq \log \nu + n \{H(\Phi^\pi | \mathbf{F}_S) + H(\Phi^T | \mathbf{F}_{S\mathcal{A}})\}, \quad (91)$$

$$\geq \log \nu + n \{H(\Gamma^\pi | \mathbf{V}) + H(\Gamma^T | \mathbf{W}) - \rho_n\}, \quad (92)$$

where (91) is obtained by the non-negativity of the divergence and (92) follows from Lemma A.2.

Analogously, the upper bound is obtained as follows:

$$-\log \Pr(\mathbf{x}) \leq -\log \mu + n \{D(\Phi^\pi \| \Gamma^\pi | \mathbf{F}_S) + D(\Phi^T \| \Gamma^T | \mathbf{F}_{S\mathcal{A}}) + H(\Phi^\pi | \mathbf{F}_S) + H(\Phi^T | \mathbf{F}_{S\mathcal{A}})\}, \quad (93)$$

$$\leq -\log \mu + n \{H(\Phi^\pi | \mathbf{F}_S) + H(\Phi^T | \mathbf{F}_{S\mathcal{A}}) + \lambda_n\}, \quad (94)$$

$$\leq -\log \mu + n \{H(\Gamma^\pi | \mathbf{V}) + H(\Gamma^T | \mathbf{W}) + \lambda_n + \rho_n\}. \quad (95)$$

Thus dividing (92) and (95) by  $n$  we have (26). ■

## B.6 Proof of Theorem 3.3

We first prove the lower bound. Using the fact that  $\mathcal{C}_{\lambda_n}^n(\Gamma) \supseteq \mathcal{C}^n(\Phi, \mathbf{F}_S, \mathbf{F}_{S\mathcal{A}})$  and (65), we get

$$|\mathcal{C}_{\lambda_n}^n(\Gamma)| \geq |\mathcal{C}^n(\Phi, \mathbf{F}_S, \mathbf{F}_{S\mathcal{A}})|, \quad (96)$$

$$\geq \frac{\exp [n \{H(\Phi^\pi | \mathbf{F}_S) + H(\Phi^T | \mathbf{F}_{S\mathcal{A}})\}]}{n^{IJ(n+1)^{IJ+I^2JK}}}, \quad (97)$$

$$= \exp \left[ n \left\{ H(\Phi^\pi | \mathbf{F}_S) + H(\Phi^T | \mathbf{F}_{S\mathcal{A}}) - \frac{(IJ + I^2JK) \log(n+1) + IJ \log n}{n} \right\} \right], \quad (98)$$

$$\geq \exp [n \{H(\Gamma^\pi | \mathbf{V}) + H(\Gamma^T | \mathbf{W}) - \zeta_n\}], \quad (99)$$

where (99) is derived from Lemma A.2 and

$$\zeta_n(I, J, K, \kappa_n, \xi_n, \lambda_n) \stackrel{\text{def}}{=} \rho_n(I, J, K, \kappa_n, \xi_n, \lambda_n) + \frac{(IJ + I^2JK) \log(n+1) + IJ \log n}{n}. \quad (100)$$

Next, we consider the upper bound. By (65) and (22), we have

$$|\mathcal{C}_{\lambda_n}^n(\Gamma)| \leq \bigcup_{\substack{\Phi \in \Lambda_n: \\ D(\Phi^\pi \| \Gamma^\pi | \mathbf{F}_S) + D(\Phi^T \| \Gamma^T | \mathbf{F}_{S\mathcal{A}}) \leq \lambda_n}} |\mathcal{C}^n(\Phi, \mathbf{F}_S, \mathbf{F}_{S\mathcal{A}})|, \quad (101)$$

$$\leq I(n+1)^{IJ+I^2JK} \exp [n \{H(\Phi^\pi | \mathbf{F}_S) + H(\Phi^T | \mathbf{F}_{S\mathcal{A}})\}], \quad (102)$$

$$\leq \exp [n \{H(\Gamma^\pi | \mathbf{V}) + H(\Gamma^T | \mathbf{W}) + \eta_n\}], \quad (103)$$

where (103) is derived from Lemma A.2 and

$$\eta_n(I, J, K, \kappa_n, \xi_n, \lambda_n) \stackrel{\text{def}}{=} \rho_n(I, J, K, \kappa_n, \xi_n, \lambda_n) + \frac{(IJ + I^2JK) \log(n+1) + \log I}{n}. \quad (104)$$

Thus we have proved the upper and lower bounds in Theorem 3.3. ■

## B.7 Proof of Theorem 4.1

When the agent's return is maximized, obviously the subset of best sequences has to be included within the  $\Gamma$ -typical set. Hence,  $\mathcal{X}_n^\dagger \subset \mathcal{C}_{\lambda_n}^n(\Gamma)$  holds and then

$$\Pr(\mathcal{X}_n^\dagger) = \Pr(\mathcal{C}_{\lambda_n}^n(\Gamma)) \Pr(\mathcal{X}_n^\dagger | \mathcal{C}_{\lambda_n}^n(\Gamma)). \quad (105)$$

From Theorem 3.1, for sufficient large  $n$ ,

$$\Pr(\mathcal{C}_{\lambda_n}^n(\Gamma)) \approx 1. \quad (106)$$

Since from Theorem 3.2 every element in  $\mathcal{C}_{\lambda_n}^n(\Gamma)$  has the same probability for sufficient large  $n$ , we have

$$\Pr(\mathcal{X}_n^\dagger | \mathcal{C}_{\lambda_n}^n(\Gamma)) \rightarrow \frac{|\mathcal{X}_n^\dagger|}{|\mathcal{C}_{\lambda_n}^n(\Gamma)|}. \quad (107)$$

Therefore, (43) holds. Then, from the definition of  $\mathcal{X}_n^\dagger$ , as  $n \rightarrow \infty$  and  $\Gamma \rightarrow \Gamma^\dagger$  clearly we obtain  $\Pr(\mathcal{X}_n^\dagger) \rightarrow 1$ .  $\blacksquare$

## B.8 Proof of Theorem 4.2

Differentiating (36) with respect to  $\beta$ , we have

$$\frac{d\pi(\beta, Q_{ij})}{d\beta} = \frac{\exp(\beta Q_{ij}) \left( Q_{ij} Z_i(\beta) - \sum_{j=1}^J \{Q_{ij} \exp(\beta Q_{ij})\} \right)}{(Z_i(\beta))^2}. \quad (108)$$

Using (108) we get

$$\frac{d}{d\beta} H(\mathbf{G}_{(i)}) = - \sum_{j=1}^J \frac{d}{d\beta} (\pi(\beta, Q_{ij}) \log \pi(\beta, Q_{ij})), \quad (109)$$

$$= - \sum_{j=1}^J (\log \pi(\beta, Q_{ij}) + 1) \frac{d\pi(\beta, Q_{ij})}{d\beta}, \quad (110)$$

$$= \frac{\beta}{(Z_i(\beta))^2} \left\{ \left( \sum_{j=1}^J Q_{ij} \exp(\beta Q_{ij}) \right)^2 - Z_i(\beta) \left( \sum_{j=1}^J Q_{ij}^2 \exp(\beta Q_{ij}) \right) \right\}, \quad (111)$$

$$= - \frac{\beta}{2(Z_i(\beta))^2} \sum_{j=1}^J \sum_{j'=1}^J (Q_{ij} - Q_{ij'})^2 \exp(\beta(Q_{ij} + Q_{ij'})), \quad (112)$$

$$\leq 0. \quad (113)$$

Therefore, by  $dH(\mathbf{\Gamma}^\pi | \mathbf{V})/d\beta = \sum_{i=1}^I v_i (dH(\mathbf{G}_{(i)})/d\beta)$  we obtain (45). Also, on the limit  $\beta \rightarrow \infty$ ,  $H(\mathbf{P}_{(i)}) = 0$  holds for all  $i$ . Hence we obtain (46).  $\blacksquare$

## B.9 Proof of Theorem 4.3

Differentiating (38) with respect to  $\epsilon$ , we have

$$\frac{d\pi(\epsilon, Q_{ij})}{d\epsilon} = \frac{1}{J_i} - \theta_{ij}. \quad (114)$$

By (114) we have

$$-\frac{d}{d\epsilon}\mathbf{H}(\mathbf{G}_{(i)}) = \sum_{j=1}^J \frac{d}{d\epsilon} (\pi(\epsilon, Q_{ij}) \log \pi(\epsilon, Q_{ij})), \quad (115)$$

$$= \sum_{j=1}^J (\log \pi(\epsilon, Q_{ij}) + 1) \frac{d\pi(\epsilon, Q_{ij})}{d\epsilon}, \quad (116)$$

$$= \frac{1}{J_i} \sum_{j=1}^J \log \left( \frac{\epsilon}{J_i} + (1 - \epsilon)\theta_{ij} \right) - \sum_{j=1}^J \theta_{ij} \log \left( \frac{\epsilon}{J_i} + (1 - \epsilon)\theta_{ij} \right) - \sum_{j=1}^J \theta_{ij} + 1, \quad (117)$$

$$= \left( 1 - \frac{1}{J_i} \right) \left( \log \frac{\epsilon}{J_i} - \log \left( \frac{\epsilon}{J_i} + 1 - \epsilon \right) \right), \quad (118)$$

$$\leq 0. \quad (119)$$

Accordingly, using  $d\mathbf{H}(\mathbf{\Gamma}^\pi|\mathbf{V})/d\epsilon = \sum_{i=1}^I v_i(d\mathbf{H}(\mathbf{G}_{(i)})/d\epsilon)$  we obtain (47). Also, on the limit  $\epsilon \rightarrow \infty$ ,  $\mathbf{H}(\mathbf{P}_{(i)}) = 0$  holds for all  $i$ . Therefore we have (46).  $\blacksquare$

## B.10 Proof of Theorem 4.4

First, let us consider the lower bound. Since the length  $l(q)$  of the program  $q$  satisfies the Kraft inequality (Cover & Thomas, 1991, p. 154), we have

$$\sum_{q:\mathcal{U}(q)\text{halts}} \exp(-l(q)) \leq 1. \quad (120)$$

We assign to each  $\mathbf{x}$  the length of the shortest program  $q$  such that  $\mathcal{U}(q, 3n) = \mathbf{x} \in \mathcal{C}_{\lambda_n}^n(\mathbf{\Gamma})$ . These shortest programs also satisfy the Kraft inequality. Since the expected codeword length must be greater than the entropy, we obtain the following lower bound,

$$\mathbf{E}_{\mathbf{\Gamma}} [\mathbf{K}_{\mathcal{U}}(\mathbf{x}|3n)] \geq \mathbf{E}_{\mathbf{\Gamma}} [-\log \Pr(\mathbf{x})], \quad (121)$$

$$\geq \log \nu + n \{ \psi(\mathbf{\Gamma}) - \rho_n \}, \quad (122)$$

by (92).

Next we consider the upper bound. Let  $c$  denote a constant value. We describe the matrix  $\mathbf{\Phi} \in \mathbf{\Lambda}_n$  with respect to the empirical sequence using  $\log |\mathbf{\Lambda}_n|$  bits. Also, to describe the index of the empirical sequence within the set of all sequences having the same matrix of conditional types,  $\log |\mathcal{C}^n(\mathbf{\Phi}, \mathbf{F}_{\mathcal{S}}, \mathbf{F}_{\mathcal{S}\mathcal{A}})|$  bits are required because the set has less than  $|\mathcal{C}^n(\mathbf{\Phi}, \mathbf{F}_{\mathcal{S}}, \mathbf{F}_{\mathcal{S}\mathcal{A}})|$  elements. Hence,

$$\mathbf{K}_{\mathcal{U}}(\mathbf{x}|3n) \leq \log |\mathcal{C}^n(\mathbf{\Phi}, \mathbf{F}_{\mathcal{S}}, \mathbf{F}_{\mathcal{S}\mathcal{A}})| + \log |\mathbf{\Lambda}_n| + c, \quad (123)$$

$$\leq n \{ \mathbf{H}(\mathbf{\Phi}^\pi|\mathbf{F}_{\mathcal{S}}) + \mathbf{H}(\mathbf{\Phi}^\top|\mathbf{F}_{\mathcal{S}\mathcal{A}}) \} + \log I + (IJ + I^2JK) \log(n+1) + c, \quad (124)$$

$$\leq n \{ \psi(\mathbf{\Gamma}) + \rho_n \} + \log I + (IJ + I^2JK) \log(n+1) + c, \quad (125)$$

where (124) follows from Lemma A.1 and (65), and (125) follows from Lemma A.2. Again, taking the expectation and dividing (122) and (125) by  $n$  yields the upper bound of (56).  $\blacksquare$

## References

Amari, S., & Han, T. S. (1989, March). Statistical inference under multiterminal rate restrictions: a differential geometric approach. *IEEE Transactions on Information Theory*, 35(2), 217–227.

- Chaitin, G. J. (1977). Algorithmic information theory. *IBM journal of research and development*, 21(4), 350–359.
- Chaitin, G. J. (1987). *Algorithmic information theory* (Vol. 1). Cambridge, UK: Cambridge University Press. (reprinted with revisions in 1988)
- Cover, T. M., & Thomas, J. A. (1991). *Elements of information theory* (1st ed.). New York: John Wiley & Sons, Inc.
- Csiszár, I. (1998, October). The method of types. *IEEE Transactions on Information Theory*, 44(6), 2505–2523.
- Csiszár, I., & Körner, J. (1997). *Information theory: coding theorems for discrete memoryless systems* (3rd ed.). Budapest, Hungary: Akadémiai Kiadó. (1st impression 1981, 2nd impression 1986)
- Davissón, L. D., Longo, G., & Sgarro, A. (1981, July). The error exponent for the noiseless encoding of finite ergodic markov sources. *IEEE Transactions on Information Theory*, 27(4), 431–438.
- Dayan, P. (1992, May). The convergence of TD( $\lambda$ ) for general  $\lambda$ . *Machine Learning*, 8(3), 341–362.
- Dearden, R., Friedman, N., & Russell, S. (1998, July). Bayesian Q-learning. In *Proceedings of the 15th national conference on artificial intelligence* (pp. 761–768). Madison, Wisconsin: AAAI Press.
- Dueck, G., & Körner, J. (1979, January). Reliability function of a discrete memoryless channel at rates above capacity. *IEEE Transactions on Information Theory*, 25(1), 82–85.
- Han, T. S., & Kobayashi, K. (2002). *Mathematics of information and coding* (Vol. 203). Providence, RI: American Mathematical Society.
- Iwata, K., Ikeda, K., & Sakai, H. (2004, July). A new criterion using information gain for action selection strategy in reinforcement learning. *IEEE Transactions on Neural Networks*, 15(4), 792–799.
- Jaakkola, T., Jordan, M. I., & Singh, S. P. (1994). On the convergence of stochastic iterative dynamic programming algorithms. *Neural Computation*, 6(6), 1185–1201.
- Kullback, S. (1967, January). A lower bound for discrimination information in terms of variation. *IEEE Transactions on Information Theory*, 13(1), 126–127.
- Kushner, H. J., & Yin, G. G. (1997). *Stochastic approximation algorithms and applications* (Vol. 35). New York: Springer-Verlag.
- Li, M., & Vitányi, P. (1997). *An introduction to kolmogorov complexity and its applications* (2nd ed.). New York: Springer-Verlag. (1st edition 1993)
- Merhav, N. (1991, May). Universal coding with minimum probability of codeword length overflow. *IEEE Transactions on Information Theory*, 37(3), 556–563.
- Merhav, N., & Neuhoff, D. L. (1992, January). Variable-to-fixed length codes provide better large deviations performance than fixed-to-variable length codes. *IEEE Transactions on Information Theory*, 38(1), 135–140.
- Shannon, C. E. (1948). A mathematical theory of communication. *Bell System Technical Journal*, 27, 379–423, 623–656.
- Singh, S., Jaakkola, T., Littman, M. L., & Szepesvári, C. (2000). Convergence results for single-step on-policy reinforcement-learning algorithms. *Machine Learning*, 39, 287–308.
- Sutton, R. S., & Barto, A. G. (1998). *Reinforcement learning: An introduction*. Cambridge, MA: MIT Press.
- Tsitsiklis, J. N. (1994, September). Asynchronous stochastic approximation and Q-learning. *Machine Learning*, 16(3), 185–202.
- Watkins, C. J. C. H., & Dayan, P. (1992). Technical note: Q-learning. *Machine Learning*, 8, 279–292.
- Wolfowitz, J. (1978). *Coding theorems of information theory* (3rd ed.). Berlin: Springer-Verlag. (1st edition 1961, 2nd edition 1964)

## List of Figures

1	Interactions between the agent and the environment. . . . .	22
2	Structure of the $\Phi$ -shell. . . . .	23
3	$\Gamma$ -typical set and $\Gamma$ -typical sequence. . . . .	24
4	Matrix $\Gamma$ on the information manifold. Figure 4(a) illustrates the trajectory, drawn by updating the estimates of the action-value function using TD methods. Figure 4(b) shows the changes of $\lambda_n$ with $n$ . . . . .	25
5	Set $\Omega$ in the softmax and $\epsilon$ -greedy methods. The matrix $\Gamma^{\pi^*}$ varies with the changes of the parameter of the AS strategy, so that the set $\Omega$ is drawn as shown here. . . . .	26

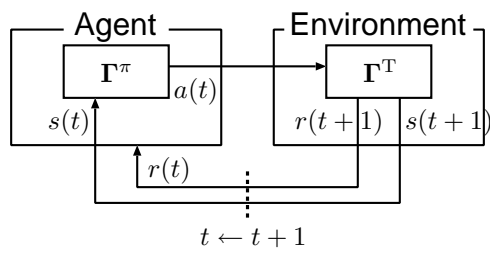


Figure 1: Interactions between the agent and the environment.



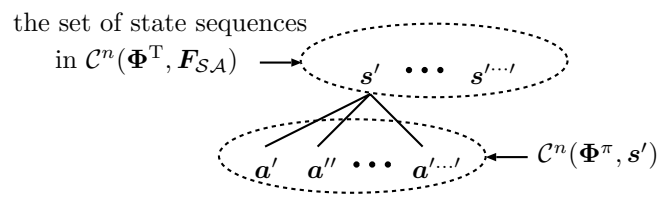


Figure 2: Structure of the  $\Phi$ -shell.

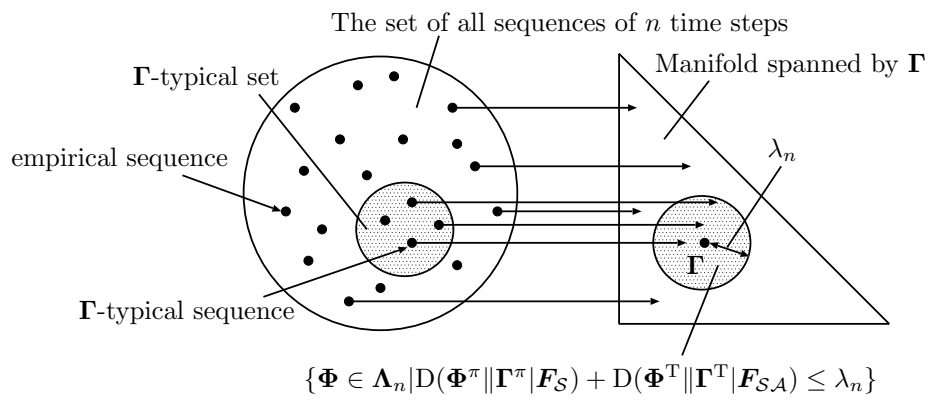
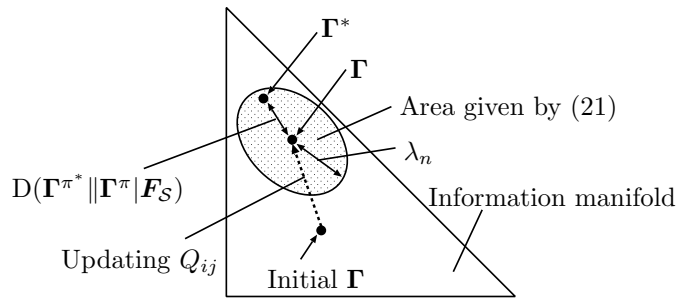
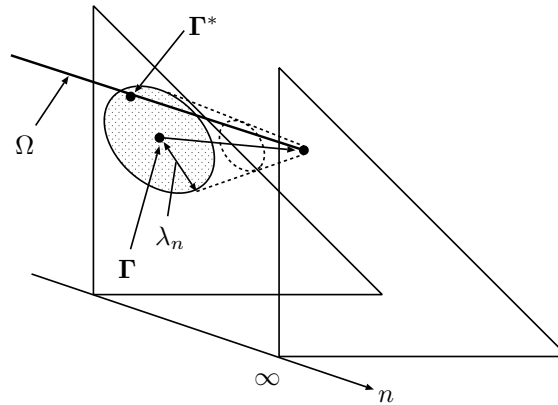


Figure 3:  $\Gamma$ -typical set and  $\Gamma$ -typical sequence.

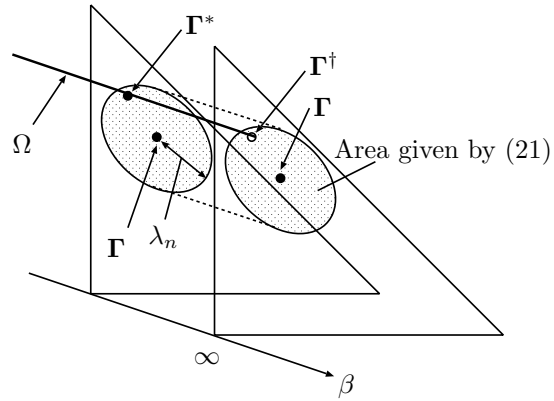


(a) Trajectory of  $\Gamma$  by updating  $Q_{ij}$ .

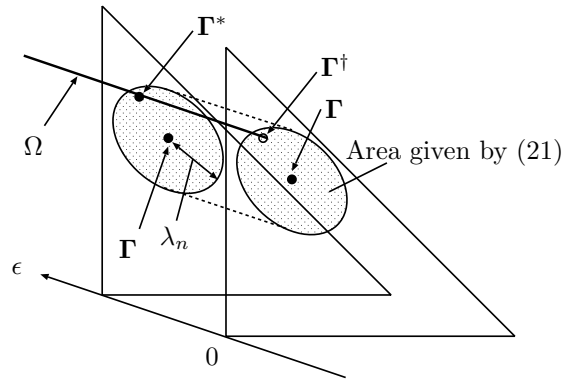


(b) Asymptotic decrease of  $\lambda_n$ .

Figure 4: Matrix  $\Gamma$  on the information manifold. Figure 4(a) illustrates the trajectory, drawn by updating the estimates of the action-value function using TD methods. Figure 4(b) shows the changes of  $\lambda_n$  with  $n$ .



(a)  $\Omega$  in the softmax method.



(b)  $\Omega$  in the  $\epsilon$ -greedy method.

Figure 5: Set  $\Omega$  in the softmax and  $\epsilon$ -greedy methods. The matrix  $\mathbf{\Gamma}^{\pi^*}$  varies with the changes of the parameter of the AS strategy, so that the set  $\Omega$  is drawn as shown here.