



The Evolution of Stochastic Automata

By

David Eric Probert

Churchill College, Cambridge

June 1976



THE EVOLUTION OF STOCHASTIC AUTOMATA

BY

DAVID ERIC PROBERT

OF

CHURCHILL COLLEGE

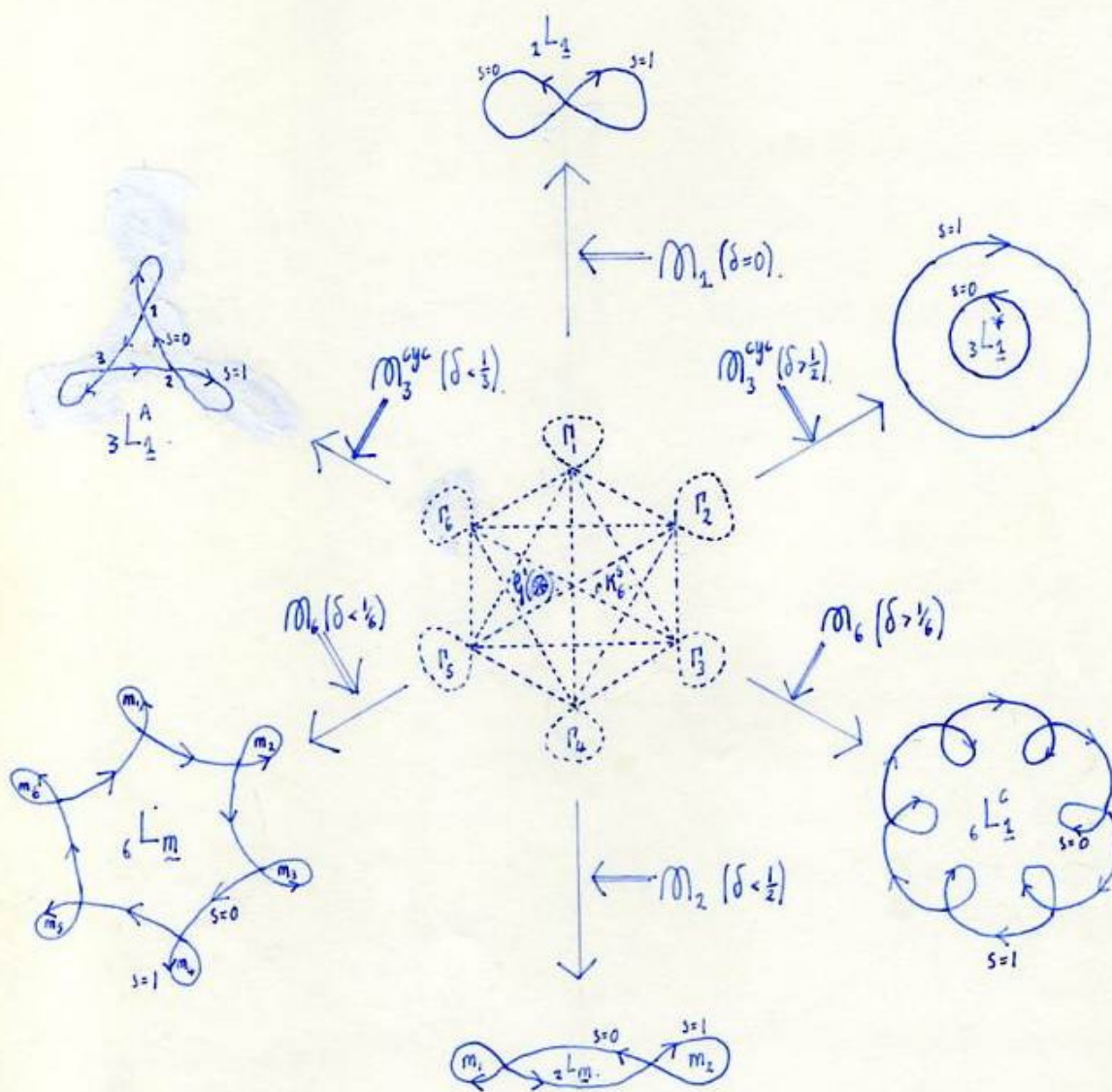
A dissertation submitted in partial fulfillment of the requirements
for the Degree of Doctor of Philosophy, at the University of Cambridge.

This dissertation is the result of my own work and includes nothing which is the outcome of work done in collaboration.

I hereby declare that this dissertation entitled,
"The Evolution of Stochastic Automata",
has not been submitted for a degree or diploma at any other university.

Signed:- D.E. Probert

Date:- 16th June 1976.



Frontispiece:-

"The Adaptation of Automaton $g'(\otimes)$
in Environments M_n ".

Contents.

0. Introduction.

- 0.1. Background, Motivation and Philosophy. 1
- 0.2. Summary. 9

1. Unstructured Automata. $\mathcal{G}^{\circ}(\otimes)$

- 1.1. The Singleton \mathbb{T} -Cell. 11
- 1.2. Uniform Learning. 12
- 1.3. Convergence. 15
- 1.4. Comparison. 18
- 1.5. Boundary Behaviour. 23
- 1.6. n-Action Extensions. 41
- 1.7. Comparison between ϵ -optimality and a.s. optimality. 48
- 1.8. The Family of \mathbb{T} -Cell Learning Rules. 51
- 1.9. Time Dependent Stimulus Probabilities. 53
- 1.10. Skeletons. 55
- 1.11. Staircases. 57
- 1.12. Dynamic Environments. 62
- 1.13. Learning Barriers. 69

2. Games between Unstructured Automata. $\prod_{i=1}^n \mathcal{G}_i^{\circ}(\otimes)$

- 2.1. The Model for \mathbb{T} -Cell Games. 72
- 2.2. Pure Saddles. 72
- 2.3. Mixed Strategies. 75
- 2.4. General Sum \mathbb{T} -Cell Games. 83
- 2.5. n-Automata Games. 85

3. Structured Automata.	$g'(\otimes)$	
3.1. The Model for \mathbb{N} -Cell Networks.		89
3.2. Static Environments.		92
3.3. Evolution in a 2-Medium.		95
3.4. Evolution in an n-Medium.		102
3.5. Relationship between A_0 and Likelihood Axis		114
	for 2-Medium.	
3.6. The n-Medium Likelihood Simplex.		116
3.7. Networks of \mathbb{N} -Cells.		119
3.8. \mathbb{N} -Cell Controllers and "Blueprint" Learning.		127
3.9. Hierarchical Automata.	$g^n(\otimes)$	130
3.10. Games between Structured Automata.	$\prod_{i=1}^n g_i'(\otimes)$	136
3.11. Concluding Remarks.		138
4. Bibliography.		

A full-page photograph of a mountain climber ascending a steep, snow-covered slope. The climber, seen from behind, wears a red beanie, a dark jacket, and a large backpack. They are using ice axes and a rope. The mountain's peak is sharp and covered in snow, set against a clear blue sky. A semi-transparent light blue box is centered over the middle of the image, containing the word "Introduction" in bold black text.

Introduction

0. Introduction.

Is there knowledge? it will vanish away; for our knowledge
and our prophecy alike are partial and the partial vanishes
when wholeness comes.

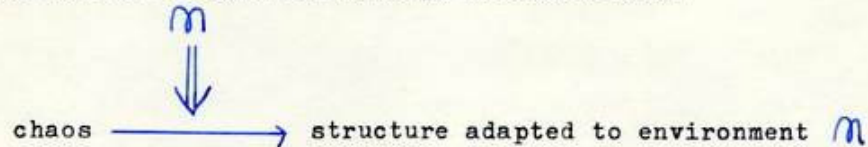
1 Corinthians 13 v 9.

0. Introduction.

1.

0.1 Background, Motivation and Philosophy.

0.1.1 During the past 25 years a large number of papers have appeared on structural adaptation through reinforcement.



Much of this work can ultimately be traced back to the ideas of Skinner (operant conditioning) in the 1930's, and Wiener (feedback and cybernetics) and Von Neumann (automata) in the 1940's, which have since developed in many disciplinary niches.

0.1.2 The approach adopted in this thesis attempts to unify the conceptual basis in a basic evolving stochastic automaton named the π -cell, defined in 1.1. This element is essentially provided by the work of Bush and Mosteller (1955), but the framework of mathematical psychology has changed little since then, and has culminated in the profound mathematical treatment of Norman (1972). Other significant contributions have been provided by Lamperti and Suppes (1960) based on the work of Luce (1959), on the β -rule. This has since been placed on a more general basis by Kanal (1962) and Marley (1967), and much of learning theory was incorporated into random systems with complete connections by Iosifescu and Theodorescu (1969).

0.1.3 Norman (1975) actually relates the Hardy-Weinberg equations of mathematical genetics to a stochastic reinforcement mechanism similar to those of psychological learning theory. Indeed, the laws of mathematical genetics of Fisher (1930) are specific reinforcement mechanisms arising from survival or extinction of genotypes, with the selective viabilities acting as the environmental stimuli.

This has a natural comparison with Skinner's (1938) operant conditioning, where reinforcement arises from the reward or penalty received on executing an action, depending stochastically on its expediency within the environment.

Both genetics and learning theory can be represented by unstructured automata, since in both fields we consider evolving distributions over genotypes or choices rather than allowing the automata to have an underlying network of transitions. In chapters 1 and 2 we consider such unstructured automata, whilst in chapter 3 we extend the theory to structured automata.

0.1.4. In 1961, the Russian mathematician Tsetlin published a pioneering paper on fixed structured automata, which opened up a new field of research linked with economic behaviour. Then in 1963 Vorontsova and Varshavskii considered the possibility of starting from an arbitrary structure, and defining reinforcement rules which give expedient adaptation to an environment. These structures evolving under reinforcement are essentially networks of n -cells with underlying digraphs (reward-penalty). It was then only possible to use computer simulations rather than theoretical techniques, since evolution by reinforcement raises many technical difficulties, which were treated by Norman (1968). The deterministic automata developed in Russia have been considered as models for biological systems, Tsetlin (1974), queuing systems and synchronisation, Varshavskii, Meleshina and Tsetlin (1965, 1968) and power regulation, Stefanyuk and Tsetlin (1967).

0.1.5. Thus, structural adaptation need not be solely biological, for we could consider the evolution of urban systems and cities, Bacon (1974), the differentiation of rôles within society, Whittle (1971), or the development of knowledge itself through research.

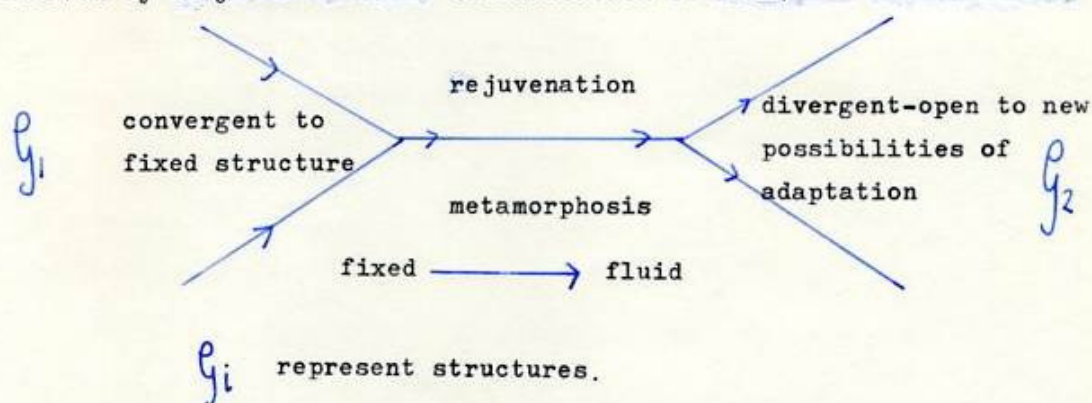
So we can consider our scientific knowledge as an evolving entity in which we abstract laws and build models from information we receive from the environment, and test these through further actions (experiments). A simple example of this is the development of the uniform polyhedra. The ancient Greeks knew of the Platonic solids, yet Kepler (1619) and subsequently Poincaré (1809) considered certain stellations and showed them to be uniform. Finally Coxeter et al (1954) enumerated 75 such uniform polyhedra with increased complexity of kernel structure and stellation depth, and Skilling (1975) showed by direct computer search that the list was indeed complete. I have taken this example because the stellation process is similar to the process of increasing "memory depth" in structured automata, whilst preserving the "SOSA" property (chapter 3), and the kernel corresponds to an action-switch (3.4.5.).

0.1.6. But all structural adaptation appears to stem ultimately from "life". The reward-penalty stimuli used in our Π -cell models are themselves derived concepts and more fundamentally we could view the actions executed by a structure, and the consequent structural reinforcement, as relating directly to its own survival. Pollution, in its broadest sense, is then that which tends to inhibit the reinforcement mechanism within the environment. Bernal (1967), Calvin (1969) and Cairns-Smith (1971) all consider how "life" could spontaneously arise on primeval earth.

Structures "feed" on -ve entropy, Schrödinger (1944), giving rise to hierarchical stability through "food" chains and ecological communities within the environment. In the absence of useful energy there is erosion and decay through the action of the 2nd law of thermodynamics. Glansdorff and Prigogine (1971) consider such an approach.

Structures may ultimately "stagnate" and become fixed, yet

they may achieve higher forms of adaptation through a metamorphosis achieved by rejuvenation. (fluidisation of form)



In chapter 3, we consider a stochastic automaton as an evolving entity with an initial random structure. Yet, through reinforcement, this probabilistic "fluid" form asymptotically becomes a deterministic "fixed" structure, acting expediently in its markovian environment. The limiting structure is not necessarily unique, so that it may be possible for the automaton to increase its payoff (adaptation) through a metamorphosis as outlined above, where fluidisation is translated mathematically as allowing state transitions to be probabilistic rather than deterministic.

0.1.7. In 1965, American control engineers Fu and Mc Murty became interested in the possibility of using learning automata as an alternative to the standard hill-climbing techniques, since virtually no prior information is required on environmental parameters and the rules themselves are simple. It is also quite easy to set several automata loose on a multimodal function in attempting to find a global maxima, as in Jarvis (1975). The paper of Fu and Mc Murty, based on the work of Vorontsova and Varshavskii (1963), initiated a new field which remained isolated from mathematical psychology until 1971. During the period 1966-71, Chandrasekaren and Shen (1967) formulated games between unstructured automata and performed

computer simulations, later extended by Viswanathan and Narendra (1974). They demonstrated oscillatory behaviour which is treated analytically in this thesis. These unstructured games are also related to the structured games of Tsetlin (1963) and the Π -cell games of chapter 2 seem a natural gaming approach.

When the engineering literature was merged with the work of Norman, several fallacious stability arguments were revealed by a counter-example of Kushner, published by Viswanathan and Narendra (1971), and are considered further by Narendra and Thathachar (1974). The stability criteria used by the engineers were only valid in a deterministic rather than a probabilistic process, when flow against a probabilistic drift is possible for all time. The optimal reinforcement rules developed in this thesis actually do only require a deterministic stability criteria which is the technique of boundary learning (1.7). So the conditions are given under which the ideas of Chandrasekaren do carry over from a probabilistic to a deterministic control theoretic framework.

0.1.8. A series of papers by Cover and Hellman (1970) on hypothesis testing by finite memory, used static automata resembling those of Tsetlin. The finite memory constraint was contested by Chandrasekaren (1970) and at present no truly satisfactory definition is to be found, apart from in the information theory of Shannon (1949).

The Π -cells themselves have ∞ memory in that a real number is held to arbitrary accuracy and similarly in the work of Cover and Hellman, a randomiser is used to generate arbitrarily small probabilities. Further, even with " ∞ computing power" we obtain undecidable propositions in logic (as with Turing's automaton), which are essentially obtained through a +ve feedback of the theory on to itself, (symbolically \odot), which gives knowledge a certain relativity.

This seems related to structural adaptation when ultimately, in setting up a mathematical framework, we may attain such intrinsic incompleteness. "Is an evolving automaton of sufficient complexity (a universal automaton) able to explain its own existence and motivate its continued evolution?" This is a deep philosophical question debated by philosophers through the ages, including Sartre and Camus in our own century.

Returning to the statistical framework, a completely independent theory of evolutionary operation was developed in America, primarily for optimising the operation of industrial chemical plants. A thorough treatment of this concept (EVOP) is covered in Box and Draper (1962 and 1969). The basic philosophy is very similar to that in this thesis in that we have operational research based on the biological theory of Darwin.



It seems possible that viewing a system, perhaps an industry, as an evolving entity in its environment, with competing systems, will prove a useful future frame of reference in operational research. This paradigm is indeed pursued in Day and Groves (1975) as a basis for future economic theory.

0.1.9. Models of brain mechanisms were initiated by Mc Culloch and Pitts in 1943, with their basic threshold neuron. This field of work has held a virtually independent existence, developing into the computer based field of pattern recognition and automatic clustering. Rosenblatt defined the α -perceptron in 1957, which was developed into the T.L.U. (threshold logic unit) used in the

monograph of Nilsson (1965). These models do exhibit structural adaptation but I believe that a probabilistic setting is more natural for abstracting the notion of environmental uncertainty.

However, this thesis does adopt an elemental approach based on the \bar{u} -cell, instead of the neuron, as the unifying element for previous work both in psychology and control engineering.

0.1.10. Cellular automata, treated rigorously by Codd (1968), have been used by Conway as a model of "life", Gardner (1971), which is deterministic yet generates unpredictable patterns.

However, in this thesis we are concerned with models that explicitly have actions executed by an automaton, with the resulting stimuli acting as the next input. In "life", we generate structural forms related to the specified rules of state transition, yet there is no structural adaptation within the environment. Such automata were formulated by Von Neumann (1948) as a model of self-reproduction. Richardson (1976) has also recently considered the self-replication of molecules, which would appear essential for structural reinforcement.

Kauffman (1969) considered cellular behaviour modelled with random genetic nets, also based on elemental "binary automata" with underlying deterministic digraph. These have been used as a practical model for learning by Aleksander (1971).

0.1.11. Recent interest in morphogenesis has been aroused by the treatise of Thom (1975), which embraces both biological and physical structural adaptation. We shall briefly consider certain similarities between boundary learning and catastrophes in 1.13, giving the theory of the evolving stochastic automaton similarities with topological morphogenesis. Indeed, in chapter 3, we shall see how a \bar{u} -cell network gives a simple model for the underlying mechanism of cellular differentiation, without any form of centralised control.

0.1.12. The work of Bush and Mosteller (1955) was based on a learning Π -cell, but the foundations have since become obscured by the analysis of reinforcement rules and parameter estimation. The Π -cell is defined to use any uniformly learning rule (1.2 and 1.8). Having returned to the basic learning entity, the concept is extended to cover:-

- a) Automata games.
- b) Adaptation in dynamic environments.
- c) Cellular differentiation.
- d) Hierarchical adaptation.

My aim in this thesis is to emphasise the new conceptual framework rather than to dwell on the rigorous mathematical derivations, which are still incomplete, particularly for Π -cell networks which require deep probabilistic ideas in their analysis.

For a deeper treatment of the background literature, I refer the reader to the excellent survey paper of Narendra and Thathachar (1974).

(1974)

0.2 Summary.

The Π -cell corresponds to an evolving unstructured automaton and gives a generalisation to the work of Norman (1972). In successive chapters we shall analyse:-

- 1). The singleton Π -cell evolving in environment \mathcal{M} .

Existing reinforcement rules mathematical psychology are either conditionally optimal or ϵ -optimal. A theory of optimal reinforcement rules is developed and their properties investigated in both static and dynamic environments. The optimality is shown to depend on a uniform learning property and behaviour near the absorbing barriers. This is in contrast to conditionally optimal rules, Luce (1959), Lamperti and Suppes (1960), which are non-uniformly learning, and ϵ -optimal rules, Norman (1968), which are centrally rather than boundary learning (1.7).

- 2). Games between Π -cells.

The Π -cell is essentially a time dependant pie graph that "adapts" to its environment \mathcal{M} . Symbolically we shall use \otimes_i to designate the i^{th} Π -cell. A game between Π -cells is now easily explained as the i^{th} Π -cell \otimes_i acting in an environment $\mathcal{M}_i = \{\otimes_{j \neq i}\}$. Each Π -cell only knows the result of its own strategy and knows nothing of the behaviour of competing Π -cells. However, 2.5.1 gives Nash point convergence when such a point exists, whilst 2.3.3 gives optimal time-averaged payoff for zero-sum games with an equilibrium point of mixed strategies, in the deterministic approximation to the automata trajectories.

- 3). Networks of Π -cells.

The network consists of a set of Π -cells with an underlying probabilistic digraph \mathcal{G} , which is described by the two markov

transition matrices σ_{ij}^s . The superscript $s \in \{0,1\}$ represents the stimulus received, and this determines which transition matrix, σ_{ij}^0 or σ_{ij}^1 to use. Initially, each automaton state is associated with a specific action, but this is developed to allow a distribution over all actions in each state. The state space is partitioned into sets of states using a particular Π -cell,

$\Pi_i = \{x_k : \otimes_i \text{ is used, where } x_k \text{ denotes state } k\}$. The mechanism of reinforcement and mode of operation of the automata is described in detail in 3.1. Limiting structures are considered in 3.2 \rightarrow 3.6.

In 3.7, the Π -cell network is considered as a model for cellular differentiation. However, it has not yet been possible to give rigorous proofs owing to the complexity of the process but it is planned to carry out a program of computer simulations in the future to guide the theoretical insight.

Finally, in 3.9, we briefly consider the concept of hierarchies of Π -cells, so that a Π -cell network is a level-1 hierarchy $\mathcal{G}^1(\otimes)$ and a Π -cell becomes a level-0 hierarchy $\mathcal{G}^0(\otimes)$. This provides a further unifying link and a basis for further research.

Another fruitful area for future research appears to be the community behaviour of automata for co-operation rather than the competition of chapter 2. This has been considered in particular by Chaikovskii (1968) and Golovchenko (1974), based on the work of Tsetlin et al (1963,1964,1965) and published in full in the collected works of Tsetlin (1974).

Community behaviour is considered in the examples of 2.5 and briefly in the "sheep effect" of 3.11.

I should like to acknowledge the help, encouragement and inspiration of my supervisor, Peter Whittle, during this period of research, which was undertaken during the tenure of an S.R.C. grant (1973-1976).



Chapter One

Unstructured Learning Automata



Chapter 1.

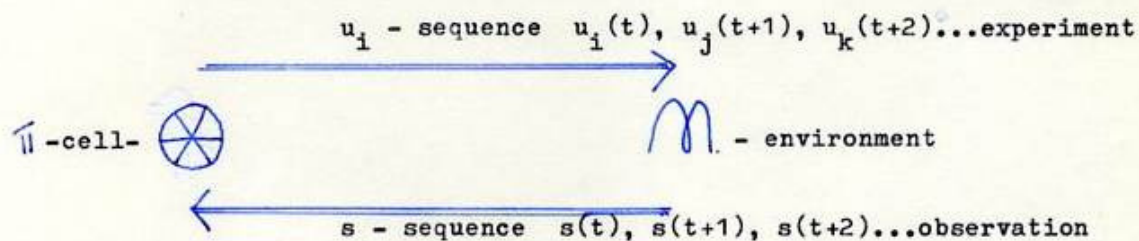
Is there a thing of which it is said, "See, this is new"?

It has been already, in the ages before us.

Ecclesiastes 1 v 10.

1. Unstructured Automata.

1.1. The Singleton Π -Cell.



Definition 1.1.1.

A Π -cell is a specific form of evolving automaton with the following properties:-



Input stimuli $s \in \{0,1\}$ $s=0$ penalty, $s=1$ reward.

Output actions u_i $1 \leq i \leq n$.

State $\Pi = (\Pi_1, \dots, \Pi_n)$, $\Pi_i(t) = \Pr(\text{output action } u_i \text{ at time } t \in \mathbb{N})$.

Transition $\left. \begin{aligned} \Pi_j(t+1) &= \Pi_j(t) + T_{ij} \text{ on receiving } s(t)=0 \\ \Pi_j(t+1) &= \Pi_j(t) + S_{ij} \text{ " " " } s(t)=1 \end{aligned} \right\} \begin{array}{l} \text{action } i \\ \text{used at} \\ \text{time } t. \end{array}$

and normalise $\sum_i \Pi_i(t+1) = 1$.

We now define the the environment $\mathcal{M}(\Delta_{\alpha\beta}, q_{u_i}^{\alpha, s})$.

$\Delta_{\alpha\beta} = \Pr(E_\alpha(t) \rightarrow E_\beta(t+1))$. where E_α is the environment state α .

$q_{u_i}^{\alpha, s} = \Pr(\text{stimulus is } s \mid u_i(t) \text{ and } E_\alpha(t))$.

and $q_{u_i}^{\alpha, 1} = q_{u_i}^{\alpha, 0} = 1 - p_{u_i}^{\alpha}$

Until section 1.12 we shall have a static $\mathcal{M} : \Delta_{\alpha\beta} = \delta_{\alpha\beta}$ (delta function)

Thus under static \mathcal{M} the environment state remains the same for all time.

Now define $R(\bar{u}(t)) = \sum_i q_i \bar{u}_i(t)$. the expected reward.

The transition rule is uniformly learning if $R(\bar{u})$ is a sub-martingale.

In 1.2 we find conditions under which T_{ij} and S_{ij} are U.L. (uniformly learning).

The transition rule is usually referred to as a reinforcement rule by the mathematical psychologists.

The rule is optimal if $\lim_{t \rightarrow \infty} \bar{u}_i(t) = 1$ only if $q_i \geq q_j$ for all j . Vorontsova (1965) considered non-linear reinforcement rules in continuous time, for 2 actions, and gave conditions for a family of rules to be optimal. However, in the survey paper of Narendra (1974), it was still unknown whether such rules gave optimal discrete time behaviour. Linear rules of the form $T_{ij} = 0$, $S_{ij} = \theta$ $S_{ii} = \theta(1 - \bar{u}_i)$ have been shown by Norman (1968) to be ϵ -optimal, that is:-

$$\lim_{\theta \downarrow 0} \lim_{t \rightarrow \infty} \bar{u}_i(t) = 1 \text{ when } q_i \geq q_j \text{ for all } j \neq i.$$

But no discrete time rules have been proven optimal independent of q .

In this chapter, we shall prove the existence of a family of optimal discrete time rules for n -actions, and show that this optimality is only dependant on boundary behaviour. This boundary learning is essential when we consider \bar{u} -cell networks and we test absorbing barriers at boundaries for probabilistic stability with respect to the optimal reinforcement rules.

1.2. Uniform Learning.

For any random variable $X(t)$, we denote the expected increment by $\Delta X(t) = (X(t+1) | X(t) \dots X(0)) - X(t)$.

Thus a rule is U.L. if $\Delta R(\bar{u}(t)) \geq 0$. Rather than state the rules that we shall be using, I shall indicate why we are restricted to a certain family, through a series of lemmas.

Theorem 1.2.1. (Whittle)

The necessary and sufficient conditions on feasible S_{ij} and T_{ij} to give $\Delta R \geq 0$ are:-

$$a). \sum_i \pi_i S_{ij} = \sum_i \pi_i T_{ij} = 0.$$

$$\text{and } b). \Delta \geq 0 \quad \text{where } \lambda_{ij} = \frac{1}{2}(\pi_i(S_{ij} - T_{ij}) + \pi_j(S_{ji} - T_{ji})).$$

Proof

$$\text{We have } \Delta R = \sum_i q_i \Delta \pi_i = \sum_i q_i (\pi_j q_j S_{ji} + \pi_j p_j T_{ji}).$$

$$= q' \Lambda q + \sum_j (\sum_i \pi_i T_{ij}) q_j. \quad i$$

$$= p' \Lambda p + \sum_j (\sum_i \pi_i S_{ij}) p_j. \quad ii$$

For necessity, we just consider a small perturbation ϵ from $q_i = q_j \forall i, j$. Put $q = c1 + \epsilon$, then $\Delta R = \epsilon' \Lambda \epsilon + \sum_j (\sum_i \pi_i T_{ij}) \epsilon_j$ as $1' \Lambda 1 = 0$ and $\sum_j S_{ij} = 0$ by normalisation. Now to prove $\Delta R \geq 0$, we just require ϵ sufficiently small and then we choose ϵ to give $\sum_j (\sum_i \pi_i T_{ij}) \epsilon_j < 0$ which is always possible if $\sum_i \pi_i T_{ij} \neq 0$. (Note that $\epsilon' \Lambda \epsilon \sim 0$ and can be neglected w.r.t. $\sum_i \pi_i T_{ij} \epsilon_j$.) Hence $\sum_i \pi_i T_{ij} = 0$ and $\Delta \geq 0$ are necessary. Similarly $\sum_i \pi_i S_{ij} = 0$ from ii above. The conditions above clearly sufficient.

Lemma 1.2.2.

$$U.L. \Rightarrow (\Delta \pi_k = 0 \text{ for all } k \text{ if } q_i = q_j \forall i, j).$$

Proof. We let $q_i = q_k$ for all i and evaluate $\Delta \pi_k$.

$$\Delta \pi_k = q_k \sum_i \pi_i S_{ik} + p_k \sum_i \pi_i T_{ik} = 0 \text{ by condition a) above, and}$$

we actually also have:- (If $q = c1$ for any $c \Rightarrow \Delta \pi_k = 0 \forall k \Leftrightarrow a$), in 1.2.1.

Lemma 1.2.3.

The linear rule $T_{ij} = 0$, $S_{ij} = -\theta_{ij} \pi_j$, $S_{ii} = \theta_{ii} (1 - \pi_i)$ is U.L. iff $\theta_{ij} = \theta \forall i, j$.

Proof.

$\sum_i \pi_i = 1 \Rightarrow \sum_i (\theta_{ii} - \theta_{ij}) \pi_j = 0$ and hence $\theta_{ii} = \theta_{ij} \forall i, j$ and lemma 1.2.2 gives $\sum_j (\theta_{jj} - \theta_{ij}) \pi_i = 0$ and so we have $\theta_{jj} = \theta_{ij} \forall i, j$. Combining the above, $\theta_{ij} = \theta = \text{constant}$.

$$\text{And } \Delta R = \sum_{i,j} \frac{1}{2} \theta \pi_i (q_i - q_j)^2 \pi_j \geq 0 \text{ gives sufficiency.}$$

We cannot extend 1.2.3 to give $U.L. \Rightarrow \theta_{ij}(\bar{u}) = \theta(\bar{u}) \forall i, j$, for non-linear rules but the final three lemmas in this section go as far as we can.

Lemma 1.2.4.

The non-linear rule $T_{ij} \equiv 0$, $S_{ij} = -\theta_i(\bar{u})\bar{u}_j$ and $S_{ii} = \theta_i(\bar{u})(1 - \bar{u}_i)$ is U.L. iff $\theta_i(\bar{u}) = \theta(\bar{u}) \forall i$.

Proof.

$U.L. \Rightarrow \sum_i \bar{u}_i S_{ij} = 0$ and hence $\sum_j \bar{u}_j (\theta_i(\bar{u}) - \theta_j(\bar{u})) = 0 \forall i$.

Hence $\theta_i(\bar{u}) = \theta_j(\bar{u}) \forall i, j$ and $\Delta R = \sum_{i,j} \frac{1}{2} \theta(\bar{u}) \bar{u}_i (q_i - q_j)^2 \bar{u}_j$ gives sufficiency.

Lemma 1.2.5.

If we restrict the no of actions to $n = 2$, then the non-linear rule $T_{ij} = 0$, $S_{ij} = -\theta_{ij}(\bar{u})\bar{u}_j$ $i \neq j$ and $S_{ii} = \theta_{ii}(\bar{u})(1 - \bar{u}_i)$ is U.L. iff $\theta_{ij}(\bar{u}) = \theta(\bar{u})$ i, j .

Proof.

We have sufficiency as in 1.2.4 and for the necessity.

$\sum \bar{u}_i = 1$ gives $\theta_{11} = \theta_{12}$ and $\theta_{22} = \theta_{21}$

whilst $\sum_i \bar{u}_i S_{ij} = 0$ gives $\theta_{11} = \theta_{21}$ and $\theta_{22} = \theta_{12}$

Hence we have

$$\theta_{ij}(\bar{u}) = \theta(\bar{u}) \quad \forall i, j.$$

Lemma 1.2.6.

The non-linear rule $T_{ij} \equiv 0$, $S_{ij} = -\theta_{ij}(\bar{u})\bar{u}_j$ $i \neq j$ and $S_{ii} = \theta_{ii}(\bar{u})(1 - \bar{u}_i)$ $n \geq 2$ is U.L. if $\theta_{ij}(\bar{u}) = \theta_{ji}(\bar{u})$ and $\theta_{ii}(\bar{u}) = \sum_j \theta_{ij}(\bar{u}) \bar{u}_j$ (normalisation).

Proof.

$$\begin{aligned} \Delta R &= \sum_i q_i \Delta \bar{u}_i(t) = \sum_i q_i \bar{u}_i \sum_j \bar{u}_j (\theta_{ii} q_i - \theta_{ij} q_j) \\ &= \sum_i \theta_{ii} q_i^2 \bar{u}_i - \sum_{i,j} \theta_{ij} q_i q_j \bar{u}_i \bar{u}_j = \sum_{i,j} \theta_{ij} q_i \bar{u}_i \bar{u}_j (q_i - q_j) \\ &= \frac{1}{2} \sum_{i,j} \theta_{ij}(\bar{u}) \bar{u}_i \bar{u}_j (q_i - q_j)^2 \geq 0 \quad 0 < \theta_{ij}(\bar{u}) < 1. \end{aligned}$$

and hence the rule is U.L. and normalisation gives $\theta_{ii}(\bar{u}) = \sum_j \theta_{ij}(\bar{u}) \bar{u}_j$

Actually it is necessary to have $\sum_i \bar{u}_i \theta_{ij} = \sum_j \theta_{ij} \bar{u}_j$ using $\sum_i \bar{u}_i S_{ij} = 0$.

Thus asymptotically as $\bar{u}_i \rightarrow 1$ we must have $\theta_{ij}(\bar{u}) \rightarrow \theta_{ji}(\bar{u})$. So 1.2.6. is really the best practical rule we can formulate.

U.L. is sufficient for later theorems but not necessary throughout $\bar{u}_i \in I = [0,1]$ and optimality of the $\theta_{ij}(\bar{u})$ rules will be characterised by certain boundary properties. In 1.6.7. we briefly discuss the possible use of T_{ij} in addition to S_{ij} . First we consider the convergence of U.L. rules, which follows from semi-martingale theorems.

1.3. Convergence.

To ensure that $\exists i$ s.t. $\lim_{t \rightarrow \infty} \bar{u}_i(t) = 1$, we shall use the symmetric $\theta_{ij}(\bar{u})$ rules which are the natural extension of the 2-action $\theta_{ij}(\bar{u}) = \theta(\bar{u})$ rules. Under the same condition, $\theta(\bar{u}) = 0$ if and only if some $\bar{u}_i = 1$, we obtain boundary convergence for $\theta(\bar{u})$ family, but apart from 2-actions, conditions for optimality are still unknown. e.g. $\theta(\bar{u}) = \prod_{i=1}^n (1 - \bar{u}_i)$. All U.L. rules converge by s/mg theorems, but the difficulty lies in obtaining a.s. absorption in boundary, $\bar{u}_i \rightarrow 1$ for some i .

On taking action i at time t , $u_i(t)$.

$$\bar{u}_j(t+1) = \bar{u}_j(t)(1 - \theta_{ij}(\bar{u})) + \delta_{ij} \theta_{ij}(\bar{u}) \quad s(t)=1$$

$$\bar{u}_j(t+1) = \bar{u}_j(t). \quad s(t)=0$$

We denote this family of rules by \mathcal{R} , where $\theta_{ij}(\bar{u})$ iff $\bar{u}_i = 0$ or $\bar{u}_j = 0$.

Theorem 1.3.1.

Under \mathcal{R} , $\lim_{t \rightarrow \infty} \bar{u}_i(t) \in \{0,1\} \forall i$.

Proof

First we must prove that the limit exists.

$\Delta R \geq 0$ and hence $\lim_{t \rightarrow \infty} R(\bar{u}(t)) \rightarrow R^*$ by s/martingale theorem.

And $\Delta \bar{u}_k = \bar{u}_k \sum \theta_{kj}(\bar{u}) (\bar{u}_j (q_k - q_j))$ for k s.t. $q_k \geq q_j \forall j$.

Since \bar{u}_k is bounded $\bar{u}_k \xrightarrow{a.s.} v_k \in [0,1]$ by s/mg theorem.

a). First suppose $\Delta \bar{u}_k > 0$ and $v_k \in (0,1)$

Now by Doob (1953), as $\{\bar{u}_k(t)\}$ is uniformly integrable, we have

$$\lim_{t \rightarrow \infty} E |\bar{u}_k(\infty) - \bar{u}_k(t)| = 0.$$

i

But $E(\bar{u}_k(t+1) | \bar{u}_k(t)) - \bar{u}_k(t) = \Delta \bar{u}_k > 0$ for $\bar{u}_k \in (0, 1)$

and so $E(\bar{u}_k(t+1) | \bar{u}_k(0)) - E(\bar{u}_k(t) | \bar{u}_k(0)) = \int_{\Omega} \Delta \bar{u}_k(t) d\mu > 0$

and $\lim_{t \rightarrow \infty} \Delta \bar{u}_k(t) = \lim_{t \rightarrow \infty} \int_{\Omega} \Delta \bar{u}_k(t) d\mu > 0$

ii

(see Sawaragi and Baba (1974) for $(\Omega, \mathcal{F}, \mu)$).

where the limit exists as $\Delta \bar{u}_k$ converges a.s. and is bounded absolutely

But $\lim_{t \rightarrow \infty} E(\bar{u}_k(t)) = E(\bar{u}_k(\infty)) = v_k$ by i and so $\lim_{t \rightarrow \infty} \Delta \bar{u}_k(t) = 0$

We have a contradiction unless $v_k \in \{0, 1\}$.

b) Now let $\Delta \bar{u}_k \equiv 0$, which for U.L. rules occurs iff $q_i = q_j \forall j$.

Again we use Doob, since $|\bar{u}_k(t)|^2$ is U.I. $\lim_{t \rightarrow \infty} |\bar{u}_k(\infty) - \bar{u}_k(t)|^2 = 0$

But the variance of increment $= E(\bar{u}_k(t+1) - \bar{u}_k(t) - \Delta \bar{u}_k(t))^2$
 $= E(\bar{u}_k(t+1) - \bar{u}_k(t))^2 > 0$ for $\bar{u}_k \in (0, 1)$.

and $E(\bar{u}_k(t+1) - \bar{u}_k(t))^2 | \bar{u}_k(t) = q_k \bar{u}_k [(q_{kk} (1 - \bar{u}_k))^2 + \sum_{i \neq k} q_{ik}^2 \bar{u}_i \bar{u}_k]$

So $\lim_{t \rightarrow \infty} E(\bar{u}_k(t+1) - \bar{u}_k(t))^2 | \bar{u}_k(0) = \lim_{t \rightarrow \infty} \int_{\Omega} q_k \bar{u}_k [(q_{kk} (1 - \bar{u}_k))^2 + \sum_{i \neq k} q_{ik}^2 \bar{u}_i \bar{u}_k] d\mu$
 > 0 as $v_k \in (0, 1)$.

Yet $\lim_{t \rightarrow \infty} E(\bar{u}_k(t+1)) = E(\bar{u}_k(\infty))$ and we have a contradiction.

Hence both conditional expectation and variance vanishing give us boundary convergence of \bar{u}_k .

c) If the rule is optimal, we have $\bar{u}_i \rightarrow 1$ only if $q_i \geq q_j \forall j$ by defn.

d) Now re-order suffices to give $q_1 \geq q_2 \geq \dots \geq q_n$, to give

$\bar{u}_i \rightarrow v_i \in \{0, 1\} \forall i$. Note that $q_n \leq q_j \quad j$ gives $\bar{u}_n \rightarrow v_n \in \{0, 1\}$.

by the same reasoning as a).

Suppose $\lim_{t \rightarrow \infty} \bar{u}_1(t) = 1$ then result is immediate

Thus let $\bar{u}_1(t) \rightarrow 0$ and consider $\text{sgn}(\Delta \bar{u}_2)$

We have $\Delta \bar{u}_2 = \sum_j \bar{u}_2 q_{2j} (q_j - q_1) \bar{u}_j$ and $\text{sgn}(\Delta \bar{u}_2)$ can only alternate infinitely often if $\lim_{t \rightarrow \infty} \bar{u}_2(t) \in \{0, 1\}$.

Suppose $\text{sgn}(\Delta \bar{u}_2)$ alternates finitely often, then for $t > t^*$ say,

we have $\bar{\pi}_i(t)$ is a s/martingale and $\bar{\pi}_i \rightarrow v_i \in [0,1]$ and by a) and b) $v_i \in \{0,1\}$.

Now by induction, if $\lim_{t \rightarrow \infty} \bar{\pi}_i(t) = 0$ for $0 < i \leq r$ we require $\lim_{t \rightarrow \infty} \bar{\pi}_{r+1}(t) \in \{0,1\}$, which is easily proved by adapting above. If $\lim_{t \rightarrow \infty} \bar{\pi}_{r+1}(t) = 1$ we are done, else continue until all suffixes are exhausted, noting if $\lim_{t \rightarrow \infty} \bar{\pi}_i(t) = 0$, $0 < i \leq n-1$, then $\bar{\pi}_n(t) \rightarrow 1$ by normalisation. //

An alternative proof of a) and b) could be based on #up-crossings over any rational interval is $\stackrel{a.s.}{\rightarrow} \infty$, Breiman (1968), and hence to prove $v_i \in \{0,1\}$, we need only note conditional variance > 0 at $v_i \in (0,1)$ and so $\bar{\pi}_i \rightarrow v_i \in (0,1) \Rightarrow \# \text{ up crossings } \stackrel{a.s.}{\rightarrow} \infty$ across any sufficiently short interval $[v_i - \epsilon, v_i + \epsilon]$ giving contradiction.

Since optimality will be shown to arise from $\theta_{ij}(\bar{\pi}) \downarrow 0$ as $\bar{\pi}_i \downarrow 0$ or $\bar{\pi}_j \downarrow 0$, sufficiently fast, it may be thought that $\theta(t)$ would also give optimal rules. However, even though $\bar{\pi}_i \rightarrow v_i \in [0,1]$ there is now no reason to restrict $v_i \in \{0,1\}$; we may only have time dependence occurring through $\bar{\pi}(t)$ to give our $\theta_{ij}(\bar{\pi}(t))$ rules. Similarly for $\theta(\bar{\pi}(t))$ rules over > 2 actions; if $\theta(\bar{\pi}(t)) \downarrow 0$ as $\bar{\pi}_i \downarrow 0$, then we need not have boundary absorption.

Corollary 1.3.2.

Under U.L. $\theta(\bar{\pi})$ rules with $\theta(\bar{\pi}) \downarrow 0$ iff $\bar{\pi}_i \uparrow 1$ for some i , then $\lim_{t \rightarrow \infty} \bar{\pi}_k(t) \in \{0,1\} \quad \forall k$.

Proof

Use a) and b) of 1.3.1, which hold since conditional variance only vanishes at boundary. Indeed b) is sufficient but it is a) which brings in the concept of probabilistic drift, which is the key to discussions of optimality. //

In the next sections, we shall initially just consider 2-actions which are extended to $n > 2$ in 1.6. Laksmirvarahan and Thathachar (1973) considered uniformly learning $\theta(\bar{\pi})$ rules over n -actions, yet seemed unaware that they are not necessarily absorbed in $\bar{\pi}_i = 1$.

Theorem 1.4.1

If $\gamma_{\theta(\bar{u})}(\bar{u}) = \Pr(\lim_{t \rightarrow \infty} \bar{u}_t(i) = 1 \mid \bar{u}_0(i) = \bar{u}_i)$ then $U_{\theta(\bar{u})} \gamma_{\theta(\bar{u})} = \gamma_{\theta(\bar{u})}$
 and if $\gamma_{\theta(\bar{u})}$ is continuous, it is unique.

Proof.

We know by 1.3.1 that $\lim_{t \rightarrow \infty} \bar{u}_t(i) \in \{0, 1\}$, and hence iterating $U_{\theta(\bar{u})}$ we obtain $\lim_{n \rightarrow \infty} U_{\theta(\bar{u})}^n \gamma_{\theta(\bar{u})} = \sum (\gamma_{\theta(\bar{u})}(\bar{u}_t(i)) \mid \bar{u}_0(i) = 1) \gamma_{\theta(\bar{u})}(\bar{u}_t(i))$ and we have the result from $\lim_{n \rightarrow \infty} U_{\theta(\bar{u})}^n \gamma = U \lim_{n \rightarrow \infty} U_{\theta(\bar{u})}^{n-1} \gamma = U \gamma = \gamma$.
 Indeed, if $\gamma(i) = i$ for boundary $i \in \{0, 1\}$, then $\lim_{n \rightarrow \infty} U_{\theta(\bar{u})}^n \gamma(\bar{u}) = \gamma_{\theta(\bar{u})}(\bar{u})$ easily, by the same argument.

Finally, let $\gamma_{\theta(\bar{u})}^*$ be another solution to $U_{\theta(\bar{u})} \gamma^* = \gamma^*$ which is continuous. Then $\lim_{n \rightarrow \infty} U_{\theta(\bar{u})}^n \gamma^* = \gamma^*$, by continuity we also have $\lim_{n \rightarrow \infty} U_{\theta(\bar{u})}^n \gamma_{\theta(\bar{u})}^* = \gamma_{\theta(\bar{u})}^*$ and this gives uniqueness. We could write this fundamental equation more concisely as $\Delta \gamma_{\theta(\bar{u})}(\bar{u}) \equiv 0$. //

Definition. 1.4.2.

A function $\gamma(\bar{u})$ on $[0, 1]$ with $\gamma(i) = i$ at $i \in \{0, 1\}$ is super-regular (sub-regular) if $\gamma(\bar{u}) \geq (\leq) U \gamma(\bar{u}) \quad \forall \bar{u} \in [0, 1]$.
 (i.e. $\Delta \gamma(\bar{u}) \leq (\geq) 0$)

We may find that $U_{\theta(\bar{u})}$ has no continuous solutions, in particular, if the rule is optimal, we have $\gamma_{\text{opt}}(\bar{u}) \equiv 1$ for $\bar{u} \in (0, 1)$, $\gamma(0) = 0$.

To prove that $U_{\theta(\bar{u})}$ is the operator for an optimal rule, we construct a continuous sub-regular family $\gamma_\delta(\bar{u})$ with $\lim_{\delta \downarrow 0} \gamma_\delta(\bar{u}) = \gamma_{\text{opt}}(\bar{u})$ the discontinuous limit. Then since $\lim_{n \rightarrow \infty} U_{\theta(\bar{u})}^n \gamma_\delta = \gamma_\delta \quad \forall \delta$, the sub-regularity $U \gamma_\delta \geq \gamma_\delta$ ensures $\gamma_{\theta(\bar{u})}(\bar{u}) = \gamma_{\text{opt}}(\bar{u})$. The sub-regular family would give us a contradiction if we asserted any other solution, apart from γ_{opt} , as giving the absorption probabilities.

Now we prove the convexity of $\gamma_\theta(\bar{u})$, for the linear rule, with $q_2 \gg q_1$.

We have $U_\theta \phi(\bar{u}) = (\bar{u}_1 q_1 + \bar{u}_2 q_2) \phi(\bar{u}) + \bar{u}_1 q_1 \phi(T_1 \bar{u}) + \bar{u}_2 q_2 \phi(T_2 \bar{u})$
 where $T_2 \bar{u} = (1 - \theta)$ and $T_1 \bar{u} = T_2 \bar{u} + \theta$.

or
$$U_\theta \phi(\bar{u}) - \phi(\bar{u}) = \bar{u}_1 q_1 (\phi(\tau_1 \bar{u}) - \phi(\bar{u})) + \bar{u}_2 q_2 (\phi(\tau_2 \bar{u}) - \phi(\bar{u}))$$
 with $\phi(i) = \bar{u} \quad i \in \{0, 1\}$.

and also define $V_{\bar{u}}$ by,

$$V_{\bar{u}} \phi(\bar{u}) = (\bar{u}_1 q_1 + \bar{u}_2 q_2) \phi(\bar{u}) + \bar{u}_1 q_1 \phi(\tau_1 \bar{u}) + \bar{u}_2 q_2 \phi(\tau_2 \bar{u})$$

with $U \phi(\bar{u}) = V_{\bar{u}} \phi(\bar{u})$.

Lemma 1.4.3.

$\gamma_\theta(\bar{u})$ is monotone increasing on $\bar{u} \in [0, 1]$

Proof.

Take ϕ as any monotone increasing function with $\bar{u}' \leq \bar{u}'' \Rightarrow \phi(\bar{u}') \leq \phi(\bar{u}'')$
 then we show the same is true for $U \phi$ and hence γ_θ , as $\lim_{n \rightarrow \infty} U^n \phi = \gamma$.

$$\begin{aligned} \text{Now } U_\theta \phi(\bar{u}'') - U_\theta \phi(\bar{u}') &= V_{\bar{u}''} \phi(\bar{u}'') - V_{\bar{u}'} \phi(\bar{u}') \\ &\geq V_{\bar{u}'} \phi(\bar{u}'') - V_{\bar{u}'} \phi(\bar{u}') \geq 0 \end{aligned} \quad (*)$$

where * follows from monotonicity of ϕ and the following:-

$$V_{\bar{u}''} \phi(\bar{u}'') - V_{\bar{u}'} \phi(\bar{u}') = (\bar{u}_1'' - \bar{u}_1') q_1 (\phi(\tau_1 \bar{u}'') - \phi(\bar{u}')) + (\bar{u}_2'' - \bar{u}_2') q_2 (\phi(\tau_2 \bar{u}'') - \phi(\tau_2 \bar{u}')) \geq 0 \quad //$$

Lemma 1.4.4.

$\gamma_\theta(\bar{u})$ is convex on $\bar{u} \in [0, 1]$

Proof.

Assume $q_1 < q_2$; as $\gamma_\theta(\bar{u}) = \bar{u}$ at $q_1 = q_2$.

Then we take ϕ convex and show that $U \phi$ and hence γ_θ is convex.

Take $\alpha, \beta \geq 0 \quad \alpha + \beta = 1$. Then for $\bar{u}' < \bar{u}''$.

$$\begin{aligned} \alpha V \phi(\bar{u}') + \beta V \phi(\bar{u}'') &= V \phi(\alpha \bar{u}' + \beta \bar{u}'') = \alpha V_{\bar{u}'} \phi(\bar{u}') + \beta V_{\bar{u}''} \phi(\bar{u}'') - V_{\alpha \bar{u}' + \beta \bar{u}''} \phi(\alpha \bar{u}' + \beta \bar{u}'') \\ &\geq \alpha V_{\bar{u}'} \phi(\bar{u}') + \beta V_{\bar{u}''} \phi(\bar{u}'') - V_{\alpha \bar{u}' + \beta \bar{u}''} (\alpha \phi(\bar{u}') + \beta \phi(\bar{u}'')) \quad \text{by convexity of} \\ &= \alpha V_{\bar{u}'} \phi(\bar{u}') + \beta V_{\bar{u}''} \phi(\bar{u}'') - (\alpha V_{\bar{u}'} + \beta V_{\bar{u}''}) (\alpha \phi(\bar{u}') + \beta \phi(\bar{u}'')) \quad \text{and linearity of } T_j. \\ &= \alpha \beta (V_{\bar{u}''} - V_{\bar{u}'})(\phi(\bar{u}'') - \phi(\bar{u}')) = \alpha \beta \delta D \quad \text{where } \delta = \bar{u}_1'' - \bar{u}_1' = \bar{u}_2' - \bar{u}_2'' > 0 \\ \text{and } D &= q_2 [\phi(\tau_2 \bar{u}'') - \phi(\bar{u}') - \phi(\tau_2 \bar{u}') + \phi(\bar{u}'')] - q_1 (\phi(\tau_1 \bar{u}'') - \phi(\bar{u}') - \phi(\tau_1 \bar{u}') + \phi(\bar{u}'')) \\ &= q_2 A_2 - q_1 A_1 \quad \text{say and } D \geq 0 \text{ if } A_2 \geq A_1 \text{ or} \\ &\quad \phi(\tau_1 \bar{u}'') - \phi(\tau_2 \bar{u}') \leq \phi(\tau_1 \bar{u}') - \phi(\tau_2 \bar{u}'') \quad (**) \end{aligned}$$

But $T_1 \bar{u} - T_2 \bar{u} = \theta = \text{const}$ and $\bar{u}' \leq \bar{u}''$ and hence (**) follows by convexity.

//

Theorem 1.4.5.

If $\theta(\bar{u}) \leq \theta$ then $U_{\theta}(\bar{u}) \leq \bar{\gamma}_{\theta}$.

Proof

We use $\theta^* = \theta(\bar{u})$ and note that $\bar{\gamma}_{\theta}$ convex gives:-

$[\bar{\gamma}(\bar{u}+d) - \bar{\gamma}(\bar{u})]/d$ increases with d and $[\bar{\gamma}(\bar{u}) - \bar{\gamma}(\bar{u}-d)]/d$ decreases with d .

$$\begin{aligned} \text{Thus } U_{\theta} \bar{\gamma}_{\theta}(\bar{u}) - \bar{\gamma}_{\theta}(\bar{u}) &= \bar{u}_1 q_1 (\bar{\gamma}_{\theta}(T_1^* \bar{u}) - \bar{\gamma}_{\theta}(\bar{u})) - \bar{u}_2 q_2 (\bar{\gamma}_{\theta}(\bar{u}) - \bar{\gamma}_{\theta}(T_2^* \bar{u})) \\ &\leq \theta^*_{\bar{u}_1 \bar{u}_2} (q_1 (\bar{\gamma}_{\theta}(T_1^* \bar{u}) - \bar{\gamma}_{\theta}(\bar{u})) / (T_1^* \bar{u} - \bar{u}) - q_2 (\bar{\gamma}_{\theta}(\bar{u}) - \bar{\gamma}_{\theta}(T_2^* \bar{u})) / (\bar{u} - T_2^* \bar{u})) \\ &\leq \theta^*_{\bar{u}_1 \bar{u}_2} (q_1 (\bar{\gamma}_{\theta}(T_1 \bar{u}) - \bar{\gamma}_{\theta}(\bar{u})) / (T_1 \bar{u} - \bar{u}) - q_2 (\bar{\gamma}_{\theta}(\bar{u}) - \bar{\gamma}_{\theta}(T_2 \bar{u})) / (\bar{u} - T_2 \bar{u})) \\ &= \theta^*_{\theta} (U_{\theta} \bar{\gamma}_{\theta} - \bar{\gamma}_{\theta}) = 0. \end{aligned}$$

//

Corollary 1.4.6.

The family of 2-action $\theta(\bar{u})$ U.L. rules is at least ϵ -optimal.

Proof.

We have $\lim_{\theta \downarrow 0} \bar{\gamma}_{\theta}(\bar{u}) = 0$ and $\lim_{n \rightarrow \infty} U_{\theta(\bar{u})}^n \bar{\gamma}_0 = \bar{\gamma}_{\theta(\bar{u})} \leq \bar{\gamma}_0$ and we can always find $\theta \geq \theta(\bar{u})$, by compactness and $\theta(\bar{u})/\bar{\gamma}_0$ at boundaries. Thus $\lim_{\theta \downarrow 0} \bar{\gamma}_{\theta}(\bar{u}) = 0$ where we suppose $\theta(\bar{u}) = \theta f(\bar{u})$ with θ the learning parameter and $f(\bar{u})$, the learning function. //

It is conjectured that if $\theta^*(\bar{u}) \leq \theta^{**}(\bar{u})$ then $\bar{\gamma}_{\theta^*} \leq \bar{\gamma}_{\theta^{**}}$, but this remains unproven. However, we do achieve monotonicity of $\bar{\gamma}_{\theta(\bar{u})}$ under quite general conditions.

Lemma 1.4.7.

- i) If $\frac{d}{d\bar{u}_i} (T_i \bar{u}) > 0$, for $i = 1$ and 2 , then $\bar{\gamma}_{\theta(\bar{u})}$ is monotone.
- ii) If $\exists \bar{u}_i$ s.t. $\frac{d}{d\bar{u}_i} (T_i \bar{u}) < 0$ for $i = 1$ or 2 , then $\exists \phi(\bar{u})$ s.t.:-
 - a) $\phi(\bar{u})$ monotone $\Rightarrow U\phi(\bar{u})$ monotone \uparrow .
 - b) $\phi(\bar{u})$ non-monotone $\Rightarrow U\phi(\bar{u})$ non-monotone.
- iii) If $\theta(\bar{u}_1) = \theta(\bar{u}_2)$ then if $\frac{d}{d\bar{u}_i} (T_i \bar{u}) > 0 \forall \bar{u}_i$, $\bar{\gamma}_{\theta(\bar{u})}$ is monotone.

Proof.

- i) We assume $\phi(\bar{u})$ is monotone increasing with $\phi(\bar{u}) = \bar{u}_i$, $i = 0, 1$.

and show that this is preserved by U and hence that $\mathcal{J}_{\theta(\bar{u})}$ is monotone.

$$\begin{aligned}
 U\phi(\bar{u}^*) - U\phi(\bar{u}) &= (\bar{u}_1^* p_1 + \bar{u}_2^* p_2) \phi(\bar{u}^*) - (\bar{u}_1 p_1 + \bar{u}_2 p_2) \phi(\bar{u}) \\
 &\quad + \bar{u}_1^* q_1 \phi(T_1 \bar{u}^*) - \bar{u}_1 q_1 \phi(T_1 \bar{u}) \\
 &\quad + \bar{u}_2^* q_2 \phi(T_2 \bar{u}^*) - \bar{u}_2 q_2 \phi(T_2 \bar{u}) \quad \text{where } \bar{u}_i^* > \bar{u}_i.
 \end{aligned}$$

Now

$$\frac{\partial}{\partial \bar{u}_i} (T_i \bar{u}) > 0 \quad i=1,2 \Rightarrow T_1 \bar{u}^* > T_1 \bar{u} \\
 T_2 \bar{u}^* > T_2 \bar{u}$$

Thus

$$\begin{aligned}
 U\phi(\bar{u}^*) - U\phi(\bar{u}) &> (\bar{u}_1^* - \bar{u}_1) (q_2 - q_1) \phi(\bar{u}) + q_1 (\bar{u}_1^* - \bar{u}_1) \phi(\bar{u}) + q_2 (\bar{u}_2^* - \bar{u}_2) \phi(T_2 \bar{u}) \\
 &= (\bar{u}_1^* - \bar{u}_1) q_2 (\phi(\bar{u}) - \phi(T_2 \bar{u})) > 0 \quad \text{by monotonicity of } \phi.
 \end{aligned}$$

Hence $\mathcal{J}_{\theta(\bar{u})}$ is monotone.

ii) I shall just prove that if $T_1 \bar{u}^* < T_1 \bar{u}$, with $\bar{u}_i^* > \bar{u}_i$ and

ϕ is monotone, then $U\phi$ is not always monotone, where $(\bar{u}_1^* - \bar{u}_1) < \epsilon$ sufficiently small such that $T_1(\bar{u}_1 - \bar{u}_1^*) > 0$ is possible with $\frac{\partial}{\partial \bar{u}_1} (T_1 \bar{u}) < 0$.

Define $0 < \phi(T_2 \bar{u}) = \phi(T_2 \bar{u}^*) = \phi(\bar{u}) = \phi(\bar{u}^*) = \phi(T_1 \bar{u}^*) < \phi(T_1 \bar{u}) = 1$ with ϕ cts.

and so $U\phi(\bar{u}^*) - U\phi(\bar{u}) < q_1 (\bar{u}_1 - \bar{u}_1^*) \phi(\bar{u}) - \bar{u}_1 q_1 + \bar{u}_1^* q_1 \phi(\bar{u}) = \bar{u}_1 q_1 (\phi(\bar{u}) - 1) < 0$.

and hence $U\phi$ is non-monotone, since $U\phi(i) = i$, $i=0,1$, holds as for ϕ .

The other assertion follows similarly, choosing $\phi(\bar{u})$ appropriately.

iii) If $\frac{\partial}{\partial \bar{u}_1} (T_1 \bar{u})|_{\bar{u}_1^*} > 0$ then $\frac{\partial}{\partial \bar{u}_1} (T_2 \bar{u})|_{\bar{u}_2^*} > 0$.

$$\begin{aligned}
 \text{since } \frac{\partial}{\partial \bar{u}_1} (\bar{u}(1 - \theta(\bar{u})) + \theta(\bar{u}))|_{\bar{u}_1} &= 1 - \theta(\bar{u}_1) + (1 - \bar{u}_1) \frac{\partial \theta}{\partial \bar{u}_1}|_{\bar{u}_1} = 1 - \theta(\bar{u}_2) - \bar{u}_2 \frac{\partial \theta}{\partial \bar{u}_1}|_{\bar{u}_2} \\
 &= \frac{\partial}{\partial \bar{u}_1} (\bar{u}(1 - \theta(\bar{u}))|_{\bar{u}_2} = \frac{\partial}{\partial \bar{u}_1} (T_2 \bar{u})|_{\bar{u}_2}.
 \end{aligned}$$

Hence centrally symmetric rules give monotone $\mathcal{J}_{\theta(\bar{u})}$ if $\frac{\partial}{\partial \bar{u}_1} (T_1 \bar{u}) > 0$. //

The condition $|\frac{\partial}{\partial \bar{u}_1} (T_1 \bar{u})| < 1$ is necessary for Norman's (1972) distance-diminishing rules, but this is only satisfied by linear rules, with $\frac{\partial}{\partial \bar{u}_1} (T_1 \bar{u}) = 1 - \theta$ in our U.L. $\theta(\bar{u})$ family. It may be that such a Lifshitz condition is required in order for the method of proving a property P carries across from ϕ to $U\phi$ and hence to \mathcal{J} , to operate.

More subtle techniques may be required to prove the more general comparison theorem.

The comparison theorem 1.4.5. says that the slower we learn, the more likely we are to follow the expected drift. This is the converse situation to gambling in which the optimum strategy is to play boldly and attempt to go against the drift as in Dubins and Savage (1965). In learning we wish $(\text{drift}/\text{diffusion}) = r(n) \uparrow \infty$, whilst in gambling we need $r \downarrow 0$, where $\phi'(n) + r(n)\phi(n) = 0$ gives the diffusion approximation to the process, with $\phi(n)$ generating absorption probabilities.

For the optimal rules in the following sections, we attempt to learn very slowly near boundaries, so that we achieve asymptotic reflection from sub-optimal boundaries. We are then absorbed only at the optimal stable boundary, where the drift is with us.

1.5. Boundary Behaviour.

We shall now partition the non-linear U.L. rules \mathcal{R} to obtain those which are actually optimal \mathcal{R}_0 , by examining the behaviour of $\theta(n)$ near to boundaries $\bar{n}_i \in \{0, 1\}$.

Definition 1.5.1.

A rule has boundary behaviour α at $\bar{n}_i = 0$ if $\theta(n) \sim O(\bar{n}_i^\alpha)$ as $\bar{n}_i \downarrow 0$.

I shall only consider rules which can be classified according to their α -dependence, thus eliminating unwanted pathological rules. The following lemmas give strong reasons to conjecture that $\alpha = 1$ is the transitional class occurring between $\alpha > 1$ optimal rules and $0 \leq \alpha < 1$ ϵ -optimal rules.

Lemma 1.5.2.

The continuous time deterministic approximation to $\theta(n)$ rules gives convergence with exponential behaviour in t iff $\alpha = 0$.

Proof.

We have $\Delta \bar{n}_i = \theta(\bar{n}) \bar{n}_i (1 - \bar{n}_i) (q_1 - q_2)$ and w.l.o.g. put $\theta(\bar{n}) = O(\bar{n}_i^\alpha (1 - \bar{n}_i)^\alpha)$.

$$= \sigma_{ii}^2 (1-\bar{u}_i)^\beta \quad \text{with} \quad \beta = \alpha + 1, \quad \sigma = \theta(q_1 - q_2) > 0.$$

Then we take deterministically

$$\frac{dx}{dt} = \sigma x^\beta (1-x)^\beta \quad \text{a)}$$

Now when $x \uparrow 1$ $x \sim \gamma(1-x)^\beta$ or $\int \frac{dx}{(1-x)^\beta} = \int \gamma dt = \gamma t + C$

Hence for $\beta > 1$ $\frac{1}{(1-x)^\beta} = \gamma t + C$ and $x = 1 - (\gamma t + C)^{-\frac{1}{\beta-1}}$ b)

Whilst for $\beta = 1$ $\log(1-x) = -(\gamma t + C)$ or $(1-x) = e^{-\gamma t}, x = 1 - e^{-\gamma t}$ c)

Now redefine time-scale:- $T = \gamma t + C$ and put $y = 1 - x$

$$\text{Thus } \left. \begin{aligned} x &= 1 - T^{-\frac{1}{\beta-1}}, \alpha > 0 \\ x &= 1 - e^{-T}, \alpha = 0 \end{aligned} \right\} \begin{aligned} y T^{\frac{1}{\beta-1}} &= 1 \\ y &= e^{-T} \end{aligned} \quad \begin{aligned} &\text{power behaviour} \\ &\text{exponential.} \end{aligned}$$

and in case $\alpha = 1$, we obtain $y T = 1$.

//

The linear case is special due to $\theta = \text{const}$ being distance-diminishing, whilst for $\alpha > 0$ $\lim_{\bar{u} \rightarrow 0,1} \theta(\bar{u}) = 0$ so such a Lifshitz condition cannot be imposed.

Lemma 1.5.3.

The continuous time stochastic diffusion approximation is optimal iff $\alpha > 1$, and is conditionally optimal for $\alpha = 1$.

Proof.

We solve $\phi'' + 2a(\bar{u})/b(\bar{u}) \phi' = 0$ where $a(\bar{u}) = \text{expected drift}$ and $b(\bar{u}) = \text{variance}$.

$$\begin{aligned} \text{Then } a(\bar{u}) &= \theta(\bar{u}) \bar{u}(1-\bar{u})(q_1 - q_2) \\ b(\bar{u}) &= q_1 \theta^2(\bar{u}) \bar{u}(1-\bar{u})^2 + q_2 \theta^2(\bar{u}) (1-\bar{u}) = (q_1(1-\bar{u}) + q_2 \bar{u}) \theta^2(\bar{u}) \bar{u}(1-\bar{u}). \end{aligned}$$

and for $|q_1 - q_2|$ small put $q_1 + \bar{u}(q_2 - q_1) = q_1 + o(1)$ and this will not affect the conclusion of optimality for $\alpha > 1$.

$$\text{Thus } r(\bar{u}) = \frac{2a(\bar{u})}{b(\bar{u})} = \frac{\text{drift}}{\text{diffusion}} \sim \frac{2(q_1 - q_2)}{q_1 \theta(\bar{u})} = \frac{2(1 - q_2/q_1)}{\theta(\bar{u})}$$

and for $q_1 > q_2$ (so that $\bar{u}_1 = 1$ is optimum) put $(1 - q_2/q_1) = \theta k$

where $\theta(\bar{u}) = \theta f(\bar{u})$ Thus $\frac{\phi''}{\phi'} = -\frac{2k}{f(\bar{u})}$ and for $\alpha = 1$ we put w.l.o.g. $f(\bar{u}) = \bar{u}(1-\bar{u})$.

and $\log \phi' = \int -2k \frac{1}{u(1-u)} = -2k \log \left(\frac{u}{1-u} \right) + C = \log \left(\frac{1-u}{u} \right)^{2k} + C$
 $\phi' = A \left(\frac{1-u}{u} \right)^{2k}$ and $\phi(i) = i \quad i=0,1$.

Then for $\alpha=1$ at both boundaries $\phi(x) = \frac{\int_0^x \left(\frac{1-u}{u} \right)^{2k} du}{\int_0^1 \left(\frac{1-u}{u} \right)^{2k} du}$.

Clearly when $2k \geq 1$, $\phi(x) \equiv 1$ for $x \neq 0$ and diffusion is optimal since $\int_0^x \frac{1}{u^{2k}} du$ is divergent iff $2k \geq 1$.

Whilst $2k < 1$ gives $\phi(x) < 1$ for $x < 1$ and ϕ is monotone \uparrow ; we have ϵ -optimality.

(Note: formal bounds for $\alpha=1$ will be proved in 1.5.17.)

Now since $\int_0^x \frac{1}{u^{2k}} du$ is divergent for $\alpha > 1$, we can easily verify $\phi(x) \equiv 1$ for $x \neq 0$ here also. Similarly $\int_0^x \frac{1}{u^{2k}} du$ is convergent for $\alpha < 1$, giving ϵ -optimality.

Where $\phi(x) = \left[\int_0^x \exp \int_{y_0}^y -2k \frac{1}{f(u)} du dy \right] / \left[\int_0^1 \exp \int_{y_0}^y -2k \frac{1}{f(u)} du dy \right]$ is easily found.

Finally, we find this diffusion limit for $\alpha=0$.

$$\phi'/\phi = -2k \quad \log \phi' = -2kx + C \quad \phi' = Ae^{-2kx}, \quad \phi = Be^{-2kx} + C.$$

And boundary conditions give $\phi(x) = (1 - e^{-2kx}) / (1 - e^{-2k})$ //

Norman (1971) proves that the discrete time $\phi_0(n)$ converges weakly to $\phi(x)$ for this linear rule. The theorem requires $r(u)$ bounded and hence breaks down for $\alpha > 0$, when $r(u) \uparrow \infty$ at boundaries. Vorontsova (1965) proved a result similar to 1.5.3. for a family of rules with penalty and reward reinforcement. However, for $\alpha < 1$, such rules cannot be normalised to give uniform learning and thus our reward-inaction scheme appears the more fundamental. I shall now state a result of Norman (1968) which gives bounds on the discrete time absorption probabilities for the linear rule.

Lemma 1.5.4.

- i) Let $\phi_{z,0}(\bar{u}) = (1 - e^{-z\bar{u}/\theta}) / (1 - e^{-z/\theta})$ then \exists +ve y and z
 such that $\phi_{y,0} \leq \gamma_{\theta}(\bar{u}) \leq \phi_{z,0} \quad \forall \theta, \bar{u} \in [0, 1]$
 ii) $\gamma_{\theta}(\bar{u}) \geq (1 - e^{-z\bar{u}/\theta})$ where $q_1 > q_2$.

Proof.

Norman (1968), proves $\gamma_{x,0}(\bar{u}) = e^{x\bar{u}/\theta}$ can be super or sub-regular depending on x , and the result i) then follows easily by noting that the class of super and sub-regular functions are closed under addition and multiplication by +ve constants.

For ii) we have $(1 - e^{-z/\theta}) < 1$ and hence also $\lim_{\theta \rightarrow 0} \gamma_{\theta}(\bar{u}) = 1$ giving ϵ -optimality.

//

In the literature, such as Sawaragi (1974), we find that the same method is applied to rules of classes $\alpha > 0$. However, we may observe that Norman's $\phi_{z,0}(\bar{u})$ arises precisely from the weak convergence limit of 1.5.3, which only has exponential behaviour for $\alpha = 0$, as in the deterministic solution 1.5.2. For $\alpha > 0$, we should only expect to be able to put tight bounds on $\gamma_{\theta, \alpha}(\bar{u})$ by using the weak convergence limit $\phi(\bar{x})$ appropriate to it. The use of the $\alpha = 0$ limit for $\alpha > 0$ rules will only give us ϵ -optimality, which has already been proved for all $\alpha > 0$ in theorem 1.4.5.

Lemma 1.5.5.

The $\lim_{n \rightarrow \infty} \Pr(\text{attain } \bar{u} > \frac{1}{2}n \text{ before absorption at } \bar{u}_1 = 0, \bar{u}_2(0) = \frac{1}{2}n, \text{ using } \theta(\bar{u}) = 1)$
 when $\alpha > 1$ and $q_1 > q_2$.

Or, concisely we write, $\lim_{n \rightarrow \infty} P_{r_{\theta}(\bar{u})}(\frac{1}{2}n \text{ BO } | \frac{1}{2}n) = 1$.

Proof.

$$P_{r_{\theta}(\bar{u})}(\frac{1}{2}n \text{ BO } | \frac{1}{2}n) \geq P_{\theta'(\bar{u})}(\frac{1}{2}n \text{ BO } | \frac{1}{2}n)$$

with

$$\begin{aligned} \theta'(\bar{u}) &= \theta(\bar{u}) \text{ on } \bar{u} \in [0, \frac{1}{2}n) \\ \theta'(\bar{u}) &= \theta(\frac{1}{2}n) \text{ on } \bar{u} \in [\frac{1}{2}n, 1] \end{aligned}$$

But $\gamma_{\theta, \frac{1}{2}n}(\bar{u}) \leq \gamma_{\theta'(\bar{u})}(\bar{u})$ by 1.4.5. as $\theta'(\bar{u}) \leq \theta(\frac{1}{2}n)$

will hold w.l.o.g. (If $\theta(\bar{u})$ oscillates near $\bar{u} = \frac{1}{2}n$ we just bound

$\theta(\bar{u})$ suitably above by some $\theta = \text{const}$)

Then by 1.5.4. $\delta_{\theta(\bar{u})}(\bar{u}) \geq 1 - e^{-z^{1/\theta(\bar{u})}}$ for +ve $z \quad \forall \bar{u}, n.$

Thus $P_{\theta(\bar{u})}(\frac{1}{n} B O | \frac{1}{2} n) \geq 1 - e^{-z^{1/\theta(\bar{u})}}$ now put $\bar{u} = \frac{1}{2} n$

and $\theta(\frac{1}{n}) = \theta(\frac{1}{n})^\alpha (1 - \frac{1}{n})^\alpha < \theta/n^{\alpha-1}$ and hence,

$$P_{\theta(\bar{u})}(\frac{1}{n} B O | \frac{1}{2} n) \geq 1 - e^{-z^{1/\theta(\bar{u})}} = 1 - \exp\left(-\frac{z n^{\alpha-1}}{2^\alpha}\right)$$

and $\lim_{n \rightarrow \infty} \exp\left(-\frac{z}{2^\alpha} n^{\alpha-1}\right) = 0$ when $\alpha > 1$, and hence result. //

This lemma shows why we may expect asymptotic reflection from sub-optimal boundaries. For $\alpha < 1$, using the diffusion limit and weak convergence for $\alpha = 0$ bounds, we can similarly show $P_{\theta(\bar{u})}(\frac{1}{n} B O | \frac{1}{2} n) = 0$ giving absorption at sub-optimal boundaries. However, for $\alpha = 1$, we can obtain no such bounds, for non-rigorously, we can put $\alpha = 1$ in $\delta_{\theta(\bar{u})}(\bar{u})$ to obtain only $P_{\theta(\bar{u})}(\frac{1}{n} B O | \frac{1}{2} n) \geq 1 - e^{-z/2^\alpha} = \text{const}$, which demonstrates how difficult this transition case really is to understand.

Definition. 1.5.6.

If at time t we choose action i and at time $t+1$ we choose action $j \neq i$, then we have an alternation.

The following lemma shows how the alternation concept is closely linked to that of optimality. However, only if we exclude the transition case $\alpha = 1$, is # alternations = nec and suff for optimality. For the transition rules, we shall prove this to be neither nec nor suff for convergence to optimum action, since 1.5.17. will show $\alpha = 1$ to be conditionally optimal as in the diffusion.

For the next lemma we consider boundary behaviour as $\bar{u} \uparrow$.

We let $\theta(\bar{u}) \sim O(1 - \bar{u})^\alpha$ as $\bar{u} \uparrow$ which determines absorption behaviour whilst 1.5.1 effectively defines reflection behaviour. However, after 1.5.7. all boundaries will be assumed to have the same α -absorption and α -reflection behaviour, so the distinction is only of mathematical interest.

Lemma 1.5.7.

At boundary 1, $\theta(\bar{u}_i) \sim O((1-\bar{u}_i)^\alpha)$ as $\bar{u}_i \uparrow 1$.

i) If $\alpha < 1$ and $\lim_{t \rightarrow \infty} \bar{u}_i(t) = 1$ then $\#abts \stackrel{q_i}{\sim} \infty$ and if $\alpha < 1$ at all boundaries, this is strengthened to $\#abts \stackrel{q_i}{\sim} \infty$. (And indeed all moments are finite.)

ii) If $\alpha > 1$ and $\lim_{t \rightarrow \infty} \bar{u}_i(t) = 1$ then $\#abts \stackrel{q_i}{\sim} \infty$.

iii) If $\alpha = 1$ and $\lim_{t \rightarrow \infty} \bar{u}_i(t) = 1$ then $\#abts \stackrel{q_i}{\sim} \infty$ iff $\theta q_i \leq 1$, where q_i = reward probability, as usual.

Proof.

We define $i_t(\bar{u}_j) = \text{Prob}(\text{We always choose action } j \mid \bar{u}_j(t) = \bar{u}_j)$.

Now we immediately have the basic recurrence equation,

$$i_t(\bar{u}_j) = \bar{u}_j (p_j i_t(\bar{u}_j) + q_j i_t(\bar{u}_j + \theta(\bar{u})/(1-\bar{u}_j))) \quad 1)$$

and hence $i_t(\bar{u}_j) = \left((1 - (1-\bar{u}_j)/(1-p_j\bar{u}_j)) \right) i_t(\bar{u}_j + \theta(\bar{u})/(1-\bar{u}_j))$ with $p_j = 1 - q_j$

I shall now omit suffixes j and t and write $e(\bar{u}) = (1-\bar{u})/(1-p\bar{u})$ and $T\bar{u} = \bar{u} + \theta(\bar{u})/(1-\bar{u})$.

Thus $i(\bar{u}) = (1 - e(\bar{u})) i(T\bar{u}) = \prod_{n=0}^{\infty} (1 - e(T^n \bar{u}))$.

where $\lim_{n \rightarrow \infty} i(T^n \bar{u}) = 1$ since $\theta(\bar{u})/(1-\bar{u}) = 0$ iff $\bar{u} = 0, 1$.

Now $(1-\bar{u})/(1-p\bar{u}) \sim 1/(1-p) \left((1-\bar{u}) - p/(1-\bar{u}) + \dots \right)$ and hence as

$\bar{u}_i \uparrow 1$ we can neglect $O(1-\bar{u})$ terms, except in the case $\alpha = 1$ when we must be much more careful.

Note that $1 - e(\bar{u}) = q\bar{u}/(1-p\bar{u})$, which is a transformed distribution function.

Now for $\alpha \neq 1$, we have $i(\bar{u}) > 0$ iff $\sum_{n=0}^{\infty} (1 - \bar{u}_n) < \infty$

We prove convergence or divergence for α by a comparison test with $\sum 1/n^\beta$, β chosen appropriately, and ii) follows immediately from i). We use the condition $i(\bar{u}) > 0$ iff $\sum_{n=0}^{\infty} e(T^n \bar{u}) < \infty$.

i) a) Put $\bar{u}_n = 1 - 1/n^\beta$, for some fixed $\beta > 1$ and so $\theta(\bar{u}) \sim \theta(1-\bar{u})^\alpha$

Define $b_n = 1 - \bar{u}_n = 1/n^\beta$ and so $b_{n+1} = 1 - \bar{u}_{n+1} = 1/n^\beta - \theta/n^{\alpha\beta}$
 $b_{n+1} = 1/n^\beta (1 - \theta/n^{\alpha\beta})$.

We take the ratio $b_n/b_{n+1} \sim 1 + \theta/n^\alpha \beta$ whilst $a_n/a_{n+1} = \frac{1/n^\beta}{1/(n+1)^\beta} \sim 1 + \beta/n$

But $\sum 1/n^\beta < \infty$ for $\beta > 1$ and if $\alpha\beta < 1$, $b_n/b_{n+1} > a_n/a_{n+1}$ for some $n > n_0$.

Hence, since the above holds for any $n \gg 1$ and we can always find

$\beta > 1$ s.t. $\alpha\beta < 1$ when $\alpha < 1$, we have (by Gauss Ratio Test,)

$$\alpha < 1 \Rightarrow \sum_{n=1}^{\infty} (1 - \bar{u}_n) < \infty \text{ and hence } i(\bar{u}) > 0$$

Now since the process is markovian and $i(\bar{u}) \uparrow$ as $\bar{u} \uparrow$, we

must a.s. have t_0 s.t. for $t > t_0$, we take only action 1 and $\#alts \stackrel{a.s.}{\rightarrow} \infty$.

b) Now we have to prove $E\#alts < \infty$ for $\alpha < 1$ at all boundaries, and indeed $E(\#alts)^n < \infty$.

Prob(at least 1 more alternation starting from next choice) $< 1 - \epsilon$

where $\min_n \sum_{i=1}^n i(\bar{u}_i) > \epsilon > 0$ which exists by continuity of $i(\bar{u})$ and $\lim_{\bar{u} \uparrow 1} i(\bar{u}) = 1$.

Now take successive alternations, starting process again at

\bar{u}_{n+1} if r^{th} alternation occurs at \bar{u}_{n-1}, \bar{u}_n .

Then $\Pr(\#alts > 2r) < (1 - \epsilon)^r / \epsilon$ (allowing possibility of alt at $n, n+1$ trials)

Clearly also $\Pr(\#alts > 2r \text{ and } \#alts < 4r) < (1 - \epsilon)^r / \epsilon$.

Then $E\#alts = \sum k p_k < 3 + \sum_{s=0}^{\infty} 2^{s+2} (1 - \epsilon)^{s/2} / \epsilon < \infty$ by comparison with G.P.

since $\Pr(2^s < \#alts < 2^{s+1}) < (1 - \epsilon)^{2^{s-1}} / \epsilon$.

Also given $\delta > 0, \exists n_0$ s.t. $\Pr(\#alts > n_0) < \delta$, follows by above, and hence

we prove rigorously the last statement in a).

Also $E(\#alts)^n = \sum k^n p_k < C + \sum_{s=0}^{\infty} 2^{n(s+2)} (1 - \epsilon)^{2^{s-1}} / \epsilon < \infty$ by comparison with G.P. where $C < 2^n + 1$.

ii) For $\alpha > 1$, we find $\beta: \alpha\beta > 1$ with $\sum 1/n^\beta < \infty$, which is always possible, and hence $\sum_{n=1}^{\infty} (1 - \bar{u}_n) = \infty$ and $i(\bar{u}) = 0 \Rightarrow \#alts \stackrel{a.s.}{\rightarrow} \infty$.

iii) We put $(1 - \bar{u}_i) = 1/x$ giving $(1 - \bar{u}_i) / (1 - p\bar{u}_i) = 1/(p+x) = 1/n$.

Note that in this lemma \bar{u}_t refers to action i at trial t , since we are considering boundary i throughout.

Define $a_t = (1 - \bar{u}_t) / (1 - p\bar{u}_t)$ Then $i(\bar{u}) = 0$ iff $\sum_t a_t = \infty$.

We prove divergence using the ratio test; comparison with $\sum 1/n$.

We just compare the tail of a_t with $b_n = 1/n$ and by integral or ratio

test; tail above $y = 1/n \Rightarrow$ divergent

and tail strictly below $y = 1/n \Rightarrow$ convergent

$$\begin{aligned} \text{Now } (1 - \bar{u}_{t+1}) / (1 - p\bar{u}_{t+1}) &= (1 - \bar{u}_t) (1 - \theta(\bar{u}_t)) / (1 - p(\bar{u}_t + (1 - \bar{u}_t) \theta(\bar{u}_t))) \\ &= (1 - \theta(\bar{u}_t)) / (p + (1 - p) - p\theta(\bar{u}_t)) \\ &= (1 - \theta(\bar{u}_t)) / (1 - p\theta(\bar{u}_t)) \end{aligned}$$

$$\text{Thus } a_{t+1}/a_t = (1 - 1/n p\theta(\bar{u}_t)) / (1 - \theta(\bar{u}_t))$$

But we have $\theta(\bar{u}_t) = 0(1 - \bar{u}_t)$ as $\bar{u}_t \uparrow 1$ and so let $\theta(\bar{u}_t) = \theta(1 - \bar{u}_t)$.

Then $\theta(\bar{u}_t) = \theta q / n(1 - p q)$ and $a_{t+1}/a_t = 1 + \theta q / n + o(1/n^2)$ which gives divergence when $\theta q < 1$

Thus $i(\bar{u}_t) > 0$ iff $\theta q > 1$ and iff $\#alt < \infty$.

//

Remarks 1.5.8

a) For the linear rule it is actually possible to express alternations in terms of $\delta(\bar{u})$ as in Norman (1968). We shall give an easier proof here.

Lemma 1.5.9.

For the linear rule $\chi(\bar{u}) = (2 - \theta(q_1, q_2))(\bar{u} - \delta(\bar{u})) / \theta(q_2 - q_1)$.

and $\chi(\bar{u}) = 2(1 - \theta q) \bar{u} (1 - \bar{u}) / \theta^2 q$ when $q_1 = q_2$

$\chi(\bar{u}) = \#$ alternations in response, $\pi(\bar{u}) = \bar{u}$.

Proof.

$\Delta \chi(\bar{u}) = -\bar{u}_1 \bar{u}_2 (2 - \theta(q_1, q_2))$ is easily found.

and $\Delta \bar{u} = \theta(q_1 - q_2) \bar{u}_1 \bar{u}_2$ enables us to substitute when $q_1 \neq q_2$.

and we have $\Delta\theta(\bar{n})=0$ and $\theta(\bar{n})$ is the unique fixed function satisfying the given boundary conditions.

So $\chi(\bar{n}) = (2-\theta(q_1+q_2))/\theta(q_2-q_1) (\bar{n}-c\theta(\bar{n}))$ where $c = 1$ from boundary values.

If $q_1 = q_2$ then we find $\Delta\bar{n}_1, \bar{n}_2 = -\theta^2 q_1 \bar{n}_2$
and $\Delta\chi = -2\bar{n}_1 \bar{n}_2 (1-\theta q)$

So again we just use uniqueness of $\chi(\bar{n})$ for result.

//

b) This method only works for $\theta = \text{const}$, else we get non-linearities.

However, if $\bar{n} \sim \frac{1}{n}$, $\theta(\bar{n}) \sim \theta(\frac{1}{n})$ then $\chi \sim n^{(\alpha-1)}$ on putting $\theta(\bar{n}) \sim 0$, $q_2 > q_1$.

Thus as $n \rightarrow \infty$ $\chi = \infty$, $\alpha > 1$ whilst at the transition case $\alpha = 1$,
 $\chi < \infty$ $\alpha < 1$

we get $\chi \sim 2/\theta(q_2-q_1)$ which gives no information. This is a non-rigorous argument and 1.5.7. is still required to give tight conditions on $\theta(\bar{n})$ in a fully rigorous manner, and to deal with the subtle case of $\alpha = 1$.

Lemma 1.5.10.

$i(\bar{n}_1) \uparrow$ when $d_{\theta\bar{n}}(T, \bar{n}) > 0 \quad \forall \bar{n} \in [0, 1]$.

Proof.

We have $i(\bar{n}_1) = (1 - \bar{n}_2/(1-p\bar{n}_1)) i(T, \bar{n})$

Thus if $\bar{n}'_1 > \bar{n}_1$, $i(\bar{n}'_1) - i(\bar{n}_1) \geq (i(T, \bar{n}'_1) - i(T, \bar{n}_1)) (1 - \bar{n}_2/(1-p\bar{n}_1))$.

and $d_{\theta\bar{n}}(T, \bar{n}) > 0 \Rightarrow T, \bar{n}'_1 > T, \bar{n}$.

Now iterate to get $i(\bar{n}'_1) - i(\bar{n}_1) \geq (i(T, \bar{n}'_1) - i(T, \bar{n}_1)) \prod_{r=0}^{n-1} (1 - \frac{\bar{n}_2}{1-p\bar{n}_r}) > 0$

Hence $i(\bar{n}_1) \uparrow$

//

In future sections, we shall consider rules with the same α -behaviour at all boundaries, so that the two definitions we have used become equivalent.

For the transition case $\alpha = 1$, we see that χ is bounded.

For the transition case iii) 1.5.7., we see that it is possible for $\theta q_1 < 1 < \theta q_2$ and yet in 1.5.17. we shall prove that the rule may still be either optimal or ϵ -optimal. Thus we may have an ϵ -opt rule with $\#alt = \infty$ with finite probability, or an opt rule with $\#alt = \infty$. If $\theta q_1 < \theta q_2 < 1$, then we shall also find ϵ -opt rules with $\#alt = \infty$.

Corollary 1.5.11.

If $i(i_i) > 0$ and $\exists j$ s.t. $q_j > q_i$ then the rule is at best ϵ -optimal, and we say it is in class R_ϵ .

Proof.

By comparison theorem 1.4.5, all U.L. $\theta(i)$ rules are at worst ϵ -opt, but $i(i_i) > 0$ with finite probability we only take sub-optimal action i. Hence all rules with $\alpha < 1$, at any boundary, (since we do not know which is optimum) are in R_ϵ , and also rules with $\alpha = 1, q_i \theta > 1$ at boundary i. //

We must now consider those rules for which $i(i_i) \equiv 0$ at every sub-optimal boundary; $q_{opt} > q_i \forall i \neq opt$. First we prove that the class of optimal rules R , includes those with $\alpha > 1$ at every boundary. Then for $\alpha = 1$, we delineate regions which give conditional optimality depending on θ and $|q_i - q_{opt}|$, and other regions which are at best, ϵ -optimal.

Although the class of rules $\alpha = 1$, with its conditional optimality (1.5.17) would not be of great practical use, its study does help us to understand the behaviour of $\alpha \neq 1$ more fully, as ϵ -opt "fades" into optimality. Both have fundamentally different workings as discussed in section 7.

Theorem 1.5.12.

- i) The class of rules with $\alpha > 1$ at each boundary is optimal.
- ii) This optimality is independent of the behaviour of $\theta(i)$ in finitely many compact subsets of $I = [0, 1]$

Proof.

We construct a sub-regular family $\mathcal{S}_{k\delta}(x)$, guided in choice by 1.5.3.

Define $\mathcal{S}_{k\delta}(x) = \frac{\gamma_{k\delta}(x)}{\gamma_{k\delta}(1)}$ where $\gamma_{k\delta}(x) = \int_0^\delta f_k(t) dt + \int_\delta^x f_k(t) dt \quad x \geq \delta$
 $= \int_0^x f_k(t) dt \quad 0 \leq x < \delta.$

and where

$$\begin{aligned} f_k(t) &= \exp \left[\frac{k}{[t(1-t)]^{\alpha-1}} \right] & t \in [0, \frac{1}{4}] \\ f_k(t) &= \exp \left[\frac{2k(1-2t)}{[t(1-t)]^{\alpha-1}} \right] & t \in [\frac{1}{4}, \frac{3}{4}] \\ f_k(t) &= \exp \left[-\frac{k}{[t(1-t)]^{\alpha-1}} \right] & t \in [\frac{3}{4}, 1] \end{aligned}$$

Our $f_k(t)$ is chosen, since at the boundaries $t=0,1$, it approximates the integrand in $\phi(x) = \int_0^x \exp \int_{y_0}^y -2k\theta_{\theta(t)} dt dy \sim \int_0^x f_k(t) dt.$

We truncate the integral at $x=\delta$, to remove the divergence, and

as $\delta \downarrow 0$, we obtain $\mathcal{S}_{k\delta}(x) \uparrow \chi_{\theta(t)}(t) \equiv 1$ on $t \neq 0$. Our task is to prove, for $\alpha > 1$, $U_{\theta(t)} \mathcal{S}_{k\delta}(t) - \mathcal{S}_{k\delta}(t) > 0, t \in [0,1]$, and hence $\chi_{\theta(t)}(t) \geq \mathcal{S}_{k\delta}(t) \quad \forall t \in [0,1]$.
 (for some $k > k^* > 0$ for every δ).

Lemma 1.5.13

$\gamma(x) = \int_0^x dt$ is sub-regular for all $\theta(t)$, when $q_1 \geq q_2$.

Proof.

$$U\gamma - \gamma = q_1 \bar{t}_1 \int_{\bar{t}_1}^{T\bar{t}_1} dt - q_2 \bar{t}_2 \int_{T\bar{t}_2}^{\bar{t}_2} dt = \theta(t) \bar{t}_1 \bar{t}_2 (q_1 - q_2) \geq 0.$$

This holds since $\chi_{\theta(t)}(t) = \bar{t}$ is the solution to $U_{\theta(t)} \chi_{\theta(t)} = \chi_{\theta(t)}$ when $q_1 = q_2$.

Lemma 1.5.14

$$\exists k > 0 \text{ s.t. } \mathcal{V}_k(t) = q_1 \bar{t}_1 \int_{\bar{t}_1}^{T\bar{t}_1} f_k(t) dt - q_2 \bar{t}_2 \int_{T\bar{t}_2}^{\bar{t}_2} f_k(t) dt > 0$$

Proof.

Here we formally prove $\mathcal{S}_{k\delta}(x)$ is sub-regular, and most of the theorem arises from this lemma. We apply the δ -truncation later so that we may write $U_{\theta(t)} \mathcal{S}_{k\delta}(t) - \mathcal{S}_{k\delta}(t) = q_1 \bar{t}_1 \int_{\bar{t}_1}^{T\bar{t}_1} f_k(t) dt - q_2 \bar{t}_2 \int_{T\bar{t}_2}^{\bar{t}_2} f_k(t) dt$ when $T\bar{t}_1 > \delta$. We partition $I = [0,1]$ into $I_1 = [0, \epsilon]$, $I_2 = [\epsilon, 1-\epsilon^*]$ and $I_3 = [1-\epsilon^*, 1]$.

i) Take $I_2 = [\epsilon, 1-\epsilon^*]$ a compact subset of $[0,1]$.

Now in $\lim_{k \downarrow 0} f_k(\bar{n}) \rightarrow 1$ and $\lim_{k \downarrow 0} q_1 \bar{n}_1 \int_{\bar{n}_1}^{\bar{n}_2} f_k d\bar{n}_1 - q_2 \bar{n}_2 \int_{\bar{n}_1}^{\bar{n}_2} f_k d\bar{n}_1 = q_1 - q_2 > 0$

By compactness, we can use a finite open covering from any open covering. So choose intervals S_i s.t. $\exists k_i$ and $\mathcal{D}_{k_i}(\bar{n})$ on S_i ; then take a finite covering of I_2 and put $k^* = \min_i k_i > 0$ for i in our open cover.

ii) Now take $I_1 = [0, \epsilon]$ with $\epsilon < \frac{1}{4}$.

$$\begin{aligned} \mathcal{D}_k(\bar{n}) &> q_1 \bar{n}_1 (\bar{n}_1 - \bar{n}_2) f_k(\bar{n}_1) - q_2 \bar{n}_2 (\bar{n}_1 - \bar{n}_2) f_k(\bar{n}_2) \\ &= (q_1 \exp[k/(\bar{n}_1 - \bar{n}_2)^{\alpha-1}] - q_2 \exp[k/(\bar{n}_2 - \bar{n}_1)^{\alpha-1}]) \theta(\bar{n}) \bar{n}_1 \bar{n}_2 \\ &= \theta(\bar{n}) \bar{n}_1 \bar{n}_2 (q_1 \exp[k/\bar{n}_1^{\alpha-1} \bar{n}_2^{\alpha-1} (1 - \theta \bar{n}_1^{\alpha-1} \bar{n}_2^{\alpha-1})^{\alpha-1} (1 + \theta \bar{n}_1^{\alpha-1} \bar{n}_2^{\alpha-1})^{\alpha-1}] \\ &\quad - q_2 \exp[k/\bar{n}_1^{\alpha-1} \bar{n}_2^{\alpha-1} (1 + \theta \bar{n}_1^{\alpha-1} \bar{n}_2^{\alpha-1})^{\alpha-1} (1 - \theta \bar{n}_1^{\alpha-1} \bar{n}_2^{\alpha-1})^{\alpha-1}]). \end{aligned}$$

Thus $\mathcal{D}_k(\bar{n}) > 0$ iff $q_1/q_2 > \exp[k/\bar{n}_1^{\alpha-1} \bar{n}_2^{\alpha-1} (1 - \theta \bar{n}_1^{\alpha-1} \bar{n}_2^{\alpha-1})] \cdot B(\bar{n})$.

where $B(\bar{n}) = (1 + \theta \bar{n}_1^{\alpha-1} \bar{n}_2^{\alpha-1})^{1-\alpha} - (1 + \theta \bar{n}_1^{\alpha-1} \bar{n}_2^{\alpha-1})^{1-\alpha}$.

Then put $x = (\bar{n}_1 \bar{n}_2)^{\alpha-1}$ to obtain $B(\bar{n}) = (1 + \theta \bar{n}_1^{\alpha-1} x)^{1-\alpha} - (1 + \theta \bar{n}_2^{\alpha-1} x)^{1-\alpha}$.

$$\begin{aligned} B(\bar{n}) &= \theta x (\bar{n}_2 - \bar{n}_1) (\alpha-1) - \theta^2 x^2 \alpha (\alpha-1) (\bar{n}_2^{\alpha-1} - \bar{n}_1^{\alpha-1}) + \dots \\ &> \theta x (\bar{n}_2 - \bar{n}_1) (\alpha-1) - \theta^2 x^2 \alpha (\alpha-1) (\bar{n}_2^{\alpha-1})_{/2}. \end{aligned}$$

by some easy analysis, where $\theta(\bar{n}) = \theta \bar{n}_1^{\alpha-1} \bar{n}_2^{\alpha-1}$ w.l.o.g.

Thus $\mathcal{D}_k(\bar{n}) > 0$ if $k \theta (\alpha-1) < (1 - \theta \bar{n}_1^{\alpha-1} \bar{n}_2^{\alpha-1}) / ((\bar{n}_2 - \bar{n}_1) - \theta x \bar{n}_2^{\alpha-1}) \cdot \log(q_1/q_2)$

or $k < g(\bar{n}) / \theta (\alpha-1) \log(q_1/q_2)$ for $\bar{n} \in [0, \epsilon]$

But r.h.s. > 0 when $q_1 > q_2, \alpha > 1$ and is strictly bounded away from zero.

Also in $\lim_{\bar{n}_1 \downarrow 0} k < \frac{1}{\theta (\alpha-1)} \cdot \log(q_1/q_2)$ is resulting constraint.

Hence we can certainly take $k^{**} = \min_{\bar{n}_1} g(\bar{n}_1) / \theta (\alpha-1) \log(q_1/q_2)$ for ϵ sufficiently small.

In fact if we put $\theta(\bar{u}) = \theta \bar{u}_1^{\alpha} \bar{u}_2^{1-\alpha}$, $\alpha > \alpha$ then we have the far weaker constraint $q_1/q_2 > \exp[k \bar{u}_1^{\alpha} \bar{u}_2^{1-\alpha} \theta(\alpha-1) g(\bar{u})]$ on k and so we can find $k > 0$ s.t. $P_k(\bar{u}) > 0$ on I_1 for $\alpha^* > \alpha$.

iii) Now we do the same calculations for I_3 and by the symmetry of \bar{u}_1, \bar{u}_2 we find $k' = \min_{\bar{u}_2} g(\bar{u}_2) / \theta(\alpha-1) \log(q_1/q_2) > 0$ for ϵ^* sufficiently small.

iv) Now take $k^s = \min(k', k^*, k^{**})$ to give $P_{k^s}(\bar{u}) > 0$ on I . //

Remarks. 1.5.15.

a) We only required the learning function for I_1 and I_3 ; the compactness argument is independent of $\theta(\bar{u})$.

b) I shall now show semi-rigorously why we achieve $k > 0$ in ii) and iii). Let $\bar{u}_1 \sim 0$ so $\theta(\bar{u}) \sim \theta \bar{u}_1^{\alpha}$

$$P_k(\bar{u}) > 0 \quad \text{iff} \quad q_1 \exp\left[\frac{k}{(T_1 \bar{u}_1)^{\alpha-1}}\right] - q_2 \exp\left[\frac{k}{(T_2 \bar{u}_1)^{\alpha-1}}\right] > 0$$

$$\text{or} \quad \exp k \left[\frac{(T_1 \bar{u}_1)^{\alpha-1} - (T_2 \bar{u}_1)^{\alpha-1}}{(T_1 \bar{u}_1)^{\alpha-1} (T_2 \bar{u}_1)^{\alpha-1}} \right] < q_1/q_2.$$

and as $(T_1 \bar{u}_1)^{\alpha} > (T_2 \bar{u}_1)^{\alpha}$ we have $0 < k < g(\bar{u}_1) \log(q_1/q_2)$
 e.g. $\alpha=2$ then $(T_1 - T_2)\bar{u}_1 = \theta(\bar{u}) = \theta \bar{u}_1^2 \bar{u}_2 \sim \theta \bar{u}_1^2$ and $(T_1 \bar{u}_1)(T_2 \bar{u}_1) \sim \bar{u}_1^2$.
 Then $\exp(k\theta) < q_1/q_2 \Rightarrow k < \frac{1}{\theta} \log q_1/q_2$.

c) The function $P_k(\bar{u}) > 0$ for $\alpha^* < \alpha$.

Lemma 1.5.16.

The family $\rho_{k\delta}(\bar{u})$ are sub-regular for k, δ , sufficiently small.

Given $\theta(\bar{u}) \exists k^*$ s.t. for $0 < k < k^*$ and $0 < \delta < \delta(k^*)$ we have $\forall \theta(\bar{u}) \rho_{k\delta}(\bar{u}) > \rho_{k\delta}(\bar{u})$

Proof.

We use 1.5.14. and note $\rho_{k\delta}(\bar{u})$ is sub-regular for $T_1 \bar{u}_1 < \delta$ by 1.5.13. We must be careful of the transition at $T_1 \bar{u}_1 = \delta$.

a) $T_1 \bar{u}_1 > \delta, \bar{u}_1 < \delta$ b) $\bar{u}_1 > \delta, T_2 \bar{u}_1 > \delta$.

Then $\int_{\bar{u}_1}^{T_1 \bar{u}_1} f_k(\bar{u}_1) d\bar{u}_1 > (T_1 \bar{u}_1 - \bar{u}_1) f_k(T_1 \bar{u}_1)$.

$$\text{and } \int_{T_2 \bar{u}_1}^{\bar{u}_1} f_h(\bar{u}) < (\bar{u}_1 - T_2 \bar{u}_1) f_h(T_2 \bar{u}_1) \quad \text{in both a) and b)}$$

So just apply 1.5.14. ii) which holds for ϵ and hence δ , sufficiently small. Finally, if $T_2 \bar{u}_1 > \delta$, then $S_{h\delta}(\bar{u})$ is sub-regular immediately by 1.5.14. since $U_{\theta(\bar{u})} \gamma_{h\delta} - \gamma_{h\delta} = \rho_h(\bar{u})$ on $T_2 \bar{u}_1 > \delta$. Put $k^* = k^s$ which we found in 1.5.14. and just choose $\delta(k^*) < \epsilon$ we used for partitioning I.

Hence $S_{h\delta}(\bar{u})$ is sub-regular on I for $0 < k < k^*$ and $0 < \delta < \delta(k^*)$.

//

i) Now if $\theta(\bar{u}) \sim O(\bar{u})$ as $\bar{u} \downarrow 0$ at each boundary, we choose a family $S_{h\delta}(\bar{u})$ which are sub-regular on I and take $\lim_{\delta \downarrow 0} S_{h\delta}(\bar{u}) = 1$ for $0 < k < k^*$ and $\bar{u} \neq 0$.

Thus $\lim_{\delta \downarrow 0} S_{h\delta}(\bar{u}) = \gamma_{\theta(\bar{u})}(\bar{u}) = \gamma_{opt}$ and hence for $\alpha > 1$, we have optimality.

ii) By remark 1.5.15 a), this optimality is only dependent on the behaviour of $\theta(\bar{u})$ at the boundaries $[0, \epsilon]$ and $(1 - \epsilon^*, 1]$.

//

Theorem 1.5.17.

If $\alpha = 1$ at each boundary and rule is of the form $\theta(\bar{u}) = \theta \bar{u}(1 - \bar{u})$, $q_1 > q_2$,

then:- a) If $1 > \theta > 6(1 - q_2/q_1)$ then the rule is ϵ -optimal.

(and with more careful analysis we obtain $1 > \theta > (1 - q_2/q_1) / (1 - \log 2)$)

b) If $\theta < 1$ and $\theta < (1 - q_2/q_1)$ then the rule is optimal.

(and with more careful analysis we obtain $\theta < (1 - q_2/q_1) 2(1 + \theta_{1/6}) / (1 + q_2/q_1)$)

Proof.

For the conditionally optimal transition case, learning is due both to boundary and central behaviour so we must assert the form of $\theta(\bar{u})$ on $I = [0, 1]$, rather than just at the boundaries, as for $\alpha > 1$.

I shall prove both a) and b) using $\gamma_k(x) = \int_0^x \left(\frac{1 - \bar{u}}{\bar{u}} \right)^k d\bar{u}$,

in a similar way to that in which Norman (1968) proves 1.5.4.

First, I consider the simpler to prove, b) for which I construct a sub-regular family $\gamma_k(x)$ with $k \uparrow 1$, and $\lim_{k \uparrow 1} \gamma_k(x)/\gamma_k(1) = \gamma_{opt}(x)$ as in 1.5.12.

$$i) \quad U_{\theta(\bar{u})} \gamma_k(x) - \gamma_k(x) = q_1 \bar{u}_1 \int_{\bar{u}_1}^{\bar{u}_2} \left(\frac{1-\bar{u}}{\bar{u}} \right)^k d\bar{u} - q_2 \bar{u}_2 \int_{\bar{u}_1}^{\bar{u}_2} \left(\frac{1-\bar{u}}{\bar{u}} \right)^k d\bar{u}$$

and since $\left(\frac{1-x}{x} \right)^k$ is monotone \downarrow on $[0, 1]$.

$$\begin{aligned} U_{\theta(\bar{u})} \gamma_k(x) - \gamma_k(x) &\geq q_1 \bar{u}_1 \theta \bar{u}_1 \bar{u}_2^2 \left[\left(\frac{1-\bar{u}_1 - \theta \bar{u}_1 (1-\bar{u}_1)}{\bar{u}_1 + \theta \bar{u}_1 (1-\bar{u}_1)} \right)^k \right] - q_2 \bar{u}_2 \theta \bar{u}_1^2 \bar{u}_2 \left[\frac{(1-\bar{u}_1 + \theta \bar{u}_1 (1-\bar{u}_1))^k}{\bar{u}_1 (1-\theta \bar{u}_1 (1-\bar{u}_1))} \right] \\ &\geq \theta \bar{u}_1^2 \bar{u}_2^2 \bar{u}_2^k / \bar{u}_1^k \left[q_1 \left[\frac{(1-\theta \bar{u}_1 \bar{u}_2)}{(1+\theta \bar{u}_1^2)} \right]^k - q_2 \left[\frac{(1+\theta \bar{u}_1^2)}{(1-\theta \bar{u}_1 \bar{u}_2)} \right]^k \right] \\ &\geq \theta \bar{u}_1^2 \bar{u}_2^2 \bar{u}_2^k / \bar{u}_1^k B(\bar{u}_1) \quad \text{with } \bar{u}_2 = 1 - \bar{u}_1. \end{aligned}$$

Now we find $\min_{\bar{u}_1} B(\bar{u}_1)$.

Lemma. 1.5.18.

$\left. \frac{d}{d\bar{u}_1} B(\bar{u}_1) \right|_{\bar{u}_1 = \frac{1}{2}} = 0$ and this is a minimum and is unique.

Proof.

$$\begin{aligned} \frac{d}{d\bar{u}_1} [q_1 (1 - \theta \bar{u}_1 \bar{u}_2)^{2k} - q_2 (1 + \theta \bar{u}_1^2)^k (1 + \theta \bar{u}_1^2)^k] \\ = 2\theta k (\bar{u}_1 - \bar{u}_2) [q_1 (1 - \theta \bar{u}_1 \bar{u}_2)^{2k-1} + q_2 (1 - \theta \bar{u}_1 \bar{u}_2) [(1 + \theta \bar{u}_1^2)(1 + \theta \bar{u}_1^2)]^{k-1}] \end{aligned}$$

and $\frac{\partial B}{\partial \bar{u}_1} = 0$ iff $\bar{u}_1 = \frac{1}{2}$ and $\left. \frac{\partial B}{\partial \bar{u}_1} \right|_{\bar{u}_1 = \frac{1}{2} + \epsilon} > 0$, $\left. \frac{\partial B}{\partial \bar{u}_1} \right|_{\bar{u}_1 = \frac{1}{2} - \epsilon} < 0$.

//

Hence we ensure $U_{\theta(\bar{u})} \gamma_k(\frac{1}{2}) - \gamma_k(\frac{1}{2}) > 0$ or $q_1 (1 - \frac{\theta}{4})^2 > q_2 (1 + \frac{\theta}{4})^2$.

Thus $\theta < \mathcal{E}(q_1 - q_2) / (q_1 + q_2)$ with $\mathcal{E} = 2(1 + \frac{\theta^2}{16})$.

Then for $0 < 1$, $2(1 + \frac{\theta^2}{16}) > 2$

and so we need satisfy only,

$$\theta < (1 - \frac{q_2 q_1}{(1 + q_2 q_1)})_{1/2}$$

But $(q_1 + q_2) < 2q_1$ and hence

we get b):- $\theta < 1$ and $\theta < (1 - \frac{q_2 q_1}{(1 + q_2 q_1)}) \Rightarrow$ rule is optimal.

The more precise inequality will allow us to give certain useful counter-examples in 1.5.20.

Lemma. 1.5.19.

Given any $\theta < 4$, and hence $\theta(\frac{1}{2}) < 1$, then $\exists q_1, q_2$ s.t. the rule is optimal.

Proof.

If $q_2 = 0, q_1 \neq 0$ then we need only satisfy $(\theta/4 - 1)^2 \geq 0$ which holds for all real θ . Clearly if $q_1 \neq 0$ and q_2 is sufficiently small, we shall still satisfy the inequality 1.5.17. b), in its strongest form, by using continuity arguments. Note $(\theta/4 - 1)^2 = 0$ iff $|\theta| = 4$.

//

ii) For a) we prove that $\exists k^* < 1$ s.t. $\forall k : 1 > k > k^*$ then $\gamma_k(h) = \int_0^h (1 - \frac{u}{h})^k du$ is super-regular. Thus $\lim_{n \rightarrow \infty} U_{\theta(h)}^n \gamma_k(h) = \gamma_{\theta(h)}(h) < \gamma_k(h)$ and the rule is at best ϵ -optimal.

Again we use the argument based on partitioning. $I_1 = [0, \delta], I_2 = [\delta, 1 - \delta^*], I_3 = [1 - \delta^*, 1]$.

Consider I_1 . We find k_1 s.t. γ_{k_1} is super-regular on I_1 , for $1 > k > k_1$.

$$U_{\theta(h)} \gamma_k - \gamma_k = q_1 \bar{u}_1 \int_{\bar{u}_1}^{\bar{u}_2} (1 - \frac{u}{h})^k du - q_2 \bar{u}_2 \int_{\bar{u}_2}^{\bar{u}_1} (1 - \frac{u}{h})^k du \\ \leq \frac{\bar{u}_1^{1-k}}{(1-k)} [q_1 \bar{u}_1 [(1 + \theta \bar{u}_2^k)^{1-k} - 1] - q_2 \bar{u}_2^{1+k} [1 - (1 - \theta \bar{u}_2^k)^{1-k}]]$$

Thus $U \gamma < \gamma$ if $q_1 \bar{u}_1 [\theta \bar{u}_2^k (1-k) - k(1-k) \theta^2 \bar{u}_2^{k+2} + k(1-k)(1+k) \theta^3 \bar{u}_2^{k+6}] < q_2 \bar{u}_2^{1+k} (\theta(1-k) \bar{u}_1 \bar{u}_2)$

or $(q_1 - q_2 \bar{u}_2^k) < k q_2 \theta [\bar{u}_2^{k+2} - (1+k) \bar{u}_2^k \theta/6]$

or $\theta > 2(1 - (q_1/q_2) \bar{u}_2^k) / (k \bar{u}_2^2 (1 - (1+k) \bar{u}_2^2 \theta/3))$

Now since δ can be made arbitrarily close to 0, and k can be arbitrarily close to 1. $\theta > 2(1 - q_1/q_2) / (1 - 2\theta/3)$ will give $U \gamma < \gamma$

and for $\theta < 1$, $\theta > 6(1 - q_1/q_2)$ for $(1-k), \delta$ sufficiently close to 0.

Further, we have $\theta > 2(1 - q_1/q_2) / (1 - 2\theta/3 + 2\theta^2/4 - 2\theta^3/5 + \dots)$ on taking further terms.

and if $\theta < 1$ we need satisfy only $\theta > (1 - q_1/q_2) / (1 - \log_e 2)$.

or for $\theta < 1$ (since our binomial expansions hold only for $\theta < 1$),

we get more precisely $\theta < (1 - \log_e(1 + \theta)) / (1 - q_1/q_2)$ gives ϵ -optimality.

It also seems intuitively clear that if $\theta^* < 1$ gives ϵ -opt, then all θ s.t.

$1 > \theta > \theta^*$ give at best only ϵ -opt, but we need the conjectured generalisation

to 1.4.6. to prove this naturally.

Consider I_3 . Again, on approximating the integrals, we get,

$$U\tau_k - \tau_k \leq \bar{u}_2^{1+k} / (1+k) (q_1 \bar{u}_1^{1-k} (1 - (1 - \theta \bar{u}_1 \bar{u}_2)^{1+k}) - q_2 \bar{u}_2 (1 + \theta \bar{u}_1)^{1+k} - 1)$$

and $U\tau < \tau$ if $q_1 \bar{u}_1^{1-k} (\theta \bar{u}_1 \bar{u}_2 (1+k)) < q_2 \bar{u}_2 (\theta \bar{u}_1^2 (1+k) + \theta^2 k(k+1) \bar{u}_1^{3/2} + \epsilon)$

where $\epsilon = (1-k) f(\theta, k, \bar{u})$ and $\lim_{k \uparrow 1} \epsilon = 0$.

Then as for I_1 , we take k sufficiently close to 1, and here also \bar{u}_1 sufficiently close to 1. This gives $\tau_k(x)$ super-regular on I_3 if

$$1+k > k_2 \text{ and } 1+\theta > 2(q_2/q_1 - 1)$$

However, comparing I_1 and I_3 conditions:- $q_1(1-2\theta/3) < q_2$ if $\theta > 3/2(1-q_2/q_1)$

Thus the stricter condition is that at I_1 , and so we may disregard that at I_3 .

Consider I_2 . Here we take $k = 1$ and use compactness.

$$U_{\theta(\bar{u})} \tau_{k=1} - \tau_1(x) = (q_2 - q_1) \theta \bar{u}_1^2 \bar{u}_2 + q_1 \bar{u}_1 \log(1 + \theta \bar{u}_2^2) + q_2 \bar{u}_2 \log(1 - \theta \bar{u}_1 \bar{u}_2)$$

and $U\tau < \tau$ when $\theta > 2(q_2 - q_1) / (q_1 \bar{u}_2 (1 - \theta \bar{u}_1^2) + q_2 \bar{u}_1)$

or more precisely $\theta > (q_2 - q_1) / (q_1 \bar{u}_2 + q_2 \bar{u}_1 (\theta \bar{u}_1^2 - \log(1 + \theta \bar{u}_2^2))) / \theta^2 \bar{u}_1^2$

or $\theta > 2(q_2 - q_1) / (q_2 \bar{u}_1 + q_1 \bar{u}_2 (1 - 2\theta \bar{u}_1^2/3 + 2\theta^2 \bar{u}_1^4/4 - 2\theta^3 \bar{u}_1^6/5 + \dots))$

Thus $U\tau < \tau$ if above inequality holds. We must now use compactness to relate $k=1$ to $k^* < 1$ on $[\delta, 1-\delta^*]$.

Since $\int_{\delta}^{1-\delta^*} (1-\bar{u})/d\bar{u}$ is bounded on $x \in I_2$ (compact), we have $|U\tau_1 - \tau_1(\bar{u}) - (U\tau_k(\bar{u}) - \tau_k(\bar{u}))| < \epsilon'$ for $1+k > k_3$ say, and $\bar{u} \in I_2$

The ϵ' is chosen s.t. $(\tau_1 - U\tau_1) > 2\epsilon'$ and hence $\tau_k - U\tau_k > \epsilon'$ and

$\tau_k(\bar{u})$ is super-regular on I_2 . Now $\min_{\bar{u}_1, \bar{u}_2=1} (q_2 \bar{u}_1 + q_1 \bar{u}_2 (1 - 2\theta \bar{u}_1^2/3 + \dots))$ exists by compactness and with some difficulty can be found to lie at the lower end-point $\bar{u}_1 = \delta$. (and θ must satisfy the given inequality at all times.)

Finally, we combine the results and obtain the required bounds.

We have $\gamma_k(u)$ is super-regular on $I: [0,1]$ if $1 > \theta > (1 - q_{1/2}) / (1 - \log_e 2)$.
 where $1 > k > \max(k_1, k_2, k_3)$. And hence we have that the rule is then ϵ -optimal.

//

To end this rather lengthy section of analysis, we shall consider how # alternations relates to optimality in the transition class. We shall denote this class by \mathcal{R}_c since we have now proved in 1.5.17. that such rules are only conditionally optimal.

And $\mathcal{R} = \mathcal{R}_\epsilon \cup \mathcal{R}_c \cup \mathcal{R}_0$.

Lemma 1.5.20.

- i) If $1 > \theta > (1 - q_{1/2}) / (1 - \log_e 2)$ and $\theta(u) = \theta u(1-u)$ then $\#alts = \infty$ and the rule is ϵ -optimal.
- ii) If $1 > \theta > 1$ and $\theta q_1 > 1 > \theta q_2$ then we may choose $\epsilon > 0$ s.t. $q_2 < \epsilon$ and the rule is optimal with $\#alt < \infty$.

Proof.

- i) We use 1.5.7. and note that $0, q_1, q_2 < 1$ gives $\#alts = \infty$, but from 1.5.17, the constraint on θ gives ϵ -optimality.
- ii) Using the finer bound of 1.5.17 b) we have $\theta < (1 - q_{1/2}) / (2(1 + \theta^{1/2}))$ gives optimal learning.
 But as $q_1 \downarrow 0$, we find $\theta \downarrow 0$ will satisfy the above for ϵ sufficiently small. Also $\theta q_1 > 1 \Rightarrow \#alt < \infty$ for the rule, since we have just seen that it is optimal.

//

Lamperti and Suppes (1960) have investigated a family of conditionally optimal learning rules, called β -rules. These give $\#alts < \infty$, but only when their behaviour is actually optimal, for convergence to the boundaries only occurs for certain q and β_{ij} . Since our conditionally optimal class is U.L., we do always converge giving at least ϵ -optimality. Optimality is determined by the relation between the learning parameter θ (or β_{ij} for β -rule) and $q_{1/2}$ (similarity between the actions).

Our bounds on conditional optimality are reasonably tight since the diffusion limit has its optimality transition at $\theta = 2(1 - q_{2,q_1})$. In discrete time 1.5.17. asserts that if $\theta < 1$, $\theta < 2(1 - q_{2,q_1}) / (1 + q_{2,q_1})$ gives optimality and $\theta > 3(1 - q_{2,q_1})$ gives ϵ -optimality. This last inequality becomes arbitrarily close to the diffusion limit as $\theta \downarrow 0$, for consider $\theta > 2(1 - q_{2,q_1}) / (1 - 2\theta/3)$. Thus it seems reasonable to conjecture $\theta = K(\theta) / (1 - q_{2,q_1})$ with $K(\theta) > \theta$ and $\lim_{\theta \rightarrow 0} K(\theta) = 2$, $\lim_{\theta \rightarrow 4} K(\theta) = 4$, $K(\theta) \uparrow$ as $\theta \uparrow$, gives the discrete time transition. Then define $F(\theta) = 1 - \theta / K(\theta)$ so that $(q_{2,q_1}) = F(\theta)$ gives the transition, with $(q_{2,q_1}) < F(\theta)$ giving optimality. As in the generalised comparison theorem, our methods of super and sub-regularity are probably not subtle enough to give us it.

1.6. n-Action Extensions.

We define $\gamma_{\theta(\bar{u})}^i(\bar{u}) = \Pr(\lim_{t \rightarrow \infty} u_i(t) = 1 \mid \bar{u}(0) = \bar{u})$ and if $(q_i - q_j) > 0$ $i \neq j$, we use $\gamma_{\theta(\bar{u})}^i = \gamma_{\theta(\bar{u})}^i$. Clearly we must have boundary conditions $\gamma_{\theta(\bar{u})}^i(e_j) = \delta_{ij}$ where e_j is the unit vector along the j -axis. Also $\gamma_{\theta(\bar{u})}^i(\bar{u} \text{ s.t. } \bar{u}_i = 0) = 0$.

Lemma 1.6.1.

$U_{\theta(\bar{u})} \gamma_{\theta(\bar{u})}^i(\bar{u}) = \gamma_{\theta(\bar{u})}^i(\bar{u})$ where $U_{\theta(\bar{u})} \gamma(\bar{u}(t)) = E(\gamma(\bar{u}(t+1)) \mid \bar{u}(t))$.

and the solution is unique if $\gamma_{\theta(\bar{u})}^i$ is continuous.

Proof.

We have convergence by 1.3.1. $\lim_{t \rightarrow \infty} \bar{u}_j(t) \in \{0, 1\}$ $\forall j$.

Now $\lim_{n \rightarrow \infty} U_{\theta(\bar{u})}^n \gamma^i(\bar{u}) = 0 \cdot \Pr(\bar{u}_i(\infty) = 0, \bar{u}_j = \bar{u}_j^* \text{ } j \neq i) + 1 \cdot \Pr(\bar{u}_i(\infty) = 1, \bar{u}_j(\infty) = \bar{u}_j^* \text{ } j \neq i)$
 $= \gamma_{\theta(\bar{u})}^i(\bar{u})$ with $\gamma^i(e_j) = \delta_{ij}$ and $\gamma^i(\bar{u} \text{ s.t. } \bar{u}_i = 0) = 0$.

And putting $\gamma^i(\bar{u}) = \gamma_{\theta(\bar{u})}^i(\bar{u})$ we have the result. Uniqueness for continuous functions is proved as in 1.4.1. Thus $\Delta \gamma_{\theta(\bar{u})}^i(\bar{u}) \equiv 0$.

//

We shall now use the methods of super and sub-regularity, as in section 4, and we shall extend our results for R_e, R_o to n -actions. The $\theta(\bar{u})$ suffix on $\gamma_{\theta(\bar{u})}$ just denotes non-linear U.L., boundary convergent rules, which includes the $\theta_{ij}(\bar{u})$ family.

Note that the discontinuous functions $\delta'_{\theta(i)}(\bar{u}) \equiv 1$, except for $\delta'_{\theta(i)}(\bar{u}_i) = 0$, and $\delta'_{\theta(i)}(\bar{u}) \equiv 0$ except for $\delta'_{\theta(i)}(\bar{u}_i) = 1$, always give $\Delta \delta'_{\theta(i)}(\bar{u}) \equiv 0$. So we always have the existence of solutions to $\Delta \delta \equiv 0$ but they will not necessarily be continuous or unique. However, we do have $\lim_{n \rightarrow \infty} U^n \delta'(\bar{u}) = \delta'(\bar{u})$ holding only for the actual absorption probabilities, so we discriminate between solutions using our super and sub-regular $\psi'(\bar{u})$. In this way we can determine whether we indeed do have a true discontinuous solution, expressing optimality, as in 1.5.12. We can actually obtain an infinity of solutions by setting $\sum_i r_i = 1$ $0 \leq r_i \leq 1$ and $\delta'(\bar{u}) \equiv r_i$, except for the usual boundary conditions of 1.6.1., so that $\Delta \delta'(\bar{u}) \equiv 0$. But these discontinuous functions do not represent actual absorption probabilities, apart from the optimality case, $r_{i_0} = 1$.

Lemma 1.6.2.

Let q_1, q_2, q'_2 and environment 1 has parameters q_1, q_2 .
 " 2 " " q_1, q'_2 .

And define $\delta'_{\theta(i),i}(\bar{u}) = \delta'_i(\bar{u})$ for environment i .

Then if δ'_1 is monotone, $\delta'_2 \geq \delta'_1$.

Proof.

We show $U_2 \delta'_1 \geq \delta'_1$ and hence $\lim_{n \rightarrow \infty} U_2^n \delta'_1 = \delta'_1 \geq \delta'_2$.

Now $U_2 \delta'_1 - \delta'_1 = \bar{u}_1 q_1 [\delta'_1(\bar{u}_1) - \delta'_1(\bar{u})] - \bar{u}_2 q'_2 [\delta'_1(\bar{u}) - \delta'_1(\bar{u}_2)]$.

and $U_1 \delta'_1 = \delta'_1$ so $U_2 \delta'_1 - \delta'_1 = \bar{u}_2 (q_2 - q'_2) (\delta'_1(\bar{u}) - \delta'_1(\bar{u}_2)) \geq 0$

by the monotonicity of $\delta'_1(\bar{u})$. //

We now extend this result to n -actions and in so doing, we prove an argument used by Narendra and Viswanathan (1972).

Lemma 1.6.3.

Let $q_1 > q_2 > \dots > q_n$ and define $\delta'_{\theta(i),i}(\bar{u}) = \delta'_i(\bar{u})$.

Now associate $\delta'_i(\bar{u})$ with parameters $q_1 > q_2 = q'_3 = \dots = q'_n$.

Then if $\delta'_1(\bar{u}) > \delta'_2(\bar{u})$ when $\bar{u}_1 > \bar{u}'_1$ (monotone in first argument as in the linear rule $\theta = \text{const}$), we have $\delta'_1(\bar{u}) > \delta'_2(\bar{u})$.

Proof.

We prove $U, \bar{x}_2 \geq \bar{x}_2$

$$U, \bar{x}_2 - \bar{x}_2 = \bar{u}_1 q_1 (\bar{x}_2(T, \bar{u}_1) - \bar{x}_2(\bar{u})) - \sum_{r=2}^n \bar{u}_r q_r (\bar{x}_2(\bar{u}) - \bar{x}_2(T, \bar{u}_r))$$

and $U, \bar{x}_2 = \bar{x}_2$

Then $U, \bar{x}_2 - \bar{x}_2 = \sum_{r=2}^n (q_2 - q_r) \bar{u}_r (\bar{x}_2(\bar{u}) - \bar{x}_2(T, \bar{u}_r)) \geq 0$.

since $\bar{x}_2(\bar{u}) \geq \bar{x}_2(T, \bar{u}_r)$ by monotonicity in first argument.

In particular $U, \bar{x}_2 - \bar{x}_2 = (\bar{x}_2(\bar{u}) - \bar{x}_2(T, \bar{u})) \sum_{r=2}^n (q_2 - q_r) \bar{u}_r \geq 0$ for $\theta_{ij}(\bar{u}) = \text{const} = \theta$.

and hence $\bar{x}_2(\bar{u}) \geq \bar{x}_2(\bar{u})$ for linear ($\alpha=0$) learning. This immediately gives the linear rule as n -action, ϵ -optimal, which is used for comparison in 1.6.6. //

Theorem 1.6.4.

In static environment \mathcal{M} with $q_1 = \dots = q_m > q_{m+1} \geq \dots \geq q_n$
with learning under $\theta_{ij}(\bar{u})$ rules with $\theta_{ij}(\bar{u}) \sim O(\bar{u}_i)^{\alpha_i}, \alpha_i > 1$ as $\bar{u}_i \downarrow 0$
 $\sim O(\bar{u}_j)^{\alpha_j}, \alpha_j > 1$ as $\bar{u}_j \downarrow 0$

Then $\exists i \in M$ s.t. $\lim_{t \rightarrow \infty} \bar{u}_i(t) = 1$ where $M = \{i : i \leq m\}$.

Proof.

Let $M^c = \{m+1, \dots, n\}$. Then under $\theta_{ij}(\bar{u})$ rules we have:-

i) Convergence to $\bar{u}_i \in \{0, 1\} \forall i$ by 1.3.1.

ii) U.L. for all \bar{u} .

If $q_m = q_j \forall j$ then we have nothing to prove.

Let $\bar{u}_m = \sum_{i \in M} \bar{u}_i$ and similarly $1 - \bar{u}_m = \sum_{i \in M^c} \bar{u}_i$

We have $\bar{u}_m(t+1) = \bar{u}_m(t) + \sum_{j=m+1}^n \theta_{ij}(\bar{u}) \bar{u}_j$ for $U_i(t)$ s.t. $i \in M$.
 $\bar{u}_m(t+1) = \bar{u}_m(t) - \sum_{j \in M} \theta_{ij}(\bar{u}) \bar{u}_j$ for $U_i(t)$ s.t. $i \in M^c$.

And $\sum_{j \in M^c} \theta_{ij}(\bar{u}) \bar{u}_j \leq \max_{\substack{i \in M \\ j \in M^c}} \theta_{ij}(\bar{u}) (1 - \bar{u}_m)$.

And $\sum_{j \in M} \theta_{ij}(\bar{u}) \bar{u}_j \leq \max_{\substack{i \in M^c \\ j \in M}} \theta_{ij}(\bar{u}) \bar{u}_m$.

We could write $\bar{u}_m \rightarrow \bar{u}_m + \left(\frac{\sum_{j \in M^c} \theta_{ij}(\bar{u}) \bar{u}_j}{\sum_{j \in M^c} \bar{u}_j} \right) (1 - \bar{u}_m)$ $i \in M$

and $\bar{u}_m \rightarrow \bar{u}_m \left(1 - \left(\frac{\sum_{j \in M} \theta_{ij}(\bar{u}) \bar{u}_j}{\sum_{j \in M} \bar{u}_j} \right) \right)$ $i \in M^c$.

Then w.l.o.g. we put $\theta_{ij}(\bar{u}) = \theta(\bar{u}_i, \bar{u}_j)^*$ and $\theta_m^* = \theta(\bar{u}_m(1 - \bar{u}_m))$.

Clearly $\max_{\substack{i \in M \\ j \in M^c}} \theta_{ij}(\bar{u}) \leq \theta_m^*$.

Thus we reduce n -actions to 2-actions with $q_1^* = q_1$, $q_2^* = q_{m+1}$, and learning rule $\theta_m^*(\bar{u}_m)$.

Intuitively, in our reduction, we are always increasing the drift, yet we prove that the final 2-action, θ_m^* process is optimal, and then by comparison, so is the original process.

$$\Delta \bar{u}_m = \sum_{i \in M} \bar{u}_i \sum_{j \in M^c} \bar{u}_j (q_i - q_j) \theta_{ij}(\bar{u}) < (q_1^* - q_2^*) \bar{u}_m (1 - \bar{u}_m) \theta_m^*(\bar{u}_m)$$

Note that the 2-action, θ_m^* process is U.L. and boundary absorbed.

Lemma 1.6.5.

Define $\mathcal{J}_1(\bar{u}) = \Pr(\bar{u}_i \rightarrow 1, i \in M \mid \bar{u}(0) = \bar{u})$ under $q_1^* = \dots = q_m^* > q_{m+1}^* = \dots = q_n^*$.

$\mathcal{J}_2(\bar{u}) = \Pr(\bar{u}_i \rightarrow 1, i \in M \mid \bar{u}(0) = \bar{u})$ under $q_1^* = \dots = q_m^* > q_{m+1}^* = \dots = q_{m+1}^*$.

and $\mathcal{J}_2^*(\bar{u}_m) = \Pr(\bar{u}_m \rightarrow 1 \text{ under } \theta_m^*(\bar{u}_m), 2\text{-action rule})$.

Then for concave \mathcal{J}_2^* we have $\mathcal{J}_2(\bar{u}_1, \dots, \bar{u}_n) \geq \mathcal{J}_2^*(\bar{u}_m, 1 - \bar{u}_m) \quad \forall \bar{u}$.

Proof.

For concave \mathcal{J}_2^* we can use the 2-action comparison theorem 1.4.5., which is easily extended to hold for such concave $\mathcal{J}(\bar{u})$. Note that here we have action 1 as optimum which gives concave $\mathcal{J}(\bar{u})$ whilst in 1.4.5. action 2 is optimum, giving convex $\mathcal{J}(\bar{u})$.

If $\theta_{ij}(\bar{u}) \leq \theta_m^*$ then $\bigcup_{i,j} \mathcal{J}_2^* \geq \mathcal{J}_2^*$, where we have partition $\{1, \dots, m\}$, $\{m+1, \dots, n\}$ and reward probabilities q_1^* and q_2^* , with learning function $\theta_{ij}(\bar{u})$. So we update as if we had n -actions with $\theta_{ij}(\bar{u})$, yet we compare with 2-action rule $\theta_m^*(\bar{u}_m)$.

But by definition $\theta_m^*(\bar{u}_m) \geq \max_{i \in M} \theta_{ij}(\bar{u}) \quad \forall \bar{u}$ and hence
 if we fix $\bar{u}(\bar{u})$, we have the result $\theta_2(\bar{u}_1, \bar{u}_0) \geq \theta_2^*(\bar{u}_m, 1-\bar{u}_m) \quad \forall \bar{u}.$
 with $\bar{u}_m = \sum_{i \in M} \bar{u}_i$. //

Now a) $\theta_m^* \sim 0$ ($\bar{u}_m \downarrow 0$) as $\bar{u}_m \downarrow 0$
 ~ 0 ($1-\bar{u}_m \uparrow 1$) as $\bar{u}_m \uparrow 1$

Hence by 1.5.12. $\theta_m^*(\bar{u}_m)$ is optimally 2-action learning, and so
 is concave.

b) By 1.6.5. we have $\theta_{ij} \geq \theta_m^* \quad \forall \bar{u}$ so that we
 have n-action optimality over $q_1 = \dots = q_m, q_{m+1} = \dots = q_{m+1}$.

c) By 1.6.3. for monotone θ_{ij} we have $\theta_i \geq \theta_j$ and
 hence $\theta_{ij}(\bar{u}) \equiv 1$ except at $\bar{u}_m = 0$.
 So $\exists i \in M$ s.t. $\lim_{t \rightarrow 0} \bar{u}_i(t) = 1$.

This proof may seem artificial since the lemmas are being
 used in the "degenerate" cases. However, as in 1.5.12., we could
 consider the discontinuous optimal limit $\theta_{ij}(\bar{u})$ as the limit
 of a sub-regular family and then obtain the n-action optimality.

//

Intuitively, the $\theta_{ij}(\bar{u})$ rules compare actions in pairs, and
 no optimal action is allowed to vanish, by a simple extension
 of asymptotic reflection from sub-optimal boundaries. This is the
 principle of boundary learning which will be treated conceptually
 in the next section. Note that $\theta_{ij}(\bar{u}) < 1$ must always hold.

We now extend our ϵ -optimality result to all n-action
 boundary absorbed rules.

Corollary 1.6.6.

The family of $\theta_{ij}(\bar{u})$ rules and the boundary absorbed $\theta(\bar{u})$ rules

are all at least ϵ -optimal.

$$\lim_{\theta \downarrow 0} \lim_{t \rightarrow \infty} \pi_m(t) = 1.$$

Proof.

We have $U_{\theta_{ij}} \mathcal{J}_{\theta_{ij}} = \mathcal{J}_{\theta_{ij}}$ and we prove $U_{\theta_{ij}} \mathcal{J}_{\theta} \geq \mathcal{J}_{\theta}$ where

\mathcal{J}_{θ} corresponds to $q_1 > q_2 = \dots = q_n$, $\theta = \text{const}$, with a 2-action or n-action rule, as identical in result.

$\mathcal{J}_{\theta_{ij}}$ corresponds to $q_1 > q_2 > \dots > q_n$ $\theta_{ij}(\bar{u})$ rule.

If there are multiple optima, then the proof proceeds as

1.6.4., with partitioning $M \vee M^c$.

$$\text{Now } U_{\theta_{ij}} \mathcal{J}_{\theta} - \mathcal{J}_{\theta} = q_1 \bar{u}_1 (\mathcal{J}(T_1^* \bar{u}_1) - \mathcal{J}(\bar{u}_1)) - \sum_{j \neq 1} q_j \bar{u}_j (\mathcal{J}(\bar{u}_1) - \mathcal{J}(T_{2,j}^* \bar{u}_1)).$$

$$\text{where } T_{2,j}^* \bar{u}_1 = \bar{u}_1 (1 - \theta_{ij}(\bar{u}))$$

And since $\mathcal{J}_{\theta}(\bar{u}_1)$ is concave, we proceed as in the 2-action proof, 1.4.6.

$$U_{\theta_{ij}} \mathcal{J}_{\theta} - \mathcal{J}_{\theta} \geq \theta_{ij}(\bar{u}) q_1 \bar{u}_1 \bar{u}_2^* \left(\frac{\mathcal{J}(T_1^* \bar{u}_1) - \mathcal{J}(\bar{u}_1)}{T_1^* \bar{u}_1 - \bar{u}_1} \right) - q_2 \bar{u}_2 \left(\sum_{j \neq 2} \theta_{ij} \bar{u}_j \right) \left(\frac{\mathcal{J}(\bar{u}_1) - \mathcal{J}(T_{2,j}^* \bar{u}_1)}{\bar{u}_1 - T_{2,j}^* \bar{u}_1} \right).$$

And choose θ s.t. $\theta_{ij}(\bar{u}) = 0 \quad \forall i, j, \bar{u}$. with $\bar{u}_2^* = 1 - \bar{u}_1$.

$$\begin{aligned} &\geq \theta_{ij}(\bar{u}) q_1 \bar{u}_2^* \bar{u}_1 \left(\frac{\mathcal{J}(T_1^* \bar{u}_1) - \mathcal{J}(\bar{u}_1)}{T_1^* \bar{u}_1 - \bar{u}_1} \right) - q_2 \bar{u}_2 \left(\sum_{j \neq 2} \theta_{ij} \bar{u}_j \right) \left(\frac{\mathcal{J}(\bar{u}_1) - \mathcal{J}(T_{2,j}^* \bar{u}_1)}{\bar{u}_1 - T_{2,j}^* \bar{u}_1} \right) \\ &\geq \theta_{ij}(\bar{u}) / \theta (U_{\theta} \mathcal{J}_{\theta} - \mathcal{J}_{\theta}) = 0. \end{aligned}$$

Then by 1.6.3. we have $\lim_{\theta \downarrow 0} \mathcal{J}_{\theta}(\bar{u}) = \mathcal{J}(\bar{u}) \neq 0$ and hence we also have $\theta_{ij}(\bar{u})$ rules are all at least ϵ -optimal.

$$\mathcal{J}_{\theta_{ij}}(\bar{u}) \geq \mathcal{J}_{\theta}(\bar{u}).$$

We now have the n-action generalisation of 1.4.5.

Clearly the above goes through with $\theta(\bar{u})$ boundary absorbed, even more easily, to give ϵ -optimality.

//

Corollary 1.6.7.

$$\text{If } \left. \begin{aligned} \theta_{ij}(\bar{u}) &\sim O(\bar{u}_i)^{\alpha} & \bar{u}_i \downarrow 0 \\ &\sim O(\bar{u}_j)^{\alpha} & \bar{u}_j \downarrow 0 \end{aligned} \right\} \forall i, j.$$

- then
- i) If $\alpha < 1$, the rule is at best ϵ -optimal.
 - ii) If $\alpha > 1$, the rule is optimal.

Proof.

By 1.6.6., all the rules are at worst ϵ -optimal, but with $i(\bar{u}) > 0$ where $i(\bar{u}) = \text{prob (just take action } i) \text{ for } \alpha < 1$ we have $E(\#alts)^n < \infty$. Whilst for $\alpha > 1$ we get the result from 1.6.4.

//

Note that it is unresolved if the $\theta(\bar{u}) = K(\bar{u}_i(1-\bar{u}_i)^\alpha)$ are optimal for certain α , but 1.6.6 gives them all ϵ -optimal with $\gamma_{\theta(\bar{u})} \geq \gamma_\theta$.

We have now considered R_ϵ, R_0 extensions, leaving the difficult R_c . Here we can use 1.6.3. to reduce the reward parameters to $q_1 > q_2 = \dots = q_2$ and for optimality we can find K s.t. $\theta < K(1-q_2)$ and proceed in reverse as in 1.6.4. with our concave $\theta_{opt}(\bar{u})$. However, for ϵ -optimality we cannot use 1.6.3. easily for we require monotonicity, and 1.6.5. needs concavity, and fails to work in reverse to give bounds away from optimality. Thus a thorough treatment would seem to require the construction of a n -action sub-regular family $\theta_k(\bar{u})$ in analogy to the 2-action fundamental techniques of R_c analysis.

It is possible that the $\theta_{ij}(\bar{u})$ are uniquely optimal family of rules. But further work remains to be done on properties such as the $E\#alts$, for $K(\bar{u}_i(1-\bar{u}_i)^\alpha) = \theta(\bar{u})$ rules, if they are to be excluded from consideration. We saw in 1.2 that the $\theta(\bar{u})$ and $\theta_{ij}(\bar{u})$ rules essentially exhaust U.L. possibilities, and this holds even if we have $R - P$ rules with T_{ij} and S_{ij} in use.

The use of T_{ij} , gives difficulties in normalisation, and the corresponding $\theta_{ij}(\bar{u})$ rules can only be defined when $\theta_{ij}(\bar{u}) \sim K(\bar{u}_i \bar{u}_j)^{\alpha \geq 1}$.

$$\bar{u}_i(t+1) = \bar{u}_i(t) \pm \theta_{ii}(\bar{u}) (1 - \bar{u}_i(t)).$$

$$\bar{u}_j(t+1) = \bar{u}_j(t) (1 - \theta_{ij}(\bar{u})).$$

under $u_i(t)$, with sign
depending on $s(t) = 1, 0$.
 $j \neq i$, with normalisation.

The corresponding diffusion is optimal when $\alpha > 1$, so a discrete time analysis should give similar behaviour to the reward-inaction rules. Difficulties arise in proving the subregularity of diffusion absorption function w.r.t. the functional equation $\Delta \tilde{\gamma}_{\theta_j}(\tilde{u}) = 0$.

1.7. Comparison between ϵ -optimality and a.s. optimality.

Definition 1.7.1.

A rule is centrally learning at $\underline{c} \in (I^n)^\circ$ if $\exists \delta > 0$ s.t. if we construct the regular n -simplex S , in homogenous co-ordinates with faces $(\tilde{u}_i, -c_i) = \delta_i$; then $\Delta R(\tilde{u}) > 0$ for $\tilde{u} \in S$ (where $\Delta R(\tilde{u}) = \sum q_i \Delta \tilde{u}_i$ as usual).

Definition 1.7.2.

A rule is boundary learning at \tilde{e}_j if $\exists \epsilon_j$ s.t. $\Delta R(\tilde{u}) > 0 \forall \tilde{u}$ s.t. $|\tilde{u} - \tilde{e}_j| < \epsilon_j$.

We first show that ϵ -optimality corresponds to central learning and so we take the diffusion approximation for the 2-action linear rule and verify that we obtain ϵ -opt without consideration of boundaries.

Example 1.7.3.

Let $0 < a < c < b < 1$, then we show $\lim_{\theta \downarrow 0} P(a B b | c) = 0$, where $P(a B b | c) = \text{Prob}(\text{reach } \tilde{u}_1 < a \text{ before } \tilde{u}_2 > b \text{ starting from } c, \text{ with } q_1 > q_2)$.

In continuous time $\phi'' + r(\tilde{u})\phi' = 0$ with $r(\tilde{u}) = 2k \div 2(1 - q_2/q_1)/\theta$

Then subject to $\phi(b) = 0$, $\phi(a) = 1$

We get $\phi(\tilde{u}) = (e^{-2k\tilde{u}_1} - e^{-2k\tilde{u}_2}) / (e^{-2ka} - e^{-2kb})$ and $\phi(c) = (e^{-2k(c-b)} - 1) / (e^{-2k(a-b)} - 1)$.

$\phi(c) < 2e^{-2k(c-b)} / e^{-2k(a-b)}$ for k large and so $\phi(c) < 2 \exp(-2k(c-a))$.

$\lim_{k \uparrow \infty} \phi(c) = 0$ for $c > a$, and we have the result. //

In 1.11.1 we show how this is linked with Wald's Identity, for the discrete time process.

From this example we see that to from an ϵ -optimal rule, we need only find some point c at which the rule is centrally learning. Then with $\bar{u}(0) = c$ we stop the rule at the first action i s.t. $|\bar{u}_i - c_i| \geq \delta$. An n-dim analogy of 1.7.3 would show

$\lim_{\delta \rightarrow 0} \Pr(|\bar{u}_i - c_i| \geq \delta, i \neq \text{opt} \mid \bar{u}(0) = c(0)) = 0$, using the methods of section 1.6. to reduce it to a 2-action problem.

In contrast, we now show that optimality corresponds to boundary learning at e_j $\forall j$ and in 1.7.5. we shall prove that U.L. at boundaries, with $\alpha > 1$, together with inter-boundary communication are sufficient for optimality.

Definition 1.7.4.

- i) The boundary point e_j is stable iff $\Delta \bar{u}_j > 0$ for $|\bar{u}_j - e_j| < \epsilon_j$ and some $\epsilon_j > 0$. (The rule must be boundary learning with drift into boundary.)
- ii) The boundary point e_j is unstable iff $\exists i \neq j$ s.t. $\Delta \bar{u}_i > 0$ $\forall \bar{u}$ s.t. $|\bar{u}_j - e_j| < \epsilon_j$, for some $\epsilon_j > 0$, and $\Delta R(\bar{u}) > 0$. (The rule is boundary learning with drift away from boundary).

In our proof of 1.5.12 we saw that we achieved optimality depending on $\bar{u}(0)$ behaviour at the boundary, and learning could be made as rapid as required in finitely many compact subsets of I . The U.L. property was only used centrally for communication.

- iii) Let the boundary neighbourhood of e_j be E_j . We say that $\bar{u}(0) \in I \setminus \bigcup_j E_j$ communicates with $\{E_j\}$ if a.s. $\exists t > 0$ s.t. $\bar{u}(t) \in E_k$ for some k and $\exists \text{Prob } p_k(\bar{u}(0)) > 0$ that we first reach $\bigcup_j E_j$ at E_k .

Theorem 1.7.5.

For a rule to give a.s. convergence to a stable boundary, it is sufficient that :-

- i) Each boundary is either stable or unstable, and at least one is stable.
- and ii) The rule is optimally boundary learning (R_0) at e_j $\forall j$.

and iii) $\bar{u} \in I \setminus \bigcup_j \mathcal{E}_j$ communicates with $\{\mathcal{E}_i\} \forall \bar{u}$.

Proof.

We prove that as in 1.5.12, we are asymptotically reflected from unstable boundaries, using $\#$ up-crossings $\stackrel{a.s.}{\rightarrow} \infty$ across any rational interval. (see Breiman 1968).

Suppose that we are in $\bar{u}(t) \in \mathcal{E}_i$, unstable; then if we were to have $\lim_{t \rightarrow \infty} \bar{u}_i(t) = 1$ with $\bar{u}_i(t^*) \in \mathcal{E}_i$, $\forall t^* > t$ we could extend the optimal rule to 1 and obtain a contradiction from its sub-optimality. (If $\exists j$ s.t. $\Delta \bar{u}_j > 0$, $\bar{u} \in \mathcal{E}_i$ then $\lim_{t \rightarrow \infty} \bar{u}_j(t) = 0$ if we use optimal rules).

Hence we leave \mathcal{E}_i and next enter \mathcal{E}_k say, and for unstable k we repeat the previous argument for the process is markovian. If \mathcal{E}_k is stable, \exists a finite probability that we never leave. Now apply $\#$ up-crossings $\stackrel{a.s.}{\rightarrow} \infty$ and we see that a.s. we do eventually remain in some stable \mathcal{E}_m for all $t > t_0$.

//

This result actually follows easily from 1.5.12 and 1.6.4., once we notice that the fundamental idea is asymptotic reflection at unstable boundaries, which leaves only the stable boundaries for absorption, if any exist.

Algorithm 1.7.6.

We have n -actions with reward probabilities q_i . To asymptotically converge to the optimum, (and attain its nbd \mathcal{E}_{opt} at t_0 a.s. and remain there for all $t > t_0$) we apply optimal boundary learning to each \mathcal{E}_j and construct nbds \mathcal{E}_j . Set $\bar{u}(0) \in \mathcal{E}_i$ some i , and if $\exists t$ s.t. $\bar{u}(t) \notin \mathcal{E}_i$ set $\bar{u}(t_0) \in \mathcal{E}_k$ with probability $1/n$, $1 \leq k \leq n$.

Proof.

Apply 1.7.5. and note, assuming optimum is unique, that $\Delta \bar{u}_{opt} > 0$ on $\bigcup \mathcal{E}_i$, and communication is immediate. Hence $\lim_{t \rightarrow \infty} \bar{u}_{opt}(t) = 1$. A similar argument follows in the case of multiple optima, when we obtain multiple stable boundaries.

//

The theorem 1.7.5. will be applied in chapter 2 to games between π -cells, in order to prove convergence to pure strategies, and a Nash Point in the case of general sum games. It allows us to investigate the limiting behaviour when the π -cells execute mixed strategies, in attempting to attain the Von Neumann value of the zero-sum game.

Remark 1.7.7.

We can relax the communication with finite probability to all \mathcal{E}_j in 1.7.5., to give a more natural algorithm. Set $\bar{u}(0) \in \mathcal{E}_1$ and if $\exists t$ s.t. $\bar{u}(t) \notin \mathcal{E}_1$, set $\bar{u}(t+1) \in \mathcal{E}_2$ and then visit boundary nbds cyclically until we are absorbed at some stable \mathcal{E}_m .

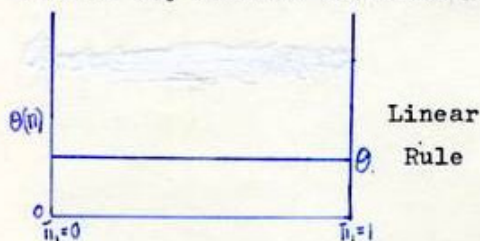
$$(\bar{u}(t^*-1) \in \mathcal{E}_r, \bar{u}(t^*) \notin \mathcal{E}_r \text{ then } \bar{u}(t^*+1) \in \mathcal{E}_{(r+1) \bmod n})$$

The proof of sufficiency for a.s. convergence to stable \mathcal{E}_m follows as in 1.7.5, by asymptotic reflection from unstable \mathcal{E}_r .

1.8. The Family of π -cell Learning Rules.

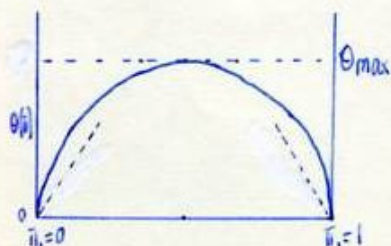
We have now derived a form of probabilistic stability theorem for discrete-time stochastic difference equations, and have shown optimality to be a boundary property.

In this short section, we shall list certain U.L. rules which will be used in chapters 2 and 3, in order that the π -cell may achieve environmental adaptation.



1. $\theta = \text{constant}$, for all \bar{u} .

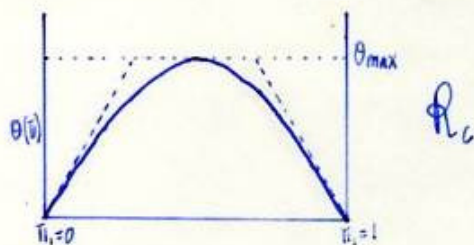
Analysed by Norman (1968).



2. ϵ -optimal rule with $\theta(\bar{u}) = K(\bar{u}_1, \bar{u}_2)^{\alpha-1}$.

$$E(\# \text{alts})^n < \infty, \forall n.$$

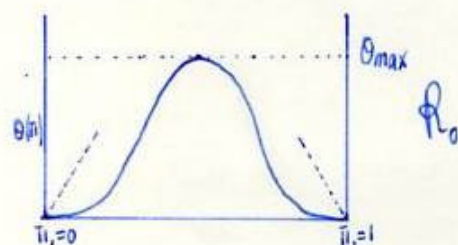
The linear rule is the case $\alpha=0$.



3. Conditionally optimal rule

$$\theta(\bar{u}) = K(\bar{u}_1, \bar{u}_2)$$

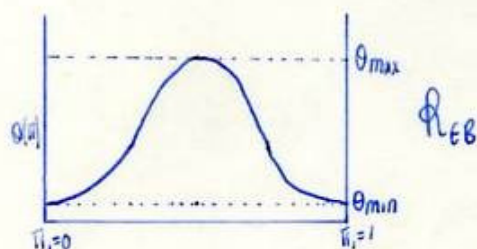
Optimality depends on $q_{u_i}^s$.



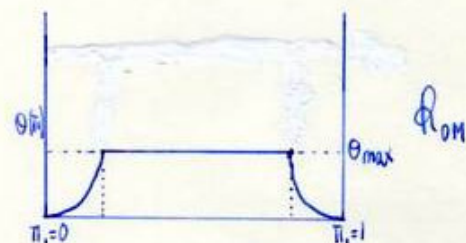
4. Optimal, with boundary learning

$$\theta(\bar{u}) = K(\bar{u}_1, \bar{u}_2)^{\alpha > 1}$$

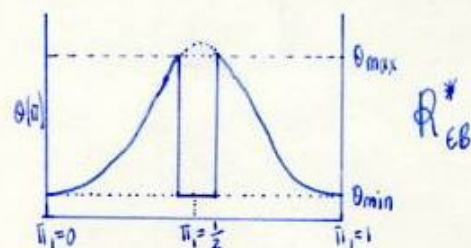
$$\#alts \stackrel{as}{=} \infty$$

5. Boundary learning, yet ϵ -optimal, always converges with $E(\#alts)^n < \infty$.

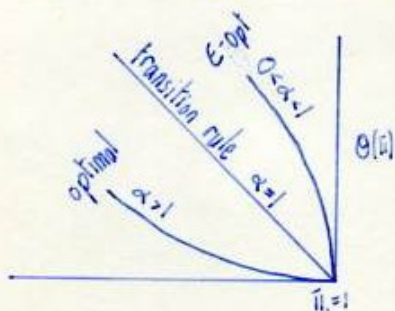
$$\theta(\bar{u}) = K(\bar{u}_1, \bar{u}_2)^{\alpha > 1} + \theta_{min}$$



6. Optimal, with both boundary and central (middle) learning. Centre is used for mixed strategies, and boundary for Nash Points.



7. The extra central "dip" may help in the cellular differentiation of 3.3.7.

8. An expanded view of the boundary, showing the essential differences between R_0 , R_{ϵ} and R_c .

1.9. Time Dependent Stimulus Probabilities.

Now that we have considered \tilde{u} -cells in static media, it is easy to extend the results to certain time varying $q_i^S(t)$.

Theorem 1.9.1.

If $\exists i$ s.t. $\limsup_{t \rightarrow \infty} q_j^{(t)}/q_i(t) < 1 \quad \forall j \neq i$ then

- under i) $\mathcal{R}_0: \tilde{u}_i \xrightarrow{a.s.} 1$
 ii) $\mathcal{R}_\epsilon: \lim_{\theta \downarrow 0} \mathcal{J}(\tilde{u}_i) = 1$

and for 2-actions under

iii) If $\limsup_{t \rightarrow \infty} q_2^{(t)}/q_1(t) = 1 - \delta < 1$ and $0 < \delta$ then $\tilde{u}_2 \xrightarrow{a.s.} 1$.

iv) If $\limsup_{t \rightarrow \infty} q_1^{(t)}/q_2(t) < 1$ and $\liminf_{t \rightarrow \infty} q_1^{(t)}/q_2(t) > \log_e 2$

then for $1 > \theta > 1 - \lim_{t \rightarrow \infty} q_1^{(t)}/q_2(t) / (1 - \log_e 2)$ we have $\lim_{\theta \downarrow 0} \mathcal{J}(\tilde{u}_2) = 1$.

Proof.

We first prove a lemma.

Lemma 1.9.2.

- i) $\lim_{t \rightarrow \infty} \frac{1}{t} \sum_{s=0}^t U^s \mathcal{Y}^i(\tilde{u}) = \mathcal{J}^i(\tilde{u})$ with $\mathcal{Y}^i(\tilde{u}: \tilde{u}_i = 1) = 1, \mathcal{Y}^i(\tilde{u}: \tilde{u}_i = 0) = 0$.
 ii) $\lim_{t \rightarrow \infty} \frac{1}{t} \sum_{s=0}^t U^s \mathcal{J}^i(\tilde{u}) = \mathcal{J}^i(\tilde{u})$ and $\mathcal{J}^i(\tilde{u}: \tilde{u}_i = 1) = 1, \mathcal{J}^i(\tilde{u}: \tilde{u}_i = 0) = 0$.

where U^t is the expectation operator related to $q_i^S(t)$.

Proof.

We just have $\lim_{t \rightarrow \infty} \frac{1}{t} \sum_{s=0}^t U^s \mathcal{Y}^i(\tilde{u}) = E(\mathcal{Y}^i(\tilde{u}(\omega)) | \tilde{u}(0)) = \mathcal{J}^i(\tilde{u})$.

and the result holds with $\mathcal{Y}^i(\tilde{u}) = \mathcal{J}^i(\tilde{u})$. //

Now we stop the process at t' s.t. $q_j^{(t)}/q_i(t) < 1 \quad \forall t > t' \quad \forall j \neq i$.
 and we find sub-regular \mathcal{Y}^i s.t. $U^t \mathcal{Y}^i > \mathcal{Y}^i$ for all $t > t'$

Then $\mathcal{J}^i(\tilde{u}(t')) > \mathcal{Y}^i(\tilde{u}(t'))$

For ii), we consider the result 1.5.4. of Norman (1968), for 2-actions.

It is clear that the result is only dependent on $q_1(t)/q_2(t)$ for we just require s to give $f(\pi, s) = \left(\frac{e^{s(1-\pi)} - 1}{1 - e^{-s\pi}} \right) \frac{\pi}{1-\pi} \geq (\leq) \frac{q_1}{q_2} \quad \forall \pi$.

Hence the same $\gamma_{s,0} = e^{s\pi/0}$ will give us bounds on $\gamma(\pi)$. Also, we may continue to prove $\gamma_{\theta(\pi)}(\pi(t'))$ is monotone and convex, to give a comparison theorem analogous to 1.4.5. Then we get ii) after noting that the result is unaffected by U^t , $0 < t < t'$ when we take the limit $\theta \rightarrow 0$. This gives us that all U.L. rules are at worst ϵ -optimal, under the given $q_i^s(t)$.

For i), under \mathcal{R}_0 , we can still use $\lim_{s \rightarrow 0} \gamma_{s,0}$ to give the result, by choosing $k_t < \frac{1}{\theta} \log \frac{q_1(t)}{q_2(t)}$, $\forall t > t'$ with $\inf_{t > t'} k_t > 0$ by $\lim_{t \rightarrow \infty} \frac{q_1(t)}{q_2(t)} < 1 \quad \forall j \neq i$.

Finally, for iii), with \mathcal{R}_ϵ rules, we again just note that our super (sub) regular γ_k were functions only of q_1/q_2 so that the results extend naturally.

For n -actions in i) and ii), we mimic the extension theorems of section 1.6., for $t > t'$. And in the slow learning limit, the trials $0 < t < t'$ do not affect limiting behaviour, for \mathcal{R}_ϵ . The result for \mathcal{R}_ϵ is ensured by the asymptotic reflection property. //

Remarks 1.9.3.

a) Sawaragi and Baba (1975), defined a variable $C + \delta$ medium which effectively gave ϵ -optimality when $q_1(t) - q_2(t) > \delta \quad \forall t$ for the linear rule. Our 1.9.1. gives the result for all U.L. rules, and we see that the significant factor is $\frac{q_1(t)}{q_2(t)}$. It is 1.9.2. which enables us to prove 1.9.1. quickly, whilst Sawaragi and Baba actually repeated Norman's analysis in full with the new time-dependence.

b) Tsuji et Al (1973) define a non-static environment which generates $q_1(t) > c > q_2(t)$ and this is said to be completely isolated

in the 0^{th} approximation. However, they use a structured 2^L_2 (see 3.3.3.) to learn in the environment, whilst a singleton \bar{u} -cell will give convergence to $\bar{u}_i = 1$ under R_0 . (1.9.1.)

The paper is based on the work of Yasui and Yajima (1970) who consider 2-state, 2-symbol automata and define k^{th} order isolation. We shall find in chapter 3 when a structured automaton is required rather than a single θ_i . The 2^L_2 is suited to a rapidly switching environment, which contrasts greatly with that of Tsuji et Al.

1.10. Skeletons.

We have noted (near 1.5.20.) that the β -rule is not U.L. We now gain further insight into such additive rules by attempting to construct a uniformly learning rule on a grid which spans $I = [0, 1]$.

Definition 1.10.1.

The set $\{\bar{u}_i; i \in \mathbb{Z}\}$ forms a grid if

$$\left. \begin{aligned} \bar{u}_{i+1} &= \bar{u}_i + \theta_i (1 - \bar{u}_i) \\ \text{and } \bar{u}_{i-1} &= \bar{u}_i (1 - \theta_i) \end{aligned} \right\} \forall i \in \mathbb{Z} \text{ (ive and -ve integers).}$$

The set spans I if $\bar{u}_{\infty} = 1, \bar{u}_{-\infty} = 0$.

Theorem 1.10.2.

If a U.L. rule is defined on a grid, then it will not span $I = [0, 1]$.

Proof.

Put $\alpha_n = 1 - \theta_n = \bar{u}_{n-1} / \bar{u}_n = (1 - \bar{u}_{n+1}) / (1 - \bar{u}_n) = \bar{u}_{n+1}^c / \bar{u}_n^c$ by "action reversal".

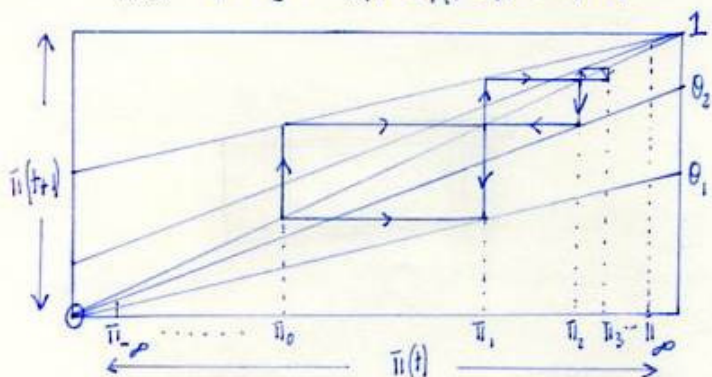
$$\text{Thus } \bar{u}_1 / \bar{u}_{\infty} = \prod_{n=1}^{\infty} \alpha_n = \bar{u}_{\infty}^c / \bar{u}_2^c$$

and if $\bar{u}_{\infty} = 1, \bar{u}_{\infty}^c = 0 \Rightarrow \bar{u}_1 = 0$; absurd.

Hence $\prod_{n=1}^{\infty} \alpha_n > 0$.

Suppose $\bar{u}_n \uparrow c$ then $(\bar{u}_n)^2 = c(1-c) = (\alpha_1 / \alpha_1)^2$ say

and $c = \frac{1}{2} \pm \frac{((1+\alpha)(1-\alpha))^{\frac{1}{2}}}{2(1+\alpha)}$



e.g. $\alpha = \frac{1}{3}$ gives $c(1-c) = \frac{1}{16}$
 $c \approx \frac{1}{15}$

Thus $c \uparrow$ as $\alpha \downarrow 0$ with $\alpha = \alpha_0 = 1 - \theta_0$.

The identity $\alpha_n = \frac{\bar{u}_{n+1}}{\bar{u}_n} = \frac{\bar{u}_{n+1}^c}{\bar{u}_n^c}$ arises by rotating the

diagram above through 180° ; we are reversing actions, with

$$\bar{u} = \text{Prob}(\text{take } u_1), \quad 1 - \bar{u} = \text{Prob}(\text{take } u_2) = \bar{u}^c.$$

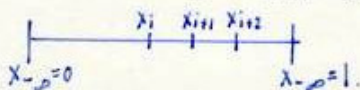
//

This result is easily seen to be closely related to the β -rule being non-U.L. The grid process gives $\left(\frac{\bar{u}_{n+1}^c}{\bar{u}_{n+1}}\right) / \left(\frac{\bar{u}_n^c}{\bar{u}_n}\right) = \frac{1}{\alpha_n \alpha_{n+1}}$,

whilst the β -rule has this ratio independent of n , in order for the additive process to span I .

If we could relate U.L. rules to a countable state space, it would help in proving results. So we shall define a spatial skeleton, just as in continuous time markov processes we may extract a τ -skeleton. This spatial skeleton will be used in an alternative method for proving optimality under \mathcal{R}_0 .

Let $P(x_i \leq x_{i+2} | x_{i+1}) = \text{Prob}(\text{reach } \bar{u} \leq x_i \text{ before } \bar{u} \geq x_{i+2} \text{ starting at } \bar{u} = x_{i+1})$.



Definition 1.10.3.

A skeleton is a set $\{x_i : i \in \mathbb{Z}\}$ which spans $I = [0, 1]$.

The probabilities $P(x_i \leq x_{i+2} | x_{i+1})$ can be calculated approximately, assuming no overshoot, using a diffusion, whilst Wald's Identity is useful when $|x_{i+2} - x_i| < \epsilon$ with ϵ sufficiently small for the learning process to be considered as a form of random walk.

In the skeleton process, we only look at the process at the trials at which the \bar{u} -cell has just crossed a point x_i in the skeleton set.

1.11. Staircases. (An alternative method of proving optimality of R_0).

We first take $[a, b] \subset I$ with $|a - b| < \epsilon$ and $\theta(\bar{u}) < \epsilon^{\delta/2}$ for $\bar{u} \in [a, b]$, so that we can approximate the process as a R.W. and use Wald's Identity to obtain the absorption probabilities. Then we construct a skeleton which resembles a staircase, when projected on function $\theta(\bar{u})$, with successive step lengths of constant ratio β . This β -staircase will approximate and asymptotically converge to the learning function $\theta(\bar{u})$ at boundaries $\bar{u} \in \{0, 1\}$. In the first lemma, we take $\theta(\bar{u}) = \theta$ and approximate the process in $[a, b]$ as a random walk with increments:-

$$\begin{aligned} \bar{u}_i(t+1) &= \bar{u}_i(t) + \theta(1-b) & \text{if } u_1(t) \text{ and } s(t) = 1. \\ \bar{u}_i(t+1) &= \bar{u}_i(t) - b\theta & \text{if } u_2(t) \text{ and } s(t) = 1 \\ \bar{u}_i(t+1) &= \bar{u}_i(t) & \text{if } s(t) = 0. \end{aligned}$$

Lemma 1.11.1.

For the R.W. on $[a, b]$, $f(\bar{u}) \geq (e^{sb/\theta} - e^{s\bar{u}/\theta}) / (e^{sb/\theta} - e^{s(a-d)})$.

where $f(\bar{u}) = \Pr(\text{absorbed in } \bar{u} \leq a \text{ starting from } \bar{u} \in (a, b) \text{ before } \bar{u} \geq b)$.

d_- = max overshoot at $\bar{u} = a$ boundary. And $q_2 > q_1$.

Proof.

Wald's Identity gives $E(e^{\sigma(\bar{u}-\bar{u})} / (A(\sigma))^n) = 1$.

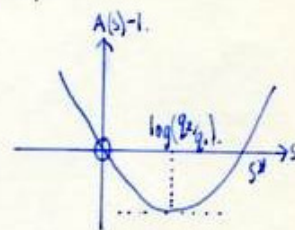
and put $A(\sigma) = 1 = bq_1(e^{(1-b)\theta} - 1) + (1-b)q_2(e^{-\sigma b\theta} - 1) + 1$.

Thus $\sigma = 0$ or $\sigma = s/\theta$ with $s > \log(q_2/q_1) > 0$ for all b .

With $\partial A / \partial \sigma = 0$ if $q_1 e^{\sigma(1-b)\theta} = q_2 e^{-\sigma b\theta}$ or $e^s = q_2/q_1$.

Let $\phi(s)$ be the stopping density; then W.I. gives

$$\int_{-\infty}^a \phi(s) e^{s\bar{u}} ds + \int_b^{\infty} \phi(s) e^{s\bar{u}} ds = e^{s\bar{u}} \quad \text{with } \sigma = s/\theta.$$



Now $q_2 > q_1$ with $s^* > 0$

$$\int_b^\infty \phi e^{s\theta/\theta} d\theta \geq (1 - f(\bar{u})) e^{sb/\theta} \\ < (1 - f(\bar{u})) e^{s(b+d)/\theta}$$

Then $e^{s\bar{u}/\theta} \geq e^{sb/\theta} (1 - f) + f e^{s(a-d)/\theta}$

or $f(\bar{u}) \geq (e^{sb/\theta} - e^{s\bar{u}/\theta}) / (e^{sb/\theta} - e^{s(a-d)/\theta})$

//

Remark 1.11.2.

i) The diffusion equation $\phi'' + r(\bar{u}) \phi' = 0$ gives:-

$$\phi(\bar{u}) = (e^{-kb} - e^{-k\bar{u}}) / (e^{-kb} - e^{-ka})$$

with $k = (q_1 - q_2) / q_1 \theta$

and so in W.I. we would expect $s \sim (q_2 - q_1) / q_1$ as $\theta \downarrow 0$.

i.e. $\theta = e^{\epsilon/\theta}$ and $|a-b| = \epsilon$ with $\epsilon \downarrow 0$.

And for $\theta(\bar{u}) = \theta(\bar{u}, \bar{u}_2)^{\alpha > 1}$ and $\bar{u} \downarrow 0$, the overshoot becomes negligibly small and the diffusion approximation becomes asymptotically, arbitrarily close to the discrete case.

ii) The bounds of Norman, 1.5.4., are just the case with

$[a, b] = [0, 1]$ since the drift is homogenous throughout the interval.

Theorem 1.11.3.

Under $\mathcal{R}_0, \theta'(\bar{u}) = 0$ for $\bar{u} \neq 1$, where the \mathcal{R}_0 family is taken as $\theta(\bar{u}) = \theta(\bar{u}, \bar{u}_2)^{\alpha > 1}$, with $q_1 < q_2$.

Proof.

This is an alternative method of proving 1.5.12. which gives more insight into the process and is easily extended to markovian environments in 1.12.9.

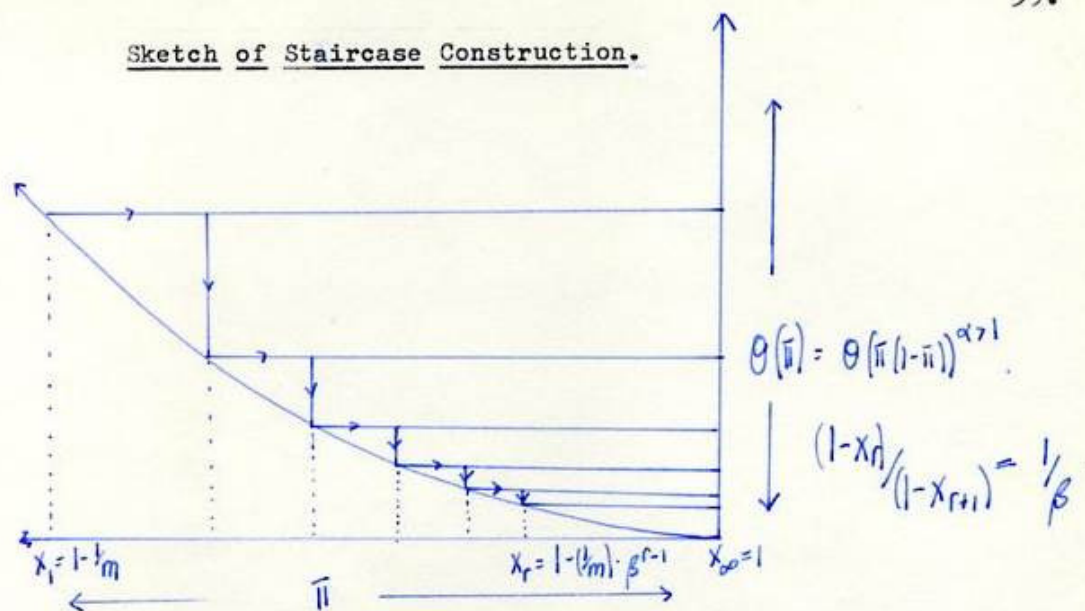
We take $X_1 = 1 - 1/m$ s.t. $m \gg 1$ and so that $\theta(x_1) \sim (1/m)^\alpha$.

Now we define the β -staircase and estimate $\prod_{i=1}^n P(X_{i-1}, \beta X_{i-1} | X_i)$

where $P(aBb|c)$ is defined in 1.10. and

$X_r = 1 - (1/m)^{\beta^{r-1}}$, with $x_2 - x_1 = (1-\beta)/m \gg \theta(x_1)$ on "first-step".

For $\theta(\bar{u}) = K(\bar{u}(1-\bar{u}))^{\alpha-1}$ we require $m \gg (K\beta^\alpha / (1-\beta))^{1/(\alpha-1)}$

Sketch of Staircase Construction.

By 1.11.1. $\Pr(x_{r-1} \leq x_{r+1} | x_r) \sim (e^{sx_{r+1}/\theta} - e^{sx_r/\theta}) / (e^{sx_{r+1}/\theta} - e^{sx_r/\theta})$

and $= (1 - e^{-(\beta/m\theta) \cdot \beta^r}) / (1 - e^{-3\beta^r/m\theta}) = (1 - e^{-\gamma_r}) / (1 - e^{-3\gamma_r})$
 $\sim (1 - e^{-\gamma_r})$

neglecting higher order
terms and overshoot.

with $\gamma_r = \beta^r/m\theta$ and $\theta(x_{r-1}) \sim \theta/m^\alpha (\beta^{r-2})^\alpha$

We use $\theta(x_{r-1})$ since this gives smaller (drift/diffusion) $= r(u) \propto 1/\theta$.

and hence is the worst possible case on $[x_{r-1}, x_{r+1}]$

For $I: [0, 1]$, we had the comparison theorem 1.4.5 to prove

this result, but as the drift for linear θ is homogenous on I ,

the result will hold on arbitrary intervals. We just require

θ small enough so overshoot can be neglected.

Now $\Pr(x_{r-1} \leq x_{r+1} | x_r) \sim 1 - \exp(-sm^{d-1}/\theta \cdot \beta^{2d}/\beta^{r(\alpha-1)}) = 1 - \exp(-\xi_r)$

and $\prod_{r=2}^n \Pr(x_{r-1} \leq x_{r+1} | x_r) \sim \prod_{r=2}^n (1 - \exp(-sm^{d-1}/\theta \cdot \beta^{2d}/\beta^{r(\alpha-1)}))$

and hence $\prod_{r=2}^{\infty} \Pr(x_{r-1} \leq x_{r+1} | x_r) > 0$ iff $\sum_{r=2}^{\infty} \exp(-\xi_r) < \infty$

But by the ratio test we get $\sum \exp(-\xi_r)$ converges extremely rapidly.

Let $a_r = \exp(-\theta_r)$

Then $a_r/a_{r+1} = \exp(-\delta/\beta^{r(\alpha-1)}) / \exp(-\delta/\beta^{(r+1)(\alpha-1)})$ with $\delta = (sm^{\alpha-1}/\theta)\beta^{2\alpha}$
 $= \exp\left(\delta/\beta^{r(\alpha-1)} \left(1/\beta^{\alpha-1} - 1\right)\right)$ and $\alpha > 1 \Rightarrow \beta^{\alpha-1} < 1$ if $0 < \beta < 1$.

Hence for $r \geq r_0$ say $a_r/a_{r+1} > \exp\left(\delta/\beta^{r(\alpha-1)} \left(1/\beta^{\alpha-1} - 1\right)\right)$.

and we have convergence by comparison with geometric series.

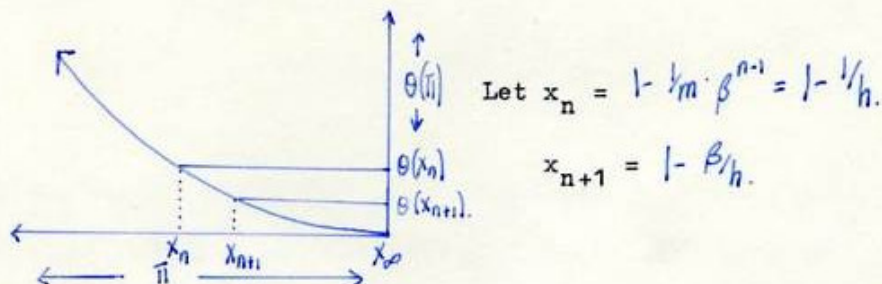
Thus $P_r(x, Bx_{n+1} | x_n) \geq \prod_{r=2}^n P_r(x_{r-1}, Bx_r | x_r) > \xi > 0$ for all $n > 1$.

Remarks 1.11.4.

i) We could use $(1-e^{-\theta_r})/(1-e^{-3\theta_r})$ and still get result easily.

ii) We are assuming m is large enough for us to neglect overshoot, which becomes negligibly small as $n \rightarrow \infty$. If we included it in $P(x_{r-1}, Bx_{r+1} | x_r)$, then its effect would be geometrically small when we take $\sum_r a_r$.

Finally we prove that $\xi=1$ must hold, in the limit $n \rightarrow \infty$.



Lemma 1.11.5.

$$P_r(x_n B | x_{n+1}) \geq 1 - \exp\left(-z\beta^{1/\theta}\right) \quad \text{for } z > 0.$$

Proof.

$$P_{\theta(n)}(x_n B | x_{n+1}) \geq P_{\theta'(n)}(0 B | x_{n+1}) \geq P_{\theta(x_n)}(0 B | x_{n+1}) \quad \text{by 1.4.5.}$$

where $P_{\theta(x_n)}$ denotes probability taken along linear $\theta(x_n)$. Similarly for

And with $\theta'(n) = \theta(n)$ for $n \geq x_n$ and $\theta'(n) = \theta(x_n)$ for $n \leq x_n$. $P_{\theta(n)}$ along non-linear $\theta(n)$.

And $P_{\theta'(n)}$ denotes the corresponding probability along $\theta'(n)$.

Thus $P_{\theta(\frac{1}{n})}(X_n B | X_{n+1}) \geq (1 - \exp(-z(1-X_{n+1})/\theta(\frac{1}{n})))$ by 1.5.4. with
super-regular $e^{-z\bar{u}/\theta}$

and $Pr_{\theta(\frac{1}{n})}(1 B X_n | X_n) \leq \exp(-z(1-X_{n+1})/\theta(\frac{1}{n})).$

//

Now denote $P(X_n B | X_{n+1}) = b_n$.

We start from x_{n+1} and with probability $1-b_n$ we are absorbed.

Else with probability b_n we reach x_n , and with probability $\geq \delta > 0$
we escape to x_1 before x_{n+1} .

We show ultimate escape is certain.

$$\begin{aligned} & \lim_{n \rightarrow \infty} Pr(\text{absorbed in } \bar{u}=1 \text{ starting at } \bar{u}=1-\frac{1}{2h} \text{ before reaching } X_1) \\ &= \lim_{n \rightarrow \infty} Pr(1 B x_1 | x_{n+1}) \leq (1-b_n) + b_n(1-\delta)(1-b_n) + (1-b_n)(b_n(1-\delta))^r \\ &= \lim_{n \rightarrow \infty} (1-b_n) / (1-(1-\delta)b_n). \end{aligned}$$

But $\lim_{n \rightarrow \infty} b_n = 1$ for $\alpha > 1$ whilst for $\alpha = 1$ this fails and we still
need the technique of 1.5.17.

Hence $\lim_{n \rightarrow \infty} Pr(x_1 B 1 | x_{n+1}) = 1$.

Now the process is semi-martingale with $\Delta \bar{u} < 0$ so # upcrossings $\stackrel{a.s.}{\leq} \infty$
and $\bar{u} \rightarrow v \in \{0,1\}$ a.s.

To prove absorption in $\bar{u}=0$ independent of s/mg theorem, we just

reverse the above and show $\lim_{n \rightarrow \infty} P(\bar{u} = \frac{1}{m} B \bar{u}=0 | \bar{u}(0) = \frac{1}{n}) = 0$

for any fixed m , (w.r.t. n) with $m \gg 1$. Then the result follows
since the boundaries communicate, as defined in 1.7.4.

//

Remark 1.11.6.

The reason for constructing a staircase is that for a dynamic
environment we can still obtain an exponentially small drift

towards the sub-optimal boundary, for intervals $[x_r, x_{r+1}]$ with

$(-\beta$ fixed, yet arbitrarily small. (The optimal boundary will

be shown to maximize the average reward.) Then with β fixed

and $x_n \uparrow 1$, we have $\theta(x_{r-1}) \ll |x_{r-1} - x_{r+1}|$ for large r and we

can use the theory developed by Miller (1962), and subsequently

by Keilson and Wishart (1964), for processes defined on markov chains.

Intuitively, letting $1-\beta$ be arbitrarily small yet $\beta < 1$ gives $\theta(n) \sim$ constant on each "step" and we can ensure that the process asymptotically becomes that of a R.W., defined on a markov chain, on each "step" as we approach $\bar{\pi} = 1$.

1.12. Dynamic Environments.

We now take $M(\Delta_{2\beta}, q_{u_i}^{x,s})$ as the environment, with $e\Delta = e \triangle$ equilibrium vector, and we show that the limiting behaviour of the reinforcement process is only dependent on the ratios $\{e q_i^{x,s} / e q_j^{x,s}\}$. The distribution over environment states is denoted by $\omega(t)$.

Theorem 1.12.1.

$$i) \lim_{n \rightarrow \infty} U^n \gamma^i(\bar{\pi}(0), \omega(0)) = \gamma^i(\bar{\pi}, \omega) \quad \text{where } \gamma^i(\bar{\pi}_i = 1, \omega) = 1, \gamma^i(\bar{\pi}_i = 0, \omega) = 0$$

$$ii) U \gamma^i(\bar{\pi}, \omega) = \gamma^i(\bar{\pi}, \omega) \quad \text{where } \gamma^i(\bar{\pi}_i = 1, \omega) = 1, \gamma^i(\bar{\pi}_i = 0, \omega) = 0.$$

with U defined as usual:- $U \gamma^i(\bar{\pi}(t), \omega(t)) = E(\gamma^i(\bar{\pi}(t+1), \omega(t+1)) | \bar{\pi}(t), \omega(t)).$

Proof.

$$\lim_{n \rightarrow \infty} U^n \gamma^i(\bar{\pi}, \omega) = 1 \cdot \gamma^i(\bar{\pi}, \omega) + 0 \cdot (1 - \gamma^i) = \gamma^i(\bar{\pi}, \omega)$$

$$\text{and } \gamma^i(\bar{\pi}, \omega) = \gamma^i = \lim_{n \rightarrow \infty} U^n \gamma^i = \lim_{n \rightarrow \infty} U(U^{n-1} \gamma^i) = U \lim_{n \rightarrow \infty} U^{n-1} \gamma^i = U \gamma^i.$$

where we can interchange U operator and limit, since $U^n \gamma^i$ converges uniformly to γ^i . //

We now actually write out the U operator explicitly, and so consider the short term behaviour of the process.

Let $\phi(\bar{\pi}, \omega)$ be linear in ω_α , $\phi = \sum_{\alpha=1}^{m_2} \phi_\alpha \omega_\alpha$ and put $\Phi = (\phi_\alpha)$.

Lemma 1.12.2.

$$U \Phi(\bar{\pi}) = \sum_{i=1}^s \bar{\pi}_i (Q_i^s \Delta \Phi(T_i^s \bar{\pi}))$$

with $T_i^0 \bar{\pi} = \bar{\pi}$

$$T_i^1 \bar{\pi} = (\bar{\pi}_1, \dots, \bar{\pi}_{i-1}, \bar{\pi}_i, \bar{\pi}_{i+1}, \dots, \bar{\pi}_n)$$

With T_i defined as:- $T_i \bar{\pi}_i = \bar{\pi}_i (1 - \theta(\bar{\pi}))$ and $T_i \bar{\pi}_i = \bar{\pi}_i + \theta(\bar{\pi}) (1 - \bar{\pi}_i).$

$$Q_i^s = \begin{pmatrix} q_i^{1,s} & & & 0 \\ & q_i^{2,s} & & \\ & & \ddots & \\ 0 & & & q_i^{m,s} \end{pmatrix}$$

Proof.

$$U\phi(\underline{i}, \underline{\omega}) = \sum_i \bar{u}_i ((\underline{\omega}, p_i) \phi(\underline{i}, L_i^s \underline{\omega}) + (\underline{\omega}, q_i) \phi(T_i^s \underline{i}, L_i^{s-1} \underline{\omega})) \\ = \sum_i \bar{u}_i ((\underline{\omega}, q_i^s) \phi(T_i^s \underline{i}, L_i^s \underline{\omega})).$$

with $(L_i^s \underline{\omega})_\alpha = \sum_\beta \omega_\beta q_{\beta i}^{s-1} A_{\beta \alpha} / (\sum_\beta \omega_\beta q_{\beta i}^{s-1})$

Now linearity gives $L_i^s \underline{\omega} = L_i^{s-1} \underline{\omega} / q_i^{s-1} \underline{\omega} = \sum_\beta L_{\alpha \beta}^{s-1} \omega_\beta / q_i^{s-1} \underline{\omega}$

Thus $\bar{\Phi}(\underline{i}) \rightarrow \sum_i \bar{u}_i (p_i \Delta \bar{\Phi} + q_i \Delta \bar{\Phi}(T_i^s \underline{i}))$
 $\rightarrow \sum_i \bar{u}_i L_i^{s-1} \bar{\Phi}(T_i^s \underline{i}).$

and $U\bar{\Phi} = \sum_{i, \alpha} \bar{u}_i (q_i^{s-1} \Delta_{\alpha \beta} \phi_\beta(T_i^s \underline{i}))$ as required. //

Now if we fix the increments $T_i^s \underline{i}$ we have the process as

a R.W. defined on the environmental markov chain. So for 2-actions, mimicing 1.11.1. we define $M(s) \triangleq bQ_1 \Delta (e^{s(1-b)} - 1) + (1-b)Q_2 (e^{-sb} - 1) + \Delta - I$ on interval $[a, b]$ with $|a - b| = \epsilon$ say.

Also with $\phi = \sum \omega_\alpha \phi_\alpha$ and $\bar{\Phi} = (\phi_\alpha)$ we let

$$W\bar{\Phi}(\underline{i}) \triangleq b(p_1 \Delta \bar{\Phi} + q_1 \Delta \bar{\Phi}(\underline{i} + \theta(1-b))) + (1-b)(p_2 \Delta \bar{\Phi} + q_2 \Delta \bar{\Phi}(\underline{i} - \theta b)).$$

and we try to solve $W\bar{\Phi} = \bar{\Phi}$ using $\bar{\Phi} = R e^{s\bar{\Phi}/\theta}$.

This gives $M(s) R = 0$ or that $|M(s)| = 0$ gives the relevant solutions.

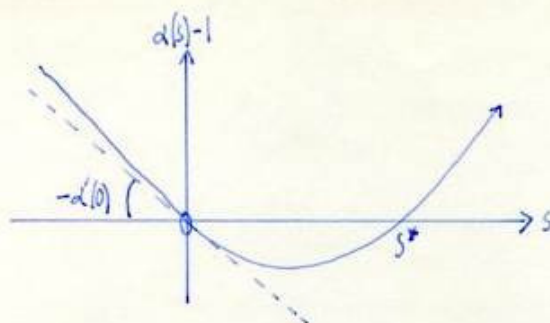
So we have obtained the correspondence between $M(s)$ and the approximation W to the U operator.

Actually $M(s) = P(-s) - I$ where $(P(s))_{ij} = \Delta_{ij} m_{ij}(s)$ with θ divided out for convenience, and with $m_{ij}(s) = \text{m.g.f. for R.W. increments if environment makes transition } i \rightarrow j$.

Miller (1962) uses $P(s)$ as the fundamental matrix and on writing

$(P(s) - \alpha(s)I) = 0$ to define $\alpha(s)$, shows that the drift constant, corresponding to $A(s) = 1$ in scalar case, is given by $\alpha(-s) = 1$ for markov processes.

We have $-m_{ij}^*(0)$ is the mean i, j , increment and $m_{jk}(t) = \int_{-\infty}^{\infty} e^{-tx} dg_{jk}(x)$ generally, where $dg_{jk}(x)$ gives the increment distribution.



The Solutions of

$$\underline{\alpha(s) = 1.}$$

Lemma 1.12.3.

a) $M(0) = \Delta - I.$

b) $\left. \frac{d}{ds} |M(s)| \right|_0 < 0$ iff $\underline{e} \cdot \underline{q}_1^{s=1} > \underline{e} \cdot \underline{q}_2^{s=1}.$

c) $\alpha'(0) > 0$ iff $\left. \frac{d}{ds} |M(s)| \right|_0 > 0.$

Proof.

a) is just observation, whilst for b) we have:-

$$\left. \frac{d}{ds} |M(s)| \right|_0 = \text{tr} [\text{adj}(M(s)) \frac{dM}{ds}] \Big|_0 = b(1-b) \text{tr} [\text{adj}(\Delta - I) N \Delta]$$

where $N = Q_1 - Q_2.$

Thus $\left. \frac{d}{ds} |M(s)| \right|_0 < 0$ iff $\prod_{j=1}^n (\lambda_j - 1) \text{tr} [S_j \eta_j' N] < 0.$

where $S_j \eta_j' = S_j$ the first spectral matrix.

Thus $-\text{tr} [S_1 N] < 0$ iff $\underline{e} \cdot \underline{q}_1^{s=1} > \underline{e} \cdot \underline{q}_2^{s=1}$

For c) $\left. \frac{d}{dt} |P(t) - \alpha(t)I| \right|_0 = \text{tr} [\text{adj}(P(0) - \alpha(0)I) (P'(0) - \alpha'(0)I)]$

and so $\text{tr} [S_1 (P'(0) - \alpha'(0)I)] = 0$ and $\alpha'(0) = \sum_{i,j} \underline{e}_i \Delta_{ij} m'_{ij}(0).$

$$= -b(1-b) \text{tr} [S_1 N] = \text{tr} [S_1 P'(0)].$$

So $\alpha'(0) = -b(1-b) \sum \underline{e}_i (\underline{q}_1^{s=1} - \underline{q}_2^{s=1}) = -\mathcal{E} \text{ drift}$ and hence c)

This is essentially the same calculation as for b) with

$P(s) - I$ replacing $M(s)$. Note also $P'(0) = M'(0)$ and $\alpha(s)=1 \Leftrightarrow |M(s)|=0.$

//

By transferring from U operator to $M(s)$ techniques, we have a close approximation to large iterates of U, when τ is kept within the small interval $[a, b]$. If we let $\theta(p) = O(\epsilon^{C\tau})$ when $|a-b| = \epsilon$ then as $\epsilon \downarrow 0$, this approximation becomes exact.

On using R_0 rules with the β -staircase, we have precisely these conditions holding as $\pi \rightarrow \pi \in \{0,1\}$ so we need only consider $\alpha(0)=1$ and $\text{sign}(\alpha'(0))$ to determine the limiting behaviour as in Miller (1962).

We shall now state a central limit theorem of Keilson and Wishart (1964) which gives us the asymptotic independence of environment and π -increment process.

Theorem 1.12.4. (Keilson and Wishart)

If $F_r(x, k) = P_r(R(k)=r, X(k) \leq x)$ then

$$\lim_{k \rightarrow \infty} F(x\sqrt{k} + km, k) = \Phi(x\sigma^{-1})e$$

where the environment has equilibrium vector e and states r .

The mean drift $m = -\alpha'(0)$ and variance $\sigma^2 = \alpha''(0) - [\alpha'(0)]^2$

which are expressed in terms of $P(s)$ matrix, Keilson and Wishart (1967).

In their notation, k is the discrete time and X is the value of the increment process variable.

Proof.

Keilson and Wishart (1964).

//

This result allows us to prove that as $\theta \downarrow 0$ we asymptotically follow the mean drift.

Lemma 1.12.5.

i). If $\theta \downarrow 0$ and $k \rightarrow \infty$ s.t. $\theta k = K = \text{constant}$, then

$$\lim_{k \rightarrow \infty} P_r(|X(k) - km| > y) \rightarrow 0 \quad \text{for all } y > 0$$

ii) This limit remains true as $km = Z \downarrow 0$ if $\theta = O(Z^{c+1})$, $y = O(Z)$, $c = \text{const.}$, and with $\theta = O(y)$ and so $K \downarrow 0$ also.

Proof.

$$P_r(X(k) - km > y) = P_r((X(k) - km)/\sqrt{k} > y/\sqrt{k}) \sim \int_{y/\sqrt{k}}^{\infty} \frac{1}{\sqrt{2\pi}} e^{-z^2/2} dz \quad \text{asymptotically.}$$

Now put $m = m'\theta$ and $km = Z = km'\theta$ with $k\theta = K = Z/m'$

Also $\sigma^2 = O(\theta^4)$ for $\theta \downarrow 0$ can be verified, say $\sigma = \theta\sigma'$

Then $y/\sigma\sqrt{k} = y/\sigma\sqrt{k} \sqrt{\theta} = y/\sigma\sqrt{k\theta}.$

Hence $\lim_{\substack{k \rightarrow \infty \\ \theta \rightarrow 0}} \Pr(X(k) - km > y) = 0$ for $y > 0.$

and result follows by a similar argument for $\Pr(X(k) - km < -y).$

Also if $y = o(Z)$ and $\theta = o(Z^{c'})$ then $\lim_{Z \downarrow 0} y/\sqrt{k\theta} = \infty.$

and we have justification for the remark after 1.12.3. We asymptotically follow the mean drift even as the interval width $\downarrow 0$, so long as $\theta \downarrow 0$ more rapidly, as defined above. Z is effectively the distance travelled by the increment process. //

We now just need to assert that $\exists s^* < 0$ s.t.

$f(\eta) = (e^{s^*b/\theta} - e^{s^*\eta/\theta}) / (e^{s^*b/\theta} - e^{s^*(a-b)/\theta})$ in the notation of 1.11.1.

Such an s^* is related to the real non-zero root of $|M(s)| = 0.$

Lemma 1.12.6.

$\exists s^* < 0$ s.t. $|M(s)| = 0 \Rightarrow s < s^* < 0$ for all b when $\alpha'(0) > 0.$

Proof.

Consider the 2nd derivative of $|M(s)|$. This has an upper bound for all b by observation. Then by Taylor's Theorem, all solutions to $|M(s)| = 0$ with $s \neq 0$ must be bounded strictly away from the origin by some $s^* < 0$. Similarly when $\alpha'(0) < 0$ and $e.q_1^1 > e.q_2^1$ we obtain $s^* > 0$, with $f(\eta)$ re-defined for "b" absorption. //

Freedman (1973) proves a useful pair of inequalities which give us the same asymptotic mean drift following and his results also hold for finite stopping times. However it has not been possible to actually construct absorption probabilities in the manner of Miller (1962) as this would require Wiener-Hopf factorizations. But since we are not concerned with environment states at absorption, such a full analysis is unnecessary and it is sufficient to just apply 1.12.3. \rightarrow 1.12.6. for reinforcement in a dynamic environment.

Theorem 1.12.7.

- a) If $\theta \downarrow 0$ with $\theta = O(\epsilon^{c-1})$ when $|b-a| = O(\epsilon)$, $c = \text{const.}$
 then $P(a \leq b | c = (a+b)/2) = O(e^{-s(b-a)/\theta})$ with $s > s^*$ for $a < b$
 and s^* as in 1.12.6. where $\alpha'(0) < 0$.
- b) $\lim_{\theta \downarrow 0} \gamma(\bar{u}_i) = 1$, $\bar{u}_i \neq 0$ iff $\underline{e} \cdot \underline{q}_1^{s=1} > \underline{e} \cdot \underline{q}_2^{s=1}$.

Proof.

The bound on s follows from 1.12.6. as $|M(s)|=0 \Leftrightarrow \alpha(s)=1$.

The drift towards the sub-optimal boundary becomes exponentially small by 1.12.4. and 1.12.5., since asymptotically, as $\theta \downarrow 0$ the mean dynamic process becomes indistinguishable from a static process. Miller's paper (1962) is also relevant, except that he considered just one absorbing barrier. However, in this case we obtain $P(\text{absorbed at the barrier}) \sim O(e^{-y^2/\theta})$ where $\alpha(s)=1$ and $y \propto$ distance from barrier, with drift away from it, and $y/\theta \gg 1$, $s > 0$.

For b) we apply a) in each interval $[a, b]$, and then piecing intervals together as in the staircase theorem.

Thus $P(1 \leq \bar{u} | \bar{u} > 0 : \bar{u}_i \in (0, 1)) \geq \lim_{n \rightarrow \infty} \prod_{i=1}^n (1 - e^{-sn^{c-1}}) = 1$.

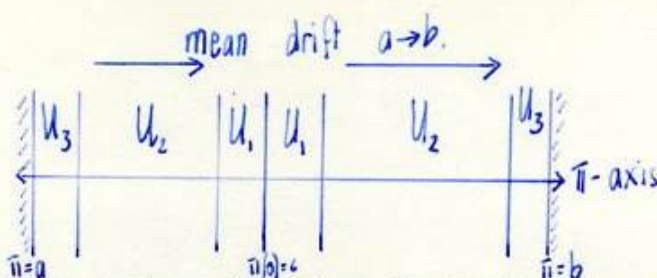
when $(b-a) = 1/n$ and $\theta = (1/n)^{c-1}$.

//

Here we are interested only in the asymptotic behaviour when the environment and \bar{u} -process become independent as $\theta \downarrow 0$, so we ignore the short term fluctuations of 1.12.2.

Note that in 1.12.7. a) that we will actually have $P(a \leq b | c)$ exponentially small for all $c \in [a, b]$ as $\theta \downarrow 0$ and the absorption probabilities asymptotically reduce to the form obtained for the static medium in 1.11.1. This is required in our β -staircase construction for dynamic environments.

Remark 1.12.8.



Diagrammatically we have short term fluctuations in region U_1 , whilst in U_2 we follow mean drift with arbitrarily high probability as $\theta \downarrow 0$. Then in U_3 we could consider the environmental state for the full analysis, whereas we just require $P(a \leq b | c)$ for our reinforcement rules.

Theorem 1.12.9.

Under R_0 , $\pi_i \rightarrow 1$ iff $\lim_{n \rightarrow \infty} q_i > \lim_{n \rightarrow \infty} q_j \quad \forall j \neq i$ with $\pi_i(b) > 0$.

Proof.

First the result for 2-actions follows from the β -staircase construction 1.11.3., since the approximation for $P(x_{i-1} \leq x_{i+1} | x_i)$ still holds with drift constant $s > 0 \quad \forall b \in [0, 1]$.

So $\lim_{n \rightarrow \infty} P(x_1 \leq x_{n+1} | x_n) > \delta > 0$ holds with $1-\beta$ arbitrarily small since 1.12.5. ii) is found to be satisfied, and also 1.12.7. a) holds for each "step" as we approach any boundary.

Then as before we get $\lim_{n \rightarrow \infty} P(x_i \leq 1 | x_n) = 1$, when drift is away from the boundary, and hence the result.

For n -actions we just use the $\theta_{ij}(\bar{\pi})$ rules to give asymptotic reflection from sub-optimal boundaries, using a comparison argument analogous to 1.6.4. to take us from n -actions to a 2-action model, when we can then apply the staircase result above. //

Remark 1.12.10.

- i) A π -cell learning under R_0 asymptotically takes that action which maximizes the average payoff in a dynamic $\pi(\Delta_{2\beta}, q_{a_i}^{n,s})$.
- ii) We use 1.12.4. and 1.12.9. in chapter 3. There we find that a network of π -cells asymptotically maximizes its payoff at the next trial w.r.t. the equilibrium environmental distribution in its present state.

1.13. Learning Barriers.

We know from 1.7. that R_b is boundary learning and R_c is centrally learning, yet in this section we show that there is still a unification in the underlying learning mechanism.

Definition 1.13.1.

A learning barrier $a, b \in I$ s.t. $|a-b| = \epsilon$ and $\theta(\bar{u}) = O(\epsilon^{c>1})$, $\bar{u} \in [a, b]$.

This section will consist of numbered remarks relating to this concept.

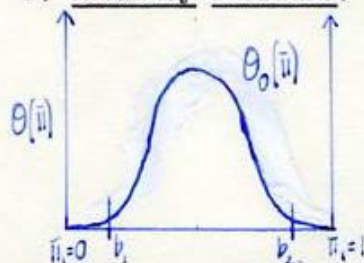
1.13.2. For boundary learning we put a barrier around each

$\bar{u}_i = 0$, whilst for central rules we place a barrier away from all boundaries. From our previous theory, when the equilibrium drift is from $a \rightarrow b$, $P(aBb | c = \frac{1}{2}(a+b)) < e^{-Z(b-a)/\theta}$, for $z > 0$, and independent of θ and $[a, b]$. This follows from weak convergence to a diffusion as $\theta \downarrow 0$ and 1.12.5., 1.12.7. for dynamic \bar{u} .

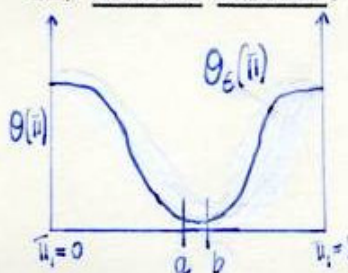
For optimality we have a learning barrier in any arbitrarily small interval of all boundaries, and this provides the mechanism for the asymptotic reflection from sub-optimal boundaries.

1.13.3. The barrier is a potential step and acts as a form of "diode", in that we can only pass through the barrier with high probability if we are travelling with the equilibrium drift. The effective "field strength" is $d\theta(\bar{u})/d\bar{u} \hat{x} = E$ and at barriers placed centrally or at sub-optimal boundaries we have $|E| \uparrow \infty$.

i) Boundary Barriers.



ii) Central Barrier.

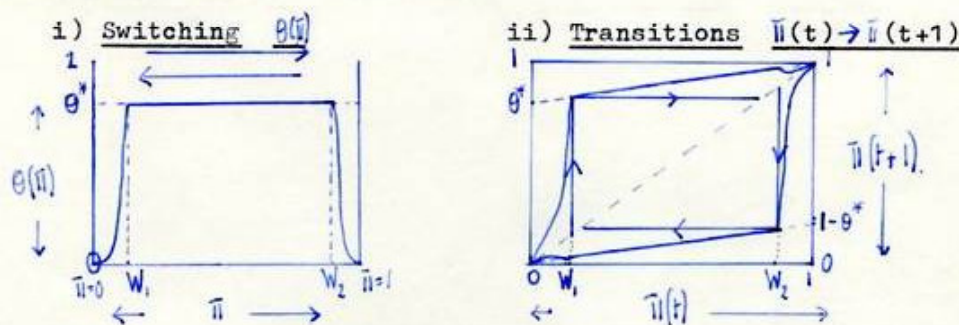


1.13.4. We could put $\theta_b(\bar{u}) + \theta_c(\bar{u}) = \text{constant}$, to see that R_b and R_c may be viewed as complementary ways of learning, yet they both operate using the learning barrier mechanism.

1.13.5. For optimal boundary learning rules, we effectively just have a deterministic stability theorem to satisfy. So we could actually reformulate the theory in terms of control theoretic terminology to give switching between attractors through a catastrophe until we find a stable attractor.

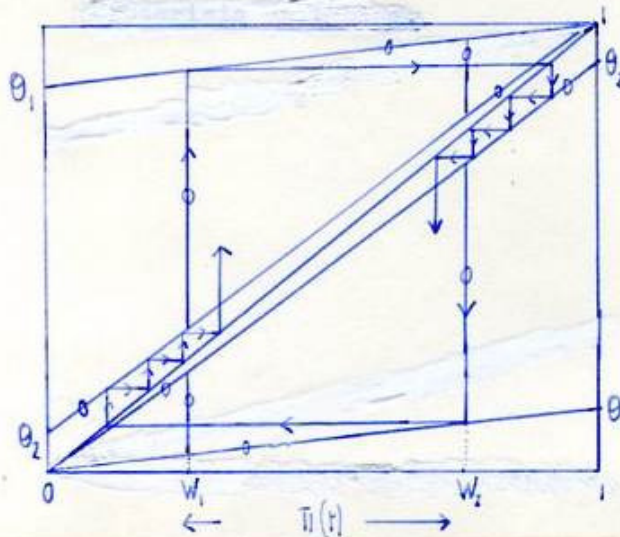
When we consider networks of $\bar{\Pi}$ -cells, this would entail an automaton increasing its environmental adaptation through structural catastrophes as in the work of Thom (1975). The automaton could then increase both its memory depth and number of actions used. (see chapter 3).

1.13.6. The advantage of stochastic automata is that we have an explicit mechanism for incorporating environmental information into the structure. Thus a $\bar{\Pi}$ -cell incorporates both a learning mechanism and boundary switching.



We can choose θ^* for the switching using the transition diagram. Now, if W_1 and W_2 are arbitrarily close to their respective boundaries the $\bar{\Pi}$ -cell just has the boundary learning barriers which are stable only if they give maximum average payoff.

1.13.7. We can also easily define $\theta(\pi)$ to give a hysteresis effect in its switching, still retaining the U.L. property, thus mimicing the switching of the cusp catastrophe. This effect is closely related to the grid 1.10.1. , and with $W_1 + W_2 = 1$, we obtain an "overshoot" on boundary switching if and only if $W_2 < 1/(2-\theta^*)$.



Hysteresis Switching.

Use θ_1 and θ_2 with

θ_1 near 1 and θ_2 near 0,
s.t. $\theta_2 < 1 - \theta_1$.

This gives R_{ES} learning rule.

We have:-

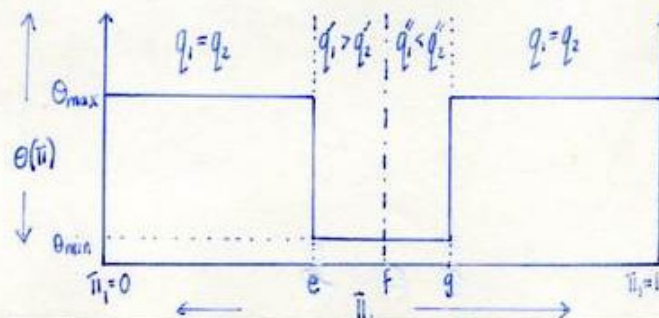
a) $\theta(\pi) = \theta_1, \pi \in [W_1, W_2]$.

b) $\theta(\pi) = \theta_2, \pi \in [0, W_1] \cup (W_2, 1]$.

—○— denotes transition line $\pi(t) \rightarrow \pi(t+1)$.

1.13.8. In π -cell games and π -cell networks, we find that the effective stimulus probabilities are dependent on the action distribution, so that $q = q(\pi)$. This can give rise to a potential barrier between strictly stable boundaries. Then, π -cell switching behaviour between stable limiting structures can be likened to the "tunnelling" of charged particles in quantum theory. In 2.4.5. we briefly examine one of the simplest cases of 2 strictly stable boundaries, which represent Nash Points of π -cell games.

A potential barrier is formed on placing two learning barriers, (potential steps) "back to back" as shown in a particular $q(\pi)$ process below.



Potential Barrier.

$$LB_1 = \{\pi, : \pi \in [e, f]\}$$

$$LB_2 = \{\pi, : \pi \in [f, g]\}$$

where LB_i denotes learning barrier.

$$PB = \{\pi, : \pi \in [e, g]\}$$

1.13.9. Our reinforcement rules R_0 depend only on the boundary learning barriers for their optimality, which could perhaps be thought of as stochastic attractors and repellers. This concept of learning barriers and boundary switching is important when we discuss structured automata, and is the basis for the main characterization theorems.



Chapter Two

Games between Automata



Chapter 2.

There is no remembrance of former things, nor will there
be any remembrance of later things yet to happen among
those who come after.

Ecclesiastes 1 v 11.

2. Games between Unstructured Automata.

2.1. The Model for \mathbb{N} -cell Games.

In the previous chapter we considered the \mathbb{N} -cell playing against an environment \mathcal{M} . It is now natural to formulate a game between \mathbb{N} -cells so that for a particular \mathbb{N} -cell \mathcal{A}_i the environment is the remaining $\{\mathcal{A}_{j \neq i}\}$. These games were introduced by Chandrasekaran and Shen (1968) and developed to a limited extent by Narendra and Viswanathan (1974), who first demonstrated, through computer simulation, the possibility of cyclic behaviour.

We shall first develop the theory for zero-sum 2-automata games to indicate the main features, before considering general sum n -automata games.

Definitions 2.1.1.

- i) Let g_{ij} = expected winnings to automaton 1 taking u_i when automaton 2 takes u_j .
and $\mathbb{N}_i^1 = \Pr(\text{automaton 1 takes } u_i)$.
ii) Let $p_{ij} = \Pr(\text{ under actions } (i,j), \text{ automaton 1 receives penalty } -1).$
 $= \Pr(\text{ " " " " 2 " reward } +1).$

$$\text{with } q_{ij} + p_{ij} = 1 \text{ and } g_{ij} = q_{ij} - p_{ij} = 1 - 2p_{ij}.$$

- iii) Let $p_i^1 = \Pr(\text{automaton 1 receives } s=0, \text{ using } u_i) = \sum_j p_{ij} \mathbb{N}_j^2$.
 $= \frac{1}{2}(1 - \sum_j g_{ij} \mathbb{N}_j^2).$
and similarly $p_j^2 = \sum_k \mathbb{N}_k^1 (1 - p_{kj}) = \frac{1}{2}(1 + \sum_k \mathbb{N}_k^1 g_{kj}).$

The above definitions hold for an arbitrary automaton, so now we restrict the study to \mathbb{N} -cells, \mathcal{A}_i , with environment determined by $\mathcal{P}_j(\{\mathbb{N}_{k \neq j}^1\})$. The game is zero-sum in expectation rather than deterministically.

2.2. Pure Saddles.

We shall use the usual game theoretic terminology and first consider the case of a pure saddle point. Using the optimal boundary

learning theory of 1.7. , the following result is almost immediate.

Theorem 2.2.1.

For 2 Π -cells acting under R_0 and if

$$\max_m \min_n g_{mn} = \min_n \max_m g_{mn}$$

then $\lim_{t \rightarrow \infty} \Pi^1_i(t) = 1$ and $\lim_{t \rightarrow \infty} \Pi^2_j(t) = 1$ only if (i,j) is a unique pure saddle of g_{ij} .

Proof.

We assume that we have a strict saddle, in that $g_{rj} < g_{ij} < g_{ik}$ with $r \neq i$, $k \neq j$, and hence that the pure saddle is unique. However, the proof is easily modified to take care of non-uniqueness, as in 1.6.4. and also in the general Nash Theorem 2.5.1.

As before we use $\Delta \Pi^k_i(t) = \mathcal{E}(\Pi^k_i(t+1) | \Pi^{\alpha}(t) \text{ for each } \alpha) - \Pi^k_i(t)$

to express the conditional increment.

$$\text{Then } \Delta \Pi^1_i(t) = \Pi^1_i(t)/2 \sum_{j,k} \theta_{ij}(\Pi^1(t)) \Pi^1_j(t) (g_{ik} - g_{jk}) \Pi^2_k(t). \quad 1)$$

$$\text{and } \Delta \Pi^2_j(t) = \Pi^2_j(t)/2 \sum_{i,k} \theta_{ij}(\Pi^2(t)) \Pi^2_i(t) (g_{ki} - g_{kj}) \Pi^1_k(t). \quad 2)$$

a) We shall first prove convergence. The process can only be absorbed at a boundary $\Pi^k_i \in \{0,1\} \forall i,k$ where $\theta_{ij}(\Pi^k) = 0$. Now we apply 1.7.5. and assume w.l.o.g. that $(1,1)$ is the saddle.

In a small neighbourhood N_{11} of $(1,1)$, using 1) we obtain:-

$$\Delta \Pi^1_i = \frac{1}{2} \Pi^1_i \sum_j \theta_{ij}(\Pi^1) \Pi^1_j (g_{11} - g_{j1}) \Pi^2_1 + O\left(\max_{k \neq 1} \sum_j \theta_{ij}(\Pi^1) \Pi^1_j (g_{jk} - g_{j1}) \Pi^2_k\right). \quad 3)$$

Then we take ϵ s.t. $\Pi^1_i > 1-\epsilon$ and $\Pi^2_j > 1-\epsilon$ } gives $\Delta \Pi^1_i > 0$ and $\Delta \Pi^2_j > 0$ in N_{11} .

Since $\theta_{ij}(\Pi^k) > 0$ away from boundaries, we are either absorbed at $(i,j) \neq (1,1)$, or else we eventually enter N_{11} where Π^1_i and Π^2_j are sub-martingales. We modify the argument of 1.7.4. so that the process is stopped if either Π^1_i or Π^2_j leaves N_{11} .

If the process were not convergent, we would enter N_{11} i.o. , but for a s/mg, # upcrossings $< \infty$, and we get a contradiction.

Hence if the process cannot be absorbed at any $(i,j) \neq (1,1)$, we converge to $(1,1)$.

b) Suppose $\lim_{t \rightarrow \infty} \bar{u}_{j \neq 1}^1 = 1$ and $\lim_{t \rightarrow \infty} \bar{u}_k^1 = 1$ with $(j,k) \neq \text{saddle}$.

Then either $\exists n$ s.t. $\varepsilon_{nk} > \varepsilon_{jk}$ or m s.t. $\varepsilon_{jm} < \varepsilon_{jk}$.

Assume w.l.o.g. that $\exists \varepsilon_{nk} > \varepsilon_{jk}$ then in neighbourhood N_{jk} of boundary point (j,k) we have:-

$$\Delta \bar{u}_n^1 = \frac{1}{2} \bar{u}_n^1 \theta_{nj}(\bar{u}) \bar{u}_j^1 (g_{nk} - g_{jk}) \bar{u}_k^1 + O(\max_{m \neq k} \sum \theta_{nj} \bar{u}_j^1 \bar{u}_m^1 (g_{nm} - g_{jm})). \quad (4)$$

and take $\bar{u}_j > 1 - \varepsilon'$ s.t. $\Delta \bar{u}_n^1 > 0$ in N_{jk} .

$\bar{u}_k^1 > 1 - \varepsilon'$

We have $q_i(\bar{u}(t)) = (1 + g_{ik})/2 + O(\varepsilon')$ in N_{jk} , and $q_n^1(\bar{u}(t)) > q_j^1(\bar{u}(t))$.

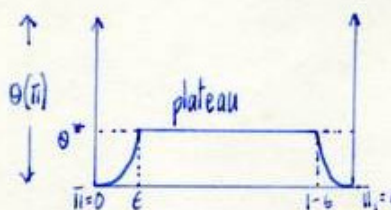
But the learning function $\theta_{ij}(\bar{u}^k)$ is optimal and hence boundary N_{jk} is probabilistically reflecting by 1.7.5., since all boundaries communicate.

Hence by a) $\lim_{t \rightarrow \infty} \bar{u}_i^1 = 1$ and $\lim_{t \rightarrow \infty} \bar{u}_j^1 = 1$ only if (i,j) is a pure saddle. //

For \bar{u} -cell games it is useful to define boundary learning rules with $\theta_{ij}(\bar{u}) = \theta = \text{constant}$ away from boundaries, giving a centrally learning plateau for mixed strategy trajectories.

Consider the case of 2-actions and R. rule $\theta(\bar{u})$.

Define R. rule s.t. $\theta(\bar{u}) = \theta(\bar{u}_1, \bar{u}_2)^{\alpha}$ for $\bar{u}_i < \varepsilon$ else



$$\theta(\bar{u}) = \theta^* = \text{constant}$$

so $\theta^* = \theta(\varepsilon(1-\varepsilon))^{\alpha}$

Similarly for n-actions with $\theta_{ij}(\bar{u})$

We then obtain $\Delta \bar{u}_i^1 = \theta^* \bar{u}_i^{1/2} (\sum_j g_{ij} \bar{u}_j^1 - V(\bar{u}^1))$ and $\Delta \bar{u}_j^1 = -\theta^* \bar{u}_j^{1/2} (\sum_i \bar{u}_i^1 g_{ij} - V(\bar{u}^1))$ away from boundaries.

where $V(\bar{u}) = \sum_{i,j} \bar{u}_i^1 g_{ij} \bar{u}_j^1$ = value of the game.

2.3. Mixed Strategies.

The analysis of games without pure saddles is both more difficult and yet more interesting, in that analogies arise naturally with population processes. Indeed, we can now reformulate such processes as games between species in an attempt to give insight to operating mechanisms. Biologists have recently been searching for such a formulation, as in the paper of Maynard-Smith (1973).

Also from the Hardy-Weinberg equations of mathematical genetics we have $p' = (p^2 w_{11} + 2pq w_{12}) / W$ where p = frequency of genotype a_1

$$q = 1 - p = \quad \quad \quad a_2$$

$$w_{ij} = \text{selective viability of } a_i a_j$$

$$\text{and } W = p^2 w_{11} + 2pq w_{12} + q^2 w_{22}.$$

$$\text{and } p = p' \text{ when } u = p'/q' = (w_{12} - w_{22}) / (w_{12} - w_{11}) = p/q.$$

This equilibrium point is precisely that obtained if the process is viewed as a game with game matrix w_{ij} . Although here each "player" is constrained to have the same distribution over the genotype strategies a_1 and a_2 .

In this thesis we shall restrict ourselves to the theoretical n -cell framework. We shall examine 2×2 games, although, as in chapter 1, many results can be extended to n -actions. The next lemma indicates the reason for cyclic behaviour, under any U.L. rule.

Lemma 2.3.1.

If there is no pure saddle, then $\Delta \bar{u}_i^1 = \Delta \bar{u}_j^2 = 0$ iff

$$(\bar{u}_i^1, \bar{u}_j^2) = \lambda_j^1$$

is the Von Neumann value of the game.

Proof.

$$\Delta \bar{u}_i^1 = \Delta \bar{u}_j^2 = 0 \quad \forall i \Rightarrow V(\bar{u}) = \sum_i \bar{u}_i^1 g_{ir} = \sum_k g_{jk} \bar{u}_k^2, \text{ but this has}$$

5.)

solution λ_j^1 , the V-N optimal mixed strategy.

We find $\lambda_1^1 = 1 / (1 + (q_{11} - q_{12}) / (q_{22} - q_{21}))$ and $\lambda_1^2 = (q_{21} - q_{11}) / (q_{22} - q_{11})$

$$\lambda_2^1 = 1 / (1 + (q_{21} - q_{11}) / (q_{22} - q_{21}))$$

with $\lambda_j^i \in (0, 1)$

//

We shall see that the λ_j^i also generate the deterministic trajectories for n-actions for a linear U.L. rule, where λ_j^i satisfy the equations 5), of 2.3.1. In practice, given an arbitrary g_{ij} , we should use Θ_{ij} with ϵ sufficiently small to give λ_j^i lying on the central plateau of $\Theta_{ij}(\bar{u})$. The $\Theta_{ij}(\bar{u})$ were chosen to give boundary optimality at pure saddles, whilst $\Theta(\bar{u})$ rules give the homogenous central learning for mixed strategies. We shall consider the form of deterministic trajectories around and we find a naturally arising diffusion of a constant. Such processes are analysed by Barbour (1973).

For the 2x2 game we use $\bar{u}_1 = x$, $\bar{u}_2 = y$, $\underline{w} = (\Delta x, \Delta y)$, $\underline{\lambda} = (\lambda_x, \lambda_y)$ with $\lambda_x = \lambda_1^1$, $\lambda_y = \lambda_1^2$.

Lemma 2.3.2.

For $\underline{\lambda} \in [0, 1]^2$ and $\underline{x} = (x, y) \in [0, 1]^2$ then $(\underline{x} - \underline{\lambda}) \cdot \underline{w} = 0$
iff $x = \lambda_x$ or $y = \lambda_y$ or $x = y$ or $x + y = 1$.

Proof.

$$\left. \begin{aligned} \Delta x &= \Theta_x x (q_{11} y + q_{12} (1-y)) - V(\underline{x}) \\ \Delta y &= \Theta_y y (V(\underline{x}) - x q_{11} - (1-x) q_{21}) \end{aligned} \right\} \text{ with } \begin{aligned} \Theta_x &= \Theta_x^\alpha (1-x)^\alpha \\ \Theta_y &= \Theta_y^\alpha (1-y)^\alpha \end{aligned}$$

and putting in λ_x and λ_y we obtain:-

$$\begin{aligned} (x - \lambda_x) \Delta x &= \beta \Theta_x x (1-x) \\ (y - \lambda_y) \Delta y &= -\beta \Theta_y y (1-y) \end{aligned}$$

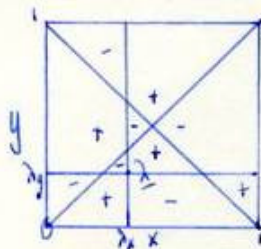
where $\beta(x, y) = (x - \lambda_x)(y - \lambda_y)((q_{11} + q_{22}) - (q_{12} + q_{21})) = \gamma(x - \lambda_x)(y - \lambda_y)$

$$(\underline{x} - \underline{\lambda}) \cdot \underline{w} = \beta (x^{\alpha+1} (1-x)^{\alpha+1} - y^{\alpha+1} (1-y)^{\alpha+1}) = \Theta \beta (x - y) (1 - (x + y)) (f(x, y))$$

with $f(x, y) > 0$ for $(x, y) \in [0, 1]^2$.

Now if $q_{11} + q_{22} = q_{12} + q_{21}$ we get $\lambda_{\lambda_x} = \lambda_{\lambda_y} = 0$.

Hence the only roots are those stated.



It is interesting to note that on the boundary ∂I^+ , symmetry gives us $\text{sgn}((x-\lambda) \cdot \underline{w})$ is +ve precisely half its length.

We also have $V(x) - V_{opt} = \beta$.

and $\nabla V = \sigma((y-\lambda_y), (x-\lambda_x))$

so $\beta = (V_x V_y / \sigma)$

Then $(x-\lambda) \cdot \underline{w} = \sigma(V_x V_y / \sigma) (x-y)(1-(x+y)) f(x, y)$.

//

With our intuition strengthened by the previous lemma, we shall consider simple properties of the deterministic trajectories.

Theorem 2.3.3.

If deterministically we set

$$\begin{aligned} dx/dt &= \sigma \theta_x (y - \lambda_y) x (1-x) = \theta_x V_x x (1-x) \\ dy/dt &= -\sigma \theta_y (x - \lambda_x) y (1-y) = -\theta_y V_y y (1-y) \end{aligned}$$

with $\theta_x = \theta_x x^{\alpha} (1-x)^{\alpha}$ and $k = \alpha + 1$

Then:- a) The trajectories are $\Phi^* = f(x) f(y) = \text{constant}$

where $\log f(x) = \int (x - \lambda_x) / (x(1-x))^k \theta_x dx$

b) The trajectories are periodic, and also if

$$(2\lambda - 1)^2 > (2k - 1)/k \quad \text{for } \lambda = \lambda_x \text{ and } \lambda = \lambda_y$$

then all Φ^* are convex.

c) $\int_{\Phi^*} V(x) dt / \int_{\Phi^*} dt = \overline{V(x)} = V_{opt}$

Proof.

a) Suppose $\exists \Phi(x, y)$ s.t. $d\Phi(x, y) = 0$ expressing "energy conservation".

On \underline{w} trajectories $d\Phi = \Phi_x \dot{x} + \Phi_y \dot{y} = 0$

and so $\Phi_x (y - \lambda_y) / \theta_y y (1-y) = \Phi_y (x - \lambda_x) / \theta_x x (1-x)$

and put $\Phi_y = (y - \lambda_y) / \theta_y^k (1-y)^k$

and similarly for x .

Then let $\Phi^* = \exp \Phi = f(x)f(y) = \text{constant}$, as required

If we put $\theta_x^* = \theta_x^{k+1} (1-x)^{k+1}$ then $\dot{x} = \theta_x^* V_x$
 $\dot{y} = -\theta_y^* V_y$

and $f(x) = \exp \int V_y / \theta_x^* dx$

In particular for $k=1$ $f(x) = x^{\lambda_x} (1-x)^{1-\lambda_x}$

and $k=2$ $\log f(x) = \lambda_x/x + (1-\lambda_x)/(1-x) + (1-2\lambda_x) \log(x/(1-x))$

b) For periodicity the result follows as in Volterra's equations as $f(x)$ is bounded. Goel et Al (1971) give a proof.

For convexity we examine 2nd derivatives.

$\dot{y} = -(x - \lambda_x) / (y - \lambda_y) \cdot (y(1-y) / x(1-x))^k$
 and $\dot{y}'' = (x - \lambda_x)^2 / (y - \lambda_y) \cdot (y(1-y) / x^2(1-x)^2)^k (h(x) + h(y))$

where $h(x) = g(x)(x(1-x))^k / (x - \lambda_x)^2$

and $g(x) = (x - \lambda_x)k(1-2x) - x(1-x)$

Now $g(x) > 0$ for some x iff $(1-2\lambda_x)^2 < (2k-1)/k$

Hence for $g(x) < 0$, $g(y) < 0$, \ddot{y} has the appropriate sign, under given condition.

In particular for $k=1$ we have convex Φ^* , whilst for $k=3$, say, all Φ^* are convex if $0.067 < \lambda_i < 0.933$.

Since ω is tangential to Φ^* , we have expected outward drift, stochastically, for convex Φ^* , for all $x \in I^2$.

c) We have $\int_{\Phi^*} \Delta \pi_i^1 = 0$ using our old notation.

Thus $\int_{\Phi^*} \Delta \pi_i^1 / \pi_i^1 \theta(\pi_i^1) = - \int_0^1 (\sum_j \pi_j^1 g_{ij} \dot{\pi}_j^1 - \sum_k g_{ik} \pi_k^1) dt = 0$

and hence $\sum_k g_{ik} \pi_k^1 = \sum_j \pi_j^1 g_{ij} \dot{\pi}_j^1 = \bar{V} = \sum_k \pi_k^1 g_{ik}$

But solutions to this are uniquely $\bar{\pi}_j^i = \lambda_j^i$, where $\bar{\pi}_j^i$ denotes time-averaged π_j^i on the deterministic, continuous time trajectories.

$$\text{Hence } \bar{V} = \sum_{i,j} \lambda_j^i g_{ij} \lambda_j^i = V_{opt}.$$

//

Corollary 2.3.4.

For n-actions and k=1, the trajectories are given by:-

$$\text{a) } \Phi^* = \prod_{i,j} (\pi_j^i)^{\lambda_j^i} = \text{constant}.$$

If the optimal strategy is completely mixed, $\lambda_j^i \in (0,1) \forall i,j$, then:-

$$\text{b) } \lim_{T \rightarrow \infty} \int_0^T V(\bar{\pi}^k) dt / \int_0^T dt = V_{opt}$$

with $\Phi^* = c$.

Proof.

a) We verify that $d\Phi^* = 0$ on the trajectories.

$$d\Phi = \sum \lambda_j^i d\pi_j^i / \pi_j^i \quad \Phi = \log \Phi^*$$

$$\Delta \pi_j^i = -\frac{1}{2} \pi_j^i \left(\sum_{i,j} \pi_j^i g_{ij} \pi_j^i - \sum_k g_{jk} \pi_k^i \right)$$

$$\Delta \pi_j^i = \frac{1}{2} \pi_j^i \left(\sum_k \pi_k^i g_{kj} - \sum_{i,j} \pi_j^i g_{ij} \pi_j^i \right)$$

Hence deterministically

$$d\Phi = -\frac{1}{2} \sum_j \lambda_j^i \left(V(\bar{\pi}^k) - \sum_k g_{jk} \pi_k^i \right) + \frac{1}{2} \sum_j \lambda_j^i \left(V(\bar{\pi}^k) - \sum_k \pi_k^i g_{kj} \right) = 0$$

since $\sum \lambda_j^i = \sum \lambda_j^i = 1$ and by definition $\sum \lambda_j^i g_{ij} = \sum g_{kj} \lambda_j^i = V_{opt}$.

b) If $\lambda_j^i \in (0,1) \forall i,j$ then $\Phi^* = \text{const}$ will be closed hypersurfaces around λ . Again see Goel et Al (1971) for the treatment of similar trajectories, using analogies with statistical mechanics.

$$\text{Now } \frac{1}{T} \int_0^T 2d\pi_j^i / \pi_j^i dt = K/T (\log \pi_j^i(T) - \log \pi_j^i(0)) \quad K = \text{const}$$

$$= \frac{1}{T} \int_0^T \left(\sum_k g_{jk} \pi_k^i - V(\bar{\pi}^k) \right) dt$$

But for $\Phi^* = c$, $(\log \pi_j^i(T) - \log \pi_j^i(0)) / T \rightarrow 0$ as $T \rightarrow \infty$

Thus by the same argument as before $\bar{\pi}_j^i = \lambda_j^i$

and $\bar{V}(\bar{\pi}^k) = V_{opt}$ where time average is taken as $T \rightarrow \infty$.

Lemma 2.3.5.

For small oscillations about λ with $\theta_x = \theta_x^\alpha (1-x)^\alpha$

a) $\int_{\Phi^*} dt \sim 2\pi / \theta \left(\prod_i \lambda_i \right)^{k/2}$

b) We have ellipses to first order in $\epsilon_x = (x - \lambda_x)$, $\epsilon_y = (y - \lambda_y)$

with ratio (major/minor) axis $\sim \left(\prod_i \lambda_i^1 / \prod_i \lambda_i^2 \right)^{k/2}$

Proof.

$d\epsilon_x/dt = \theta'_x \epsilon_y$ with $\theta'_x = \theta \lambda_x^k (1-\lambda_x)^k$

and similarly for y, giving $\epsilon_x'' = -\theta'_x \theta'_y \epsilon_x$

So period $= 2\pi / (\theta'_x \theta'_y)^{1/2}$ and $\epsilon_x = A \cos((\theta'_x \theta'_y)^{1/2} t)$

with $A/B = (\theta'_x / \theta'_y)^{1/2}$ $\epsilon_y = B \sin((\theta'_x \theta'_y)^{1/2} t)$

For $\lambda_x = \lambda_y$ or $\lambda_x = 1 - \lambda_y$ we have circles.

We also see that $\int_{\Phi^*} dt \propto 1/\theta$ for all trajectories and k. //

The pursual of the analysis of mixed strategy games leads to many difficulties, unless we are to take "large population approximations" about the $\Phi = \text{const.}$ However, I shall sketch why I believe superfluous strategies may vanish deterministically.

Sketch 2.3.6.

Let $\bar{\pi}_j = \lambda_j^i$ be optimal mixed strategies with $\lambda_s^r = 0$ for some r,s including, say, $\lambda_1^1 = 0$. Then if $\bar{\pi}_s^r < \epsilon$ $\forall r,s$ s.t. $\lambda_s^r = 0$ then for large T we obtain $\frac{1}{T} \int_{\text{trajectory}} 2d\bar{\pi}_i^1 / \bar{\pi}_i^1 \theta = -(\bar{V}(\bar{\pi}) - \sum_j g_{ij} \bar{\pi}_j^1) < 0$.

By definition, a superfluous strategy has $\sum_j g_{ij} \bar{\pi}_j^1 < V$ at the saddle and $\bar{\pi}_j^1 \neq \lambda_j^i$ $\bar{V}(\bar{\pi}) \neq V_{opt}$ near $\Phi^* = \text{const.}$ Thus $\bar{\pi}_i^1(T) < \bar{\pi}_i^1(0)$ and $\bar{\pi}_i^1 = 0$ is stable. A more rigorous treatment may show this naive argument to be false. //

Although we may begin with $\Phi^* = \bar{\pi}(\bar{\pi}_j^1)^{(\lambda_j^i)^*}$ as trajectories with some i,j, s.t. $(\lambda_j^i)^* \notin [0,1]$ we hope asymptotically to obtain the optimal

(λ_j^i) of the reduced square matrix. It seems plausible that \bar{u} -cells may automatically proceed through an algorithm of the type described by Karlin (1959 p50-51), in their attempt to maximize payoff at each trial.

It is precisely the U.L. property which enables \bar{u} -cells to discriminate between arbitrarily small differences in payoff (as $\theta \downarrow 0$) and hence give such natural trajectories around λ_{opt} .

Lemma 2.3.7.

In a pure saddle with 2-strategies, at least one of $\bar{u}_1^i(t)$ and $\bar{u}_2^i(t)$ are semi-martingales, throughout their domain.

Proof.

We have $\lambda \notin I^2$ so at least one of $(y - \lambda y)$ or $(x - \lambda x)$ has constant sign. //

Corollary 2.3.8.

Under R_ϵ with saddle (i, j) s.t. $\lambda_j^2 \notin [0, 1]$.
then $\lim_{\theta \downarrow 0} \lim_{t \rightarrow \infty} \bar{u}_i^1(t) = 1$.

Proof.

This result is due to the domination of strategy i.

$\lambda_j^2 \notin [0, 1] \Rightarrow \Delta \bar{u}_i^1(t) > 0$ for all \bar{u}^k away from the boundary
and now apply 1.4.5. for the result. //

This corollary indicates the type of difficulty we encounter for $n > 2$ actions, using R_ϵ when there are usually no dominating strategies. Boundary learning is then essential to achieve optimality.

Lemma 2.3.9.

Under $R_\epsilon \exists (i, j)$ s.t. $\lim_{t \rightarrow \infty} (\bar{u}_i^1, \bar{u}_j^1) = (1, 1)$

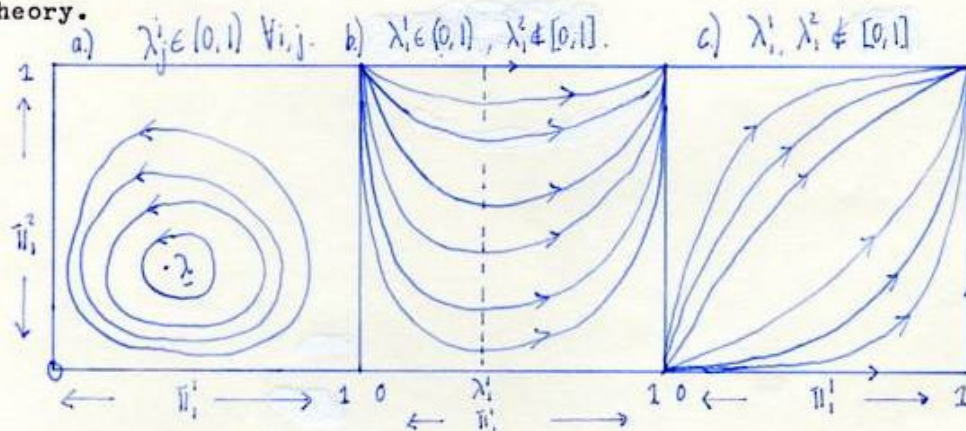
Proof.

This is an application of 1.5.7. since $i(\bar{u}_i^k) > 0$ so eventually each \bar{u}_k will take some strategy infinitely often. //

We know that 2.3.9. is false for optimal rules by 2.2.1. , but this does not prevent $\lim_{t \rightarrow \infty} \Phi^{\lambda^*} = 0$, so that we cycle arbitrarily close to the boundary. We could prove this by showing that Φ^{λ^*} is a super-martingale, which certainly seems to be true for convex Φ^{λ^*} . However, a rigorous treatment would involve much manipulation, so we shall leave it as a conjecture.

Whilst the mixed strategy behaviour gives insight into the operation of U.L. rules, the Π -cells are more suitable when optimal strategies are pure. To preserve cyclic trajectories as a deterministic structure we could use networks of Π -cells, which are considered in the next chapter, or the structured automata of Krinskii (1963, 1966).

We can actually avoid dealing with superfluous strategies by choosing a class of games which admits a unique completely mixed saddle, such as the Minkowski-Leontief form of economic theory.



The basic forms of $\Phi^{\lambda^*} = \text{constant}$ for zero-sum games are sketched above. It is to be emphasised that the Π -cell achieves optimal deterministic time-averaged behaviour, without any knowledge of the opponents strategy; only if its own strategy is successful. For convergence to λ Brown (see Karlin (1959)) required both players to possess complete knowledge of the past, and even then convergence is by no means easy to prove. Note that the theory in this section is easily extended to cover the cases with

$\lambda_j^i \in \{0,1\}$ for some i,j , which occurs when $g_{rt} = g_{mn}$ for some $(r,t) \neq (m,n)$.

2.4. General Sum n -cell Games.

Definition 2.4.1.

i) Let g_{ij}^k = expected reward for k^{th} n -cell if \otimes_i uses u_i
 \otimes_j uses u_j
 as before $g_{ij}^k = q_{ij}^k - p_{ij}^k$ and $q_i^1 = \frac{1}{2}(1 + \sum_j g_{ij}^1 g_{ij}^2)$.

ii) A Nash Point is defined as (i, j) s.t. $g_{ij}^1 \geq g_{mj}^1 \quad m \neq i$
 $g_{ij}^2 \geq g_{in}^2 \quad n \neq j$

We proceed as before, omitting the convergence proof until we treat n -automata games, with the most general Nash convergence theorem.

Lemma 2.4.2.

For 2×2 matrix games the deterministic, continuous time trajectories are given by $\Phi^* = \prod_{i,j} (\pi_{ij})^{\mu_j^i}$ = constant, for linear U.L. rules.

where $\mu_j^i = (-1)^i \lambda_j^i / \gamma^i$ $\gamma^i = \theta_i (\sum_j q_{ij}^i - \sum_{j \neq k} g_{jk}^i)$.

and $\sum_j \lambda_j^i g_{ji}^i = V^i$ $\sum_j g_{ij}^i \lambda_j^i = V^i \quad \forall i$.

Proof.

Put $\Delta x = \Delta \pi_1^1 = \gamma^x x(1-x)(y - \lambda_y)$

$\Delta y = \Delta \pi_1^2 = \gamma^y y(1-y)(x - \lambda_x)$

Thus find $\bar{\Phi}$ s.t. $d\bar{\Phi}$ on trajectories.

$\bar{\Phi}_x \dot{x} + \bar{\Phi}_y \dot{y} = 0$ if $\bar{\Phi}_x = (x - \lambda_x) / \gamma^x x(1-x)$

$\bar{\Phi}_y = -(y - \lambda_y) / \gamma^y y(1-y)$

or $\bar{\Phi} = -\frac{1}{2} \gamma^x (\lambda_x \log x + (1 - \lambda_x) \log(1 - x)) + \frac{1}{2} \gamma^y (\lambda_y \log y + (1 - \lambda_y) \log(1 - y))$

and hence result. //

Corollary 2.4.3.

Under non-linear U.L. rules $\theta(\pi) = \theta(\pi)^\alpha (1 - \pi)^\alpha$
 $\Phi^* = f(x)/f(y)$ = constant with $f(x) = \exp \left(\int (\lambda - \lambda_x) dx / \gamma^x (x(1-x))^{k-1} \right)$.

Proof.

We have $\Phi_x = (x - \lambda_x) / \sigma^x (\lambda(1-x))^k$ and similarly for Φ_y

and now integrate.

For zero-sum games $\sigma^x + \sigma^y = 0$ so we obtain $f(x) f(y) = \text{const}$
as in 2.3.3. a). //

Lemma 2.4.4.

For 2x2 games Φ^* is periodic iff:-

a) $\lambda_j^i \in (0,1) \forall i,j$ and $(\sigma^1 > 0, \text{ iff } \sigma^2 < 0)$.

or b) there are no Nash points.

Proof.

We shall assume a linear U.L. rule so that we can see the form of Φ^* explicitly. However the result is true for all U.L. rules.

a) If $\lambda_j^i \in (0,1)$ and $(\sigma^1 > 0 \text{ iff } \sigma^2 < 0)$ then
$$\Phi^* = (x^{\lambda_x} (1-x)^{1-\lambda_x})^{1/\sigma^x} (y^{\lambda_y} (1-y)^{1-\lambda_y})^{1/\sigma^y} = \text{const.}$$

or $g(x)h(y) = \text{constant}$ and is periodic since

$g(x), h(y)$ are bounded on $[0,1]$, assuming $\sigma^x > 0$ w.l.o.g.

If $\sigma^x > 0$ and $\sigma^y > 0$ say, then either $g(x)$ or $h(y)$ will be unbounded, giving hyperbolic trajectories.

Similarly if $\lambda_j^i \notin (0,1)$ for some i,j .

Now using a) we prove b).

Suppose \exists a Nash point at $(1,1)$ $g_{21}^1 < g_{11}^1$, $g_{11}^2 > g_{12}^2$

Then $\sigma^1 > 0, \sigma^2 > 0$ if $\lambda_j^i \in (0,1)$ and we have a contradiction.

And if Φ^* is periodic we get $g_{11}^1 < g_{21}^1$ as contradiction. //

Corollary 2.4.5.

i) If \exists 2 Nash Points then $\Phi^* = c$ forms a saddle.

ii) If \exists only one Nash Point then $\exists \lambda_j^i \notin (0,1)$.

iii) $\int_{\Phi^*=c} V^k(\tilde{u}^k; \tilde{v}^k) dt / \int_{\Phi^*=c} dt = V^k$, the equilibrium winnings.

Proof.

i) We take (1,1) and (2,2) as Nash points w.l.o.g.

Then $\varepsilon_{21}^1 < \varepsilon_{11}^1$, $\varepsilon_{11}^2 > \varepsilon_{12}^2$, $\varepsilon_{22}^1 > \varepsilon_{12}^1$, $\varepsilon_{22}^2 > \varepsilon_{21}^2$.

and hence $\delta^1 > 0$, $\delta^2 > 0$ and $\lambda_j \in (0,1)$

We get the result for (1,2) and (2,1) by symmetry.

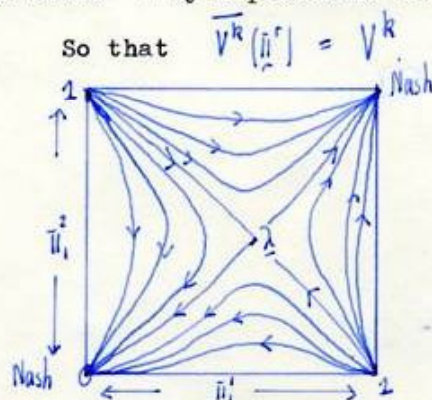
Now construct the trajectories around λ as in the sketch below.

ii) This is immediate. Again, as in zero-sum games, we have

either a) $\lambda_j \notin (0,1) \forall i,j$ with $\bar{\pi}_1^1, \bar{\pi}_2^2$ s/mg, giving process at least ε -optimal.

or b) $\lambda_j \in (0,1), \lambda_i \notin (0,1)$ say, and again a result analogous to 2.3.8 will hold.

iii) Just index V by superfix k in 2.3.3. c) for the result.



All $\bar{\pi}^k \in \{0,1\}$
and are essentially
hyperbolic around the
saddle.

This is the only other basic form we need to consider. Again the theory is easily extended to cover the special cases of

$\lambda_j \in \{0,1\}$ for some i,j .

//

2.5. n-Automata Games.

Definition 2.5.1.

i) Let g_{β}^k = expected reward to \otimes_k if \otimes_i takes β_i .

ii) Let $\bar{\pi}_{\beta}^k = \prod_{i \neq k} \Pr(\otimes_i \text{ takes } \beta_i)$.

iii) We put $g_{\beta}^k = g_{\beta}^k - p_{\beta}^k$ and to emphasise β_k we shall

use g_{β}^k if $\beta_k = i$, similarly for r_{β}^k . (Note, β_k is put only as a left suffix in this notation with $\bar{\ell} = (\beta_1, \dots, \beta_{k-1}, \beta_{k+1}, \dots, \beta_n)$.)

$$\text{iv) } p_i^k = \frac{1}{2} \left(1 - \sum_{\beta \in \mathcal{B}} g_{\beta}^k \bar{u}_{\beta}^k \right) \quad \text{as } \sum_{\beta \in \mathcal{B}} \bar{u}_{\beta}^k = 1. \quad \forall r.$$

$$\text{Then } \Delta \bar{u}_i^k = -\bar{u}_i^k / 2 \sum_{j \in \mathcal{B}} \bar{u}_j^k \theta_{ij}(\bar{u}^k) (j g_{\beta}^k - i g_{\beta}^k) \bar{u}_{\beta}^k.$$

$$\text{with } \bar{u}_i^k \text{ as before and } \Delta \bar{u}_i^k(t) = \mathbb{E}(\bar{u}_i^k(t+1) | \bar{u}^k(t), \forall r) - \bar{u}_i^k(t).$$

Theorem 2.5.2.

If \exists a Nash Point then we converge a.s. to some Nash β_* under \mathcal{R}_0 .

Proof.

Again using our boundary learning theory, the result is just a generalization of 2.2.1.

Assume w.l.o.g. that $\bar{u}_1^k = \frac{1}{2}$ is Nash.

$$\text{so } j g_{\beta}^k = \frac{1}{2} \leq i g_{\beta}^k \quad \forall j, k.$$

a) We first prove convergence by examining the probabilistic boundary stability. In a small ϵ -neighbourhood of $\frac{1}{2}$ we have:-

$$\Delta \bar{u}_1^k = \frac{1}{2} \bar{u}_1^k \sum_{j \in \mathcal{B}} \theta_{1j}(\bar{u}^k) \bar{u}_j^k (j g_{\beta}^k - i g_{\beta}^k) \bar{u}_{\beta}^k + O(\max_{\beta \in \mathcal{B}} \sum_{j \in \mathcal{B}} \theta_{1j}(\bar{u}^k) \bar{u}_j^k (i g_{\beta}^k - j g_{\beta}^k) \bar{u}_{\beta}^k).$$

and we can take ϵ s.t. for $\bar{u}_1^k > 1-\epsilon \quad \forall k$, $\Delta \bar{u}_1^k > 0$.

We effectively have a boundary learning process with time dependent g_{β}^k .

We "stop" the process if any \bar{u}_i^k leaves the ϵ -nbd N_{ϵ} . Then we obtain convergence by 1.7.5. as in 2.2.1.

Martingale theory is not essential here, for we can always use a staircase form of argument (1.11.3.) at boundaries to give Nash convergence, and also to give more intuitive insight.

b) Suppose $\bar{u}_j^1 \uparrow 1, j \bar{u}_{\beta}^1 \uparrow 1$ with β not Nash.

Then as in 2.2.1. $\exists n, k$ s.t. $\Delta \bar{u}_{n+\beta_k}^k > 0$ in N_{ϵ} ϵ -nbd and then under \mathcal{R}_0 we have instability and the boundary learning theorem 1.7.5. gives $\bar{u}_n^k \rightarrow 0$ so β is not an absorbing boundary.

Hence, since boundaries communicate,

$$\lim_{t \rightarrow \infty} \bar{u}_i^k(t) = 1 \quad \text{only if } \beta_k = i \text{ is component of Nash Point, } \beta_*$$

So $\bar{u}_i^k \uparrow 1, j \bar{u}_{\beta}^k \uparrow 1$, only if β is Nash Point.

//

In this most general case we have the difficulty of the non-uniqueness of Nash Point, but in many simple games we actually find the Nash Point is unique, as in 2.5.4., 2.5.5. below.

We apply 2.5.2 to games considered by Russian authors, initiated by Tsetlin (1963), in the context of structured automata. I believe that it is more natural to use unstructured automata, in particular, Π -cells, whenever we wish autonomous tracking to probabilistically stable boundaries. We only require structure to fix optimal non-boundary behaviour (as in cyclic $\Phi^* = 0$), in a new boundary formulation.

Tsetlin called the following homogenous games and his papers give computer simulations rather than convergence proofs.

Example 2.5.3. The Investment Game.

Let $\alpha_1, \dots, \alpha_n \geq 0$ and $\forall \eta = \text{no of players}$
and $\alpha_i / m_i = \Pr(\text{reward using } i \text{ if } m_i \text{ automata use it}),$

so $q_{\beta}^k = \alpha_i / m_i$ if $\beta_j = i$ occurs m_i times.

Then a Nash Point occurs if $\alpha_i / m_i > (\alpha_j / m_{j+1}) \quad \forall i, j.$

We can prove this exists by induction; if $\alpha_{1/2} < \alpha_r$ then result, else put $\alpha_r < \alpha_{1/2} < \alpha_{r-1}$ and continue re-ordering until stable.

It is unclear whether the result is unique w.r.t. $\underline{m} = (m_1, \dots, m_n)$ at a Nash Point.

Now if we allow Π -cells to play this game under Φ_r , then 2.5.2. gives Nash convergence with $q_{\beta}^k = 2\alpha_i / m_i - 1$. //

Example 2.5.4. Investment Game with Common Bank.

Now put $q_{\beta}^k = \frac{1}{n} \sum_j (\alpha_j / m_j)$ to give the investment game

with common bank - we share out winnings.

Then $\beta^* = (\alpha_1, \alpha_2, \dots, \alpha_n)$ or a permutation gives a unique Nash

Point in $m=1$. Each \bar{u} -cell asymptotically takes a different action from the set of the n best actions. If a θ_i fails then the remaining $\theta_{j \neq i}$ converge to Nash Point $(\alpha_1, \dots, \alpha_{n-1})$. Tsetlin calls this feature of behaviour "reliability".

So put $g_{\beta}^k = 2 \left(\frac{1}{n} \sum_j \alpha_j / m_j \right) - 1$ where m_j \bar{u} -cells take

action j , as given by β .

Besides being a Nash Point, this solution β^* is also a More Point in that we achieve maximum payoff.

//

Example 2.5.5. The Gur Game.

Let there be 2 strategies 0, and 1 say, and reward probability = $p\left(\frac{m}{n}\right)$ for all players.

m = number of \bar{u} -cells using 1.

n = " " " " .

And let $p\left(\frac{m}{n}\right)$ have a unique maximum at $\frac{m^*}{n}$ say.

This is the Gur game described by Borovikov and Bryzgalov (1965).

Hence $\exists m^*$ s.t. $p\left(\frac{m^*-1}{n}\right) < p\left(\frac{m^*}{n}\right) > p\left(\frac{m^*+1}{n}\right)$ and $p\left(\frac{m^*}{n}\right) > p\left(\frac{i}{n}\right) \quad i \neq m^*$

and $g_{\beta}^k = 2 \left(p\left(\frac{m}{n}\right) \right) - 1 \quad \forall k$ iff β s.t. m \bar{u} -cells take 1.

We see that this is really a simple case of the previous example in that winnings are shared, with a choice of just 2 actions.

However, it is a fundamental form and more recently Schmukler (1970) has considered it in more depth.

Using 2.5.2. we achieve a.s. convergence to the Nash Point.

//

\bar{u} -cells can only achieve competitive solutions (Nash) since U.L. ensures maximum reward at the next trial. (asymptotically)
So a \bar{u} -cell will always try to "double cross" its opponents unless we reformulate the game as in 2.5.5., to ensure that this would feed back on to its payoff, converting perhaps a Pareto optimum to a Nash Point.



Chapter Three

Structured Learning Automata



Chapter 3.

" Organic life develops away from the concentric unicellular phase as the evolution of the species develops, and progresses along an axis, taking a direction and discovering aims."

Le Corbusier.

(La Ville Radieuse, 1935)

3. Structured Automata.

3.1. The Model for \mathbb{N} -cell Networks.

We have seen that the \mathbb{N} -cell responds asymptotically only to the environmental equilibrium probabilities. (1.12.9. and 2.3.3.) We shall now introduce a structure which responds to environmental fluctuations through utilizing optimal boundary learning (1.5.12.) to ensure, asymptotically a deterministic graph. Tsetlin (1961) introduced non-evolving structured automata in \mathbb{N} and Vorontsova and Varshavskii (1964) conducted computer simulations on structures evolving under the β -rule of Luce (1959). These simulations indicated that an initially random graph will converge to a quasi-linear graph resembling Tsetlin's automata, with near optimal payoff.

Since then, both Fu, in Mendel and Fu (1969), and Feichtinger (1970) have carried out an analysis similar to that of Vorontsova and Varshavskii, for static \mathbb{N} , repeating the same errors when considering increments in transition probabilities. This is almost certainly due to their being unaware of the work of Norman (1966 - 1974). Yet even then there is a considerable amount of additional analysis required if we wish to apply uniform learning rules to \mathbb{N} -cell networks. The theory developed in chapter 1 now enables us to provide such an analytical basis for the evolution of stochastic automata with structure.

Definitions 3.1.1.

i) We define environment $\mathbb{N}(\Delta_{\alpha\beta}, q_{u_i}^{\alpha, s_a})$ with

$\Delta_{\alpha\beta} = \Pr(E_\alpha(t) \rightarrow E_\beta(t+1))$, where E_α is the environment state α .

$q_{u_i}^{\alpha, s_a}(t) = \Pr(\text{stimulus } s_a(t) = s_i \in \{0,1\} \mid u_i(t) \text{ and } E_\alpha(t))$.

and $\tilde{s}(t) = (s_1(t), s_2(t))$, each s_a independent with probability $q_i^{\alpha, s_a}(t)$, as above.

This vector stimulus $\underline{s}(t)$ is used to enable the structure characterization theorems to be proved as a generalization of the method used in 1.12.9. We shall see that this definition prevents the process from "overlapping", just as for 1-cell games we had a stochastic zero-sum, rather than deterministic zero-sum game.

ii) Let the automaton transition probabilities be given by:-

$$\sigma_{ij}^s(t) = \Pr(x_i(t) \rightarrow x_j(t+1) \mid s_2(t) = s \in \{0,1\}), \text{ where } x_i \text{ denotes state } i.$$

iii) Let $\Gamma_i = \{x_k : \pi\text{-cell } \Theta_i \text{ is always used}\}$

$$\text{and } \pi_j^i(t) = \Pr(\text{for } \Theta_i \text{ we use } u_j(t) \text{ at time } t)$$

iv) The reinforcement rules are uniformly learning for

both σ_{ij}^s and π_j^i .

$$\sigma_{ij}^s(t+1) = \sigma_{ij}^s(t) + \Theta_{ij}(\sigma_{ij}^s(t)) (1 - \sigma_{ij}^s(t)) \quad \text{if } s_2(t) = s$$

in $x_i(t)$ and $x_i \rightarrow x_j$ with $s_1(t+1) = 1$.

$$\sigma_{ij}^s(t+1) = \sigma_{ij}^s(t) (1 - \Theta_{jr}(\sigma_{ij}^s(t))) \quad \text{if } s_2(t) = s$$

in $x_i(t)$ and $x_i \rightarrow x_r \neq j$ with $s_1(t+1) = 1$.

$$\sigma_{ij}^s(t+1) = \sigma_{ij}^s(t) \quad \text{if } s_1(t+1) = 0.$$

$$\text{and with normalization } \sum_j \sigma_{ij}^s(t+1) = 1 \quad \forall i, s.$$

v) For Θ_i we have $s_a(t) = s$ with probability $\sum_j \pi_j^i q_j^{a,s}$, if E_α , so that each pair $s_1(t), s_2(t)$ are independent once we are in a certain $x_j(t)$.

vi) If $i \in \Gamma_k$ in definition iv) above then we reinforce

$\pi_j^k(t)$ as below.

$$\pi_j^k(t+1) = \pi_j^k(t) + \Theta_{ij}'(\pi_j^k(t)) (1 - \pi_j^k(t)) \quad \text{if } s_1(t) = 1, u_j(t).$$

$$\pi_j^k(t+1) = \pi_j^k(t) (1 - \Theta_{jm}'(\pi_j^k(t))) \quad \text{if } s_1(t) = 1, u_m \neq j(t).$$

$$\pi_j^k(t+1) = \pi_j^k(t) \quad \text{if } s_1(t) = 0.$$

with normalization $\sum_j \pi_j^k(t+1) = 1 \quad \forall k$

The state transition is made independent of $s_1(t)$, even though this would give us more information on Π . However, if we just used $s_1(t)=s_2(t)=s(t)$, then although it seems likely that we would obtain similar limiting structures, we could not split the process with respect to conditional expectation of increments in σ_{ij}^s . Thus the vector stimulus $s(t)$ is more a technical device for proving theorems rather than an intrinsic feature of the learning process.

The same Π -cell \otimes_k can be reinforced from any state $i \in \Gamma_k$. Only through this formulation can we obtain evolving Π -cell networks with "memory", and we shall discuss the partitioning of actions in section 3.7.

We denote this network of Π -cells by $g'(\otimes)$. We shall first consider the case in which we have $\pi_{ij} \in \{0,1\} \quad \forall i,j$ so that only σ_{ij}^s evolves.

$$\longrightarrow x_i(t) \mid s_1(t) \mid s_2(t) \xrightarrow{\sigma_{ij}^s(t)} x_j(t+1) \mid s_1(t+1) \mid s_2(t+1) \longrightarrow$$

where $s_a(t)$ are received with probability $\sum_j \pi_j^i(t) q_j^{\alpha, s_a}$ if in E_α .

This does impose some initial structure, but we shall consider the complete case of simultaneous σ_{ij}^s and π_k^h reinforcement later. It will then be seen that this adds very little to the analysis of limiting structures, and yet we can then immediately generalise to hierarchies of Π -cells, $g'(\otimes)$, in 3.9.

If $\pi_{ij} \in \{0,1\}$ then we have a Π_1 -cell which is just a single action taken with probability 1. Rather than refer to this as a Π_1 -cell network we shall use the term structured automaton.

In the next section we begin by showing the relationship between structures and Π -cells. The structured automaton is denoted $g'(\otimes)$

with $\textcircled{1}$ as a π_1 -cell, or equivalently, an action.

3.2. Static Environments.

We shall prove that an evolving structured automaton asymptotically maximizes the payoff in a static environment \mathcal{M} . Thus it performs the same function as a singleton π -cell.

Definitions 3.2.1.

i) Let $\mathcal{E}_i(t) = \sigma_{ij}^s(t)$ for $s_2(t)=s$ and $x_j(t)$.

Indeed, here we could put $s_1=s_2=s$, but we shall keep to $s(t)$ in order to avoid confusion.

ii) Let $\Delta\sigma_{ij}^s(t) = \mathcal{E}(\sigma_{ij}^s(t+1) | \mathcal{F}_t) - \sigma_{ij}^s(t)$.

where \mathcal{F}_t is the field of events for $0 \leq t' \leq t$ and also includes $\sigma_{ij}^s(t), s_1(t)$.

We could define $s_1(t)$ and $s_2(t+\epsilon)$ say, so that \mathcal{F}_t actually contains all events at times $t' \leq t$, but since we are more concerned with the concepts rather than complete rigour, we shall just split the process at $s_1(t)$, as in the diagram of 3.1.

For static environments:- $\Delta\sigma_{ij}^s(t) = \mathcal{E}(\sigma_{ij}^s(t+1) | \mathcal{E}_i(t), \sigma_{ij}^s(t), s_1(t)) - \sigma_{ij}^s(t)$.

iii) A state i will be called +ve recurrent if

$$\underline{\text{not}} \quad \lim_{t \rightarrow \infty} \mathcal{E}_i(t) = 0.$$

Theorem 3.2.2.

In a static environment under \mathcal{R}_0 :-

a) $\lim_{t \rightarrow \infty} \sigma_{ij}^s(t) = v_{ij}^s \in [0,1] \quad \forall i,j$.

b) If i is +ve recurrent then $\lim_{t \rightarrow \infty} \sigma_{ij}^s(t) = 1$ for some j
s.t. $q_j^1 \geq q_k^1 \quad k$.

Proof.

We have $\Delta\sigma_{ij}^s(t) = q_i^s \mathcal{E}_i(t) \sigma_{ij}^s(t) \leq \sigma_{ik}^s(t) \mathcal{E}_{jk}^s(t) (q_j^s - q_k^s)$.

and for $q_j^1 \geq q_k^1$ we find $\Delta\sigma_{ij}^s(t) \geq 0 \quad \forall t$

a) Now by semi-martingale convergence theorem, $\sigma_{ij}^s \xrightarrow{a.s.} v_{ij}^s$.

Now proceed by induction on q_i^1 in order of magnitude as in 1.3.1. to obtain the result for a).

b) If $\Delta\sigma_{ij}^s > 0$ except at $\sigma_{ij}^s \in \{0,1\}$ then $v_{ij}^s \in \{0,1\}$ and by n-action optimality theorem for R_0 , 1.6.4. we obtain

$$v_{ij}^s = 1 \quad \text{only if } q_j^1 \geq q_k^1, \forall k.$$

Similarly, if $\Delta r_{ij}^s \equiv 0$ with $q_j = q_k \forall k$, then we obtain the result since the conditional variance must vanish at v_{ij}^s , as for the singleton $\bar{\pi}$ -cell.

It is easier to put:-

$$\Delta^* \sigma_{ij}^s(t) = \mathbb{E}(\sigma_{ij}^s(t+1) | \sigma_{ij}^s(t) \text{ and } x_i(t) = \sigma_{ij}^s(t)).$$

for we now obtain a $\bar{\pi}$ -cell if i is +ve recurrent. But for dynamic environments we gain more insight by using the definition of 3.2.1. and so we have adopted it here also.

//

So the +ve recurrent states form a deterministic, connected structure with $q_i = q_j = \max_k q_k \forall i, j$ s.t. $\gamma_i > 0$ where

$$\gamma_i = \sum_{j,s} \gamma_j q_j^s v_{ji}^s = \lim_{t \rightarrow \infty} \bar{E}_i(t) = \text{limiting time-averaged } E_i(t).$$

Thus a structured automaton maximizes its payoff in a static and on considering q_i as a $\bar{\pi}_1$ -cell, \emptyset , we have that the absorbing class is a subset of Π_{opt} .

Also for each (i,s) , σ_{ij}^s acts as a $\bar{\pi}$ -cell with links $i \rightarrow j$ as actions.

Corollary 3.2.3.

For $q'(\emptyset)$ in \mathcal{M} under R_0 :-

$$a) \lim_{t \rightarrow \infty} \bar{\pi}_j^i = v_{ij}^s \in [0,1] \quad \forall i,j \quad \text{and} \quad \lim_{t \rightarrow \infty} \sigma_{ij}^s = v_{ij}^s \quad \forall i,j,s.$$

b) If i is +ve recurrent, $i \in \Gamma_k$ then

$$i) \lim_{t \rightarrow \infty} \bar{\pi}_j^k = 1 \quad \text{for some } j \text{ s.t. } q_j \geq q_m \quad \forall m.$$

$$ii) \lim_{t \rightarrow \infty} \sigma_{ij}^s = 1 \quad \text{for some } j \text{ s.t. } q_j \geq q_m \quad \forall m.$$

Proof.

Now $\Delta \bar{u}_j^i \triangleq E(\bar{u}_j^i(t+1) | \mathcal{F}_t) - \bar{u}_j^i(t)$

$$\Delta \bar{u}_j^k(t) = \sum_{i \in P_k} E_i(t) \bar{u}_j^k(t) \sum_m \bar{u}_m^k(q_j - q_m) \theta_{jm}'(\bar{u}^k(t)).$$

and as in 3.2.2. $\bar{u}_j^k \rightarrow v_j^k$ as $t \rightarrow \infty$ and then by induction.

And for $\Delta \sigma_{ij}^s(t)$ we must consider boundary behaviour, so let N_j^i

be a neighbourhood of v_j^i s.t. $|\bar{u}_j^i - v_j^i| < \epsilon_j^i$ say.

$$\Delta \sigma_{ij}^s(t) = E_i(t) \sigma_{ij}^s(t) \left(\sum_h \bar{u}_h^k q_h^s \right) \sum_{\substack{a: ik \\ \text{with } j \in P_m}} \theta_{jk} \left(\sum_a v_a^m q_a^i - \sum_b v_b^m q_b^i \right) + O(F(\epsilon_j^i)).$$

and $F(\epsilon_j^i)$ is arbitrarily small, compared with the leading terms, as in 2.5.2.

Now order $\sum_a v_a^m q_a^i$ and take m^* s.t. $\sum_a v_a^{m^*} q_a^i \geq \sum_a v_a^m q_a^i \quad \forall m$.

If $\sum_a v_a^m q_a^i = \sum_a v_a^{m^*} q_a^i \quad \forall m$ then $\Delta \sigma_{ij}^s \rightarrow 0$ and we approach a martingale

form, from which we can obtain convergence $\sigma_{ij}^s \rightarrow v_j^i$.

If the limiting \bar{u} -cell payoffs are unequal, then we successively form semi-martingales for the result, as in 1.3.1.

If i is +ve recurrent, then the conditional variance > 0 gives $v_b^i, v_a^i \in \{0, 1\}$ since $\lim_{t \rightarrow \infty} E_i(t) = 0$. (Note that this limit will only exist in the non +ve recurrent case).

In neighbourhood N_j^i of v_j^i we can apply optimal boundary learning theory for sem-martingale $\sigma_{ij}^s(t)$ as in 1.7.4., since all boundaries communicate. (Note that the +ve recurrent set of states is non-empty for all finite automata).

Again, as in 2.5.2., we can use a staircase construction instead on martingale theory to give us convergence only to stable boundaries. //

A \bar{u} -cell network thus has the same limiting payoff as a singleton \bar{u} -cell in static \mathcal{M} , so we certainly lose nothing (apart from simplicity) by this extension.

Our aim is to allow $g'(\theta)$ to adapt to a dynamic \mathcal{M} through

use of the equilibrium distribution for σ_{ij}^s . This adaptation will characterise a certain family of structures related to the work of Tsetlin (1961) and Stratonovich (1964) in that we discretize bayesian updating in the environmental likelihood simplex. First we strengthen our intuition by considering some simple properties of deterministic structures in dynamic \mathcal{M} .

3.3. Evolution in a 2-Medium.

Definitions 3.3.1.

- i) Let $\rho_i^\alpha = \text{prob}(E_\alpha, \text{ in equilibrium } | \text{ in state } i)$.
 so $\rho_i^\alpha = \gamma_i^\alpha / \gamma_i$ with $\gamma_i^\alpha = \text{prob}(E_\alpha \text{ and } x_i, \text{ in equilibrium})$.

- ii) The symmetric 2-medium $\mathcal{M}_2(A_{2p}, q_{ij}^{\alpha, s})$ is defined as:-

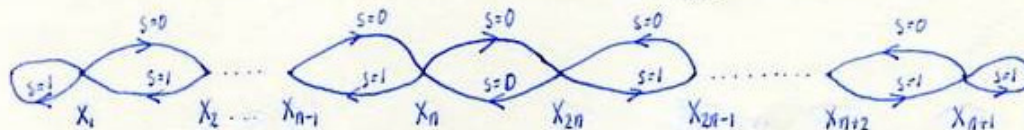
$$A_{2p} = \begin{pmatrix} 1-\delta & \delta \\ \delta & 1-\delta \end{pmatrix}$$

$$\text{and } q_1^{11} = q_0^{01} = q_0^{10} = q_1^{00} = q.$$

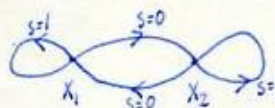
$$q_1^{10} = q_0^{00} = q_0^{11} = q_1^{01} = p = 1 - q.$$

The results from the analysis here have natural extensions to many other \mathcal{M} and in the following section 3.4. we consider the symmetric n-medium \mathcal{M}_n .

- iii) Let the structured automaton $L_{n,n}$ be defined as:-



and so that $L_{1,1}$ is just:-



and for $1 \leq i \leq n$ $i \in \mathcal{P}_1$.

$n+1 \leq i \leq 2n$ $i \in \mathcal{P}_0$.

This $L_{n,n}$ is the basic linear symmetric automaton.

Lemma 3.3.2.

- a) $\rho_1^0 = q - \delta(q-p)$ for $L_{1,1}$.

b) $\{r_{r,n+1}\}^0, \{r_{r,n}\}^1 \uparrow$ as $r \downarrow$ for $L_{n,n}$.

c) $\{r_{n-1}\}^0 = q$ iff $f(ny)/f(y) = 1/\delta$ for $L_{n,n}$ where
 $f(y) = \cosh y - 1$ and $\cosh y = \frac{1}{2pq} \left(\frac{1-\delta}{1-2\delta} \right) - 1$.

d) $\bar{R}(L_{2,2}) = \bar{R}(L_{1,1})$ iff $\{r_{n-1}\}^0 = q$ with $n=2$, where

$$\bar{R} = \text{average payoff} = \sum_{i,j} \gamma_i^x q_j^x.$$

Proof.

a) Solve $\gamma_j^B = \sum_{i,j} \gamma_i^x q_j^x \delta_{ij} \Delta_{\alpha\beta}$ 1)

b) Solve 1) for $L_{n,n}$ as in the paper of Tsetlin (1961) and then show that $\gamma_r^0/\gamma_r^1 > \gamma_{r+1}^0/\gamma_{r+1}^1$ and hence that $\{r_{r,n}\}^0 \uparrow$ as $r \downarrow$ and similarly for $\{r_{r,n+1}\}^1 \uparrow$ as $r \downarrow$.

c) Again, as in b), the result arises from manipulation and its truth is closely linked with our next result d).

d) It can be shown that $\bar{R}(L_{n,n}) = \frac{\tanh(ny/2)}{\coth(y/2) + n \coth(ny)}$ 2)

with $\delta = 2\delta/(1-2\delta)(q-p)^2$ and $\cosh(y)$ as before.

Using c) we obtain $\{r_2\}^0 = q$ if $(\cosh 2y - 1)\delta = \cosh y - 1$

or $\cosh(y) = (\frac{1}{2}\delta - 1)$

and substituting for $\cosh(y)$ gives $\delta(1-\delta)/(1-2\delta) = pq$.

But Tsetlin (1963) gives this as the condition for $\bar{R}(L_{2,2}) = \bar{R}(L_{1,1})$ and indeed it is not too difficult to verify this using 2).

//

The result d) is almost certainly false for the general case when $\bar{R}(L_{n,n}) = \bar{R}(L_{n-1,n-1})$, and this is related to the optimal automaton not possessing the SOSA property (self-one-step-ahead) except for $L_{1,1}$. This property is defined in 3.3.8. and considered

in detail in subsequent sections.

The manipulation required to achieve the above results is most tedious and so apart from a), they have not been extended to M_n even though they almost certainly exist.

We shall now see the relevance of 3.3.2. in giving rise to the family of limiting structures denoted by A_0 . First we consider the fine structure of the process before taking the equilibrium values which R_0 utilizes.

Definitions 3.3.3.

- i) Let $\Theta_j^\beta(t) = \Pr(\text{in } E_\beta \mid \text{in state } j \text{ at time } t)$.
- ii) Let $\omega_i^\alpha(t) = \Pr(\text{in } E_\alpha \text{ and state } i \text{ at } t)$.
- iii) Let $e_i(t) = \Pr(\text{in state } i \text{ at time } t)$.
- iv) We define $\Delta \sigma_{ij}^s(t) = \mathbb{E}(\sigma_{ij}^s(t+1) \mid s_1(t), \{\sigma_{ij}^s(t)\}, \mathcal{F}_t) - \sigma_{ij}^s(t)$.

Now by bayesian rules let $(\omega_i^\alpha)^*(t+1) = \sum_\beta \Theta_j^\beta(t) q_j^{\alpha, s_2} \sigma_{ji}^{s_2} \Delta \beta^\alpha / (\sum_\beta \Theta_j^\beta q_j^{\alpha, s_2})$ 3)

and where $s_2(t)$ and $u_j(t)$ are used in $x_j(t)$.

Now we receive $s_1(t+1)$ and $\omega_i^\alpha(t+1) = q_i^{\alpha, s_1} (\omega_i^\alpha)^*(t+1) / (\sum_\beta (\omega_i^\beta)^*(t+1) q_i^{\beta, s_1})$ 4)

So using 3) and 4) we have $\Theta_i^\alpha \xrightarrow{s_2(t), s_1(t+1)} \omega_i^\alpha(t+1)$.

Now $e_i(t+1) = \sum_\alpha \Theta_j^\alpha(t) q_j^{\alpha, s_2} \sigma_{ji}^{s_2} / (\sum_\alpha q_j^{\alpha, s_2} \Theta_j^\alpha)$ 5)

$$\text{so } \Theta_i^\alpha(t+1) = \omega_i^\alpha(t+1) / e_i(t+1)$$

Hence we generate our conditional probabilities by inductively using

3) \rightarrow 5), and $\Delta \sigma_{ij}^s(t) = \mathbb{E}(\sigma_{ij}^s(t+1) \mid \Theta_i^\alpha(t), \{\sigma_{ij}^s(t)\}, e_i(t)) - \sigma_{ij}^s(t)$.

Note that we can take $(\omega_i^\alpha)^* = \Theta_i^\alpha q_i^{\alpha, s_1} q_i^{\alpha, s_2} / \sum_\beta \Theta_i^\beta q_i^{\beta, s_1} q_i^{\beta, s_2}$ 6)

for the \sum updating and the result is the same as if we had two successive updatings, for the $s_2(t)$ and $s_1(t+1)$ taken individually.

Lemma 3.3.4.

i) If $\{o_j^s\}$ is constant then $\bar{w}_i^\alpha = \gamma_i^\alpha$, $\bar{\theta}_i^\alpha = \xi_i^\alpha$, $\bar{e}_i = \gamma_i$

Where $\gamma_i^\alpha = \sum_{\beta, j} \gamma_j^\beta \Delta_{\beta\alpha} q_j^{\beta, s} \cdot \sigma_{ji}^s$, $\gamma_i = \sum_{\alpha} \gamma_i^\alpha$ and $\xi_i^\alpha = \gamma_i^\alpha / \gamma_i$

ii) $\Delta \sigma_{ij}^s(t) = c \sum_{\alpha, k} (w_k^\alpha / t + 1) R_{jk}^\alpha$ where $c = \sum_{\beta} q_j^{\beta, s} w_i^\beta(t)$

and $R_{jk}^\alpha = \theta_{jk}(\{o_j^s\}) (q_j^{\alpha, 1} - q_k^{\alpha, 1})$ = reward comparison function.

Proof.

i) This is a consequence of bayes' rule and we indeed also have $\bar{w}_i^\alpha \cdot (\bar{w}_i^\alpha)^* = \gamma_i^\alpha$ and this is proved by writing out all possibilities.

ii) $\Delta \sigma_{ij}^s(t) = \sum_{\alpha} w_i^\alpha(t) q_j^{\alpha, s} \sum_{k, \beta} \sigma_{ik}^s \Delta_{\alpha\beta} \theta_{jk}(\{o_j^s(t)\}) (q_j^{\beta, 1} - q_k^{\beta, 1})$

but by i) we have the result immediately since we just update w.r.t. the stimulus $s_2(t)$ to determine our expected payoffs on taking an action at the next trial.

//

Lemma 3.3.5.

The maximum average payoff is achieved by the bayesian rule:-

Let a) $w_\alpha(t+1) = \sum_{\beta} w_\beta(t) q_j^{\beta, s} \Delta_{\beta\alpha} / \sum_{\beta} w_\beta(t) q_j^{\beta, s}$ if $s(t)=s$ and use $u_1(t)$.

where $w_\alpha(t) = \Pr(\text{in } E_\alpha \text{ at time } t)$.

Then b) At time $t+1$ take u_j only if $\sum_{\alpha} w_\alpha (q_j^{\alpha, 1} - q_k^{\alpha, 1}) \geq 0 \quad \forall k$.

(Here we have the simplified model with single scalar stimulus s , and we assume that the environmental parameters $\{\Delta_{\alpha\beta}, q_j^{\alpha, s}\}$ are known, so that only E is unknown. For vector stimulus we modify a) in a similar way to eqn 6).

Proof.

Stratonovich (1964) proves this.

//

Thus the O.S.A. policy is optimal here since environmental information is "gained equally from all states". It is for this reason that our discrete automaton will achieve environmental

adaptation, using reinforcement only over successive trials.

We now consider evolving automata in M_2 and our next theorem gives both existence and characterization of the limiting family of structures.

Theorem 3.3.6.

In M_2 under R_0 with $q, 1-\delta > \frac{1}{2}$.

- a) $\sigma_{ij}^{s=0} \rightarrow 1$ only if $\rho_i^\alpha(v_{ij}^s) > q$ with $i \in P_\alpha, j \in P_i$.
 b) $\sigma_{ij}^{s=1} \rightarrow 1$ only if $\rho_i^\alpha(v_{ij}^s) > p$ with $j \in P_i$.

where $q_i^{11} = q$ for $i \in P_1$ say $1 \leq i \leq n$

and $q_i^{01} = q$ for $i \in P_0$ say $n+1 \leq i \leq 2n$

Thus we have $q_i^{11} = q$, and $q_i^{00} = p$ for $i \in P_\alpha$.

- c) $\sigma_{ij}^s \rightarrow v_{ij}^s \in \{0, 1\}$ for +ve recurrent i .

Proof.

We write out the expected increments:-

$$\Delta \sigma_{ij}^0(t) = s_{vt} K \rho_i(t) \sigma_{ij}^0(t) (\theta_i^0 - q) \sum_{k \in P_1 - i} \sigma_{ik}^0 \theta_{jk}(\{\sigma_{ij}^0\}).$$

$$\Delta \sigma_{ij}^1(t) = s_{vt} K \rho_i(t) \sigma_{ij}^1(t) (\theta_i^1 - p) \sum_{k \in P_1 - i} \sigma_{ik}^1 \theta_{jk}(\{\sigma_{ij}^1\}).$$

where $s_{vt} = +1$ if $v=i$ and $i \in P_v, j \in P_v$.

$s_{vt} = -1$ if $v \neq i$ $K = (1-2p)(1-2q)$.

Consider now $\{\sigma_{ij}^s\}$ fixed so that we can consider the process for each state as a random walk defined on the environment markov chain. Thus near the boundary absorbing set we can apply the techniques of section 1.12. which were used to prove 1.12.9.

Here for each state i , we have an underlying markov process with equilibrium distribution $\pi_i^a(\sigma_{ij}^s)$, which is an example of the type of process considered by Miller (1962).

$$\text{Then } \Delta \sigma_{ij}^0 = s_{vt} K \pi_i^a(\sigma_{ij}^0) (\theta_i^0 - q) \sum_{k \in P_1 - i} \sigma_{ik}^0 \theta_{jk}(\{\sigma_{ij}^0\}). \quad 8)$$

$$\overline{\Delta \sigma_{ij}^s} = s_{vz} K \tau_i \sigma_{ij}^s (q_i^v - p) \sum_{k=1}^s \sigma_{ik}^s (\theta_{jk} (\sigma_{ij}^s)).$$

9)

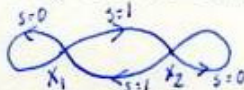
where $\overline{\Delta \sigma_{ij}^s} \triangleq$ mean increment in σ_{ij}^s w.r.t. the equilibrium process.

We can do this since the process is split up by the vector stimulus $s(t)$ so that successive increments do not "overlap", as shown in 3.1.

We now apply optimal boundary learning in a dynamic medium for the result (1.7.5, 1.12.9.) which is now immediate from 7) and 8), noting that:-

i) A limiting structure always exists since ${}_2L_1$ is stable with $q_i^v = q - \delta(q - p) > \frac{1}{2}$. Convergence to stable boundaries follows by a β -staircase construction as in 1.12.9. for each +ve recurrent state. This gives us c) and a non-empty limiting set of structures.

ii) If $\delta > \frac{1}{2}$ then the only stable limit is ${}_2L_1^c$ as below



iii) If $q < \frac{1}{2}$ then the only stable structure is ${}_2L_1$, assuming also $\delta < \frac{1}{2}$.

iv) Although $\overline{\Delta \sigma_{ij}^s} > 0$ say for all $j \in P_i, i \in P_\alpha$ if $q_i^v > q$, we get convergence to just one $j \in P_i$, as in the n-action theorem with multiple optima 1.6.4. The dynamic medium adds nothing for if we have 2-actions with $q_i^{\alpha, s} = q_i^{\alpha, s} \forall \alpha$ then $\delta(\bar{u}) = \bar{u}$ for all \bar{u} , as we have a martingale.

Thus in contrast with bayesian rules which oscillate if we try to decide between 2 identical hypotheses, our U.L. rules still give boundary convergence. This property is important in their application to \bar{u} -cell networks.

v) The +ve recurrent states i are those with $\gamma_i / (v_{ij}^s) > 0$.

//

Corollary 3.3.7.

In M_2 a) ${}_2L_r$ are stable if $\frac{1}{r-1} > q$.

b) If $\delta(1-\delta)/q(1-q) > (1-2\delta)$

then ${}_2L_1$ is the unique and optimal limit, when $\delta < \frac{1}{2}$.

Proof.

a) We have $s_i^1 > q$ gives $i \rightarrow j \in P_i$ and by 3.3.2.
 q_i^1 as $r \downarrow$ with $s_{r-1}^1 > q$ iff $f(q)/f(q) < 1/\delta$.

b) This is just case d) of 3.3.2. where ${}_2L_1$ is stable
 iff it is optimal, and uniqueness follows easily. //

Theorem 3.3.6. actually gives $i \xrightarrow{s} j$ if and only if we
 maximize the expected reward at the next trial w.r.t. the equilibrium
 distribution γ_i^s in our present state. So we essentially have
 each σ_{ij}^s as a π -cell $\otimes_{(i,s)}$ in medium with equilibrium γ_i^s and
 transitions $(\alpha, i) \rightarrow (\beta, j)$ with probability $A_{\alpha\beta} q_i^s \sigma_{ij}^s$, and with
 actions u_j . Thus we have σ_{ij}^s as a process defined on a markov
 chain in the manner of Miller (1962) and Keilson and Wishart (1965),
 in the slow learning limit $\theta \downarrow 0$. The σ_{ij}^s processes then interact
 through $\gamma_i^s(\sigma_{ij}^s)$ to give a certain family of stable solutions
 which all have the SOSA property. Thus we just need apply 1.12.9. to
 each $\otimes_{(i,s)}$ for 3.3.6.
Definition 3.3.8.

- i) A deterministic automaton has the SOSA property if for
 every state i , we have $v_i^s = 1$ only if u_j maximizes $\sum_{\beta, j} \gamma_i^s A_{\alpha\beta} q_i^s \sigma_{ij}^s$
 where $\gamma_i^s = \gamma_i^s / \gamma_i$ $= \text{Pr}(\text{ in } E_4 | \text{ in } x_1)$. (SOSA \sim self-one-step-ahead).
 ii) We denote the family of automata with the SOSA property by \mathcal{A}_0 .

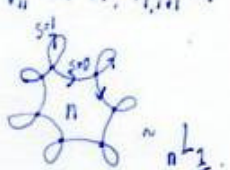
Thus the automaton maximizes the payoff at the next trial
 w.r.t. its own equilibrium distribution over states and environments.

3.4. Evolution in an n-Medium.

Definition 3.4.1.

- i) We define the symmetric n-medium $\mathcal{M}_n(\Delta_{\alpha\beta}, q_{ij}^{\alpha\beta})$ as:-
- a) $\Delta_{\alpha\beta} = \begin{pmatrix} 1-(n-1)\delta & \delta & \dots & \delta \\ \delta & 1-(n-1)\delta & \dots & \delta \\ \vdots & \vdots & \ddots & \vdots \\ \delta & \delta & \dots & 1-(n-1)\delta \end{pmatrix}$ b) $q_{ij}^{\alpha\beta} = q = q_{j\neq i}^{\alpha\beta}$, else $1-q=p$,
and $i \in P_\alpha$ if $q_{i,i}^{\alpha\beta} = q$.

So $\Delta_{\alpha\beta} = \delta \quad \alpha \neq \beta$ with $\sum_{\beta} \Delta_{\alpha\beta} = 1$
 $\Delta_{\alpha\alpha} = (1-(n-1)\delta) \quad \alpha = \beta$

- ii) The automaton ${}_nL_1$ has transitions $v_{ii}^{s=1} = 1 \quad v_{i,i+1}^{s=0} = 1$ with
 $v_{ni}^0 = 1$ so that we have  with
 n indicating the number of states.

We now prove a lemma similar in nature to 3.3.2.

Lemma 3.4.2.

If $\delta < \frac{1}{n}$ and $q > \frac{1}{2}$ then:-

- a) $\gamma_i^r \downarrow$ in r for ${}_nL_1$.
- b) $\bar{R}({}_nL_1) > p + ((q-p)/n)$, which is the \bar{n} -cell limit in \mathcal{M}_n under \mathcal{R}_0 .
- c) ${}_nL_1$ has the SOSA property whilst ${}_m({}_nL_1)$ is not SOSA in \mathcal{M}_n .
- d) $\bar{R}({}_nL_1) > \bar{R}(\lim_{\theta \rightarrow 0} L_{R-P})$ in \mathcal{M}_n .
- e) $\bar{R}(L_{R-P}) = \bar{R}({}_nL_1)$, in any static environment.

Proof.

a) Solve $\gamma_j^{\beta} = \sum_{i=1}^n \gamma_i^{\alpha} q_{ij}^{\alpha\beta} \Delta_{\alpha\beta}$ to get

$$\gamma_1^{r+1} = \frac{1}{n^2} \left(1 - \left(\frac{1}{p} - \frac{1}{q} \right) (1-n\delta)^{r-1} \left(\gamma^n \left(\frac{1+nq\frac{1}{p}}{1+nq\frac{1}{q}} \right) - 1 \right)^{r-2} \right) < \frac{1}{n^2}$$

with $\gamma = \left(\frac{q+nq\delta}{q(1-n\delta)} \right) > 1$.

and $\gamma_i^1 = \frac{1}{n^2}(p+q\delta) \left[(q+np\delta) - \left(\frac{q}{p} - 1 \right) n\delta (\delta^n \theta^{-1} - 1)^{-1} \right] > \frac{1}{n^2}$

with $\theta = (1+np\frac{q}{p})/(1+nq\delta p)$ and $\sum \gamma_i^2 = \frac{1}{n}$.

and for $\delta=0$, $\gamma_i^{r+1} = \frac{p}{q} \gamma_i^r$, $\delta = \frac{1}{n}$, $\gamma_i^r = \gamma_i^s \forall r, s$.

Clearly then $\gamma_i^r \downarrow$ by combining the results above.

b) $\bar{R}_{nL_1} = p + \gamma_1^1(q-p)n$ with $\gamma_1^1 > \frac{1}{n^2}$ and hence result.

c) We need to check $\omega_i^{s=1}/\omega_i^{s=0} > 1$ and $\omega_2^{s=0}/\omega_1^{s=0} > 1$.

where $\omega_r^s = \Pr(\text{in } E_r \mid \text{use } u_1 \text{ and receive stimulus } s, \underline{\omega})$.

and $\omega_r = \Pr(\text{in } E_r \text{ when in state } x_1) = \frac{p}{q}$.

$$\omega_1/\omega_r = \frac{(w_1 q (1-(n-1)\delta) + (1-w_1)p\delta)}{(w_1 q \delta + w_r (1-(n-1)\delta)p + (1-(w_1+w_r))\delta p)} = \frac{w_1 q}{w_r p}$$

and $\omega_1/\omega_r > 1$ iff $w_1/w_r > p/q$.

Similarly $\omega_2/\omega_1 > 1$ iff $w_2/w_1 > p/q$.

and as $w_2/w_r > 1$ we need only consider w_2 .

But as $\delta \downarrow 0$ $\gamma_i^{r+1}/\gamma_i^r \rightarrow p/q$ $\delta \uparrow \frac{1}{n}$ $\gamma_i^{r+1}/\gamma_i^r \rightarrow 1$

and by monotonicity we have $p/q < \gamma_i^{r+1}/\gamma_i^r < 1$.

Hence nL_1 is SOSA.

To show that $m \setminus nL_1$ is not SOSA let us take actions u_1, \dots, u_m with states x_1, \dots, x_m and actions $u_{r,m}$ with no states.

Then $\frac{q^{n \setminus s \leq m}}{q_r} = \frac{1}{n}$ by symmetry for all states $r \leq m$.

yet $\frac{q^{r \neq r}}{q_r^{t \leq m}} < \frac{1}{n}$. So $\frac{q^{s \geq m}}{q_r} / \frac{q^{t \neq r}}{q_r^{t \leq m}} > 1$

So if we wish to maximize payoff at the next trial given a penalty we should take an action $u_{r,m}$ which the structure will not allow.

Hence $m < n$ L_1 is not SOSA in M_n .

d) Here, L_{R-P} is the only known reinforcement rule for unstructured automata "tracking" in a dynamic environment.

$$\left. \begin{aligned} \bar{u}_i(t+1) &= \bar{u}_i(t) + \theta(1 - \bar{u}_i(t)) \\ \bar{u}_j(t+1) &= \bar{u}_j(t)(1 - \theta) \end{aligned} \right\} \begin{aligned} &u_i(t) \text{ and } s(t) = 1. \\ &j \neq i. \end{aligned}$$

and

$$\left. \begin{aligned} \bar{u}_i(t+1) &= \bar{u}_i(t)(1 - \theta) \\ \bar{u}_j(t+1) &= \bar{u}_j(t) + \theta \bar{u}_j(t)/(n-1) \end{aligned} \right\} \begin{aligned} &u_i(t) \text{ and } s(t) = 0. \\ &j \neq i. \end{aligned}$$

Norman (1972) obtains a limiting normal distribution for sufficiently small θ about the point $E(\bar{u}_i(\infty)) = 1/p_i / \sum_j 1/p_j$, with variance $\sim O(\theta^2)$.

Now from b) we have $\bar{R}(L_1) > p + ((q-p)/n)$, but with L_{R-P} we have

$\bar{R}(L_{R-P}) > p + ((q-p)/(n + (q-p)/p))$, which is the case when the

environment has just switched. But in the limit $\theta \rightarrow 0$ of slow

learning we shall asymptotically have $\bar{R}(L_{R-P}) = p + ((q-p)/n)$, by letting $\theta \ll \delta$.

Thus we have the result d).

e) For static environments we have:-

$$\bar{R}(L_1) = \sum_i \gamma_i q_i \quad \text{where} \quad \gamma_i p_i = \gamma_{i-1} p_{i-1}$$

$$\text{Thus } \gamma_i = 1/p_i / \sum_j 1/p_j \quad \text{So } \bar{R}(L_{R-P}) = \sum_i E(\bar{u}_i(\infty)) q_i = \sum_i \gamma_i q_i = \bar{R}(L_1).$$

Now by d) as $\theta \rightarrow 0$ $\bar{R}(L_1) > \bar{R}(L_{R-P})$

also as $\delta \rightarrow 0$ $\bar{R}(L_{R-P}) \rightarrow \bar{R}(L_1)$

yet $\bar{R}(L_1) > \bar{R}(L_{R-P}) \rightarrow \bar{R}(L_{R-P})$ which gives justification

for a structural strategy rather than unstructured tracking as used by Chandrasekaren (1967). In addition, an evolving \mathcal{G} in M (static) achieves optimal payoff.

//

The result e) has the interpretation that the L_{R-P} rule has the same average payoff as the strategy ; reward - stay put, penalty- move on. This strategy gives precisely the probability matching of Norman and Yellott (1966) and in the later work of Norman (1972).

Remarks 3.4.3.

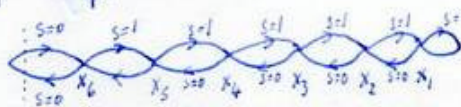
i) It is possible that for some θ and δ , $\bar{R}(L_1) < \bar{R}(L_{R-P})$ yet in many cases we could use an alternative structure nL_r with larger memory. The L_{R-P} rule in dynamic \mathcal{M} has not yet been analysed in detail.

ii) Just as the $2L_1$ was the kernel structure for \mathcal{M}_1 we shall now show that nL_1 is a stable limit in \mathcal{M}_n by considering $\bar{\Delta}\sigma_{ij}^s$.

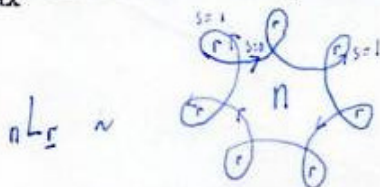
iii) We shall assume that \uparrow_i on "arms" $i \in \mathbb{N}_\alpha$ as in 3.3.2. b)

Certainly as $\delta \downarrow 0$ $\sum_{r=0}^{\alpha} \frac{r}{q} \rightarrow \frac{q}{p}$ with $r \downarrow$ as we proceed down the "arm".

e.g. $1 \leq i \leq 6$ gives $i \in \mathbb{N}_\alpha$, and \uparrow_i as $i \downarrow$.



We then obtain nL_r as SOSA, and for given δ we shall obtain some r_{\max} with r indicating the memory depth. (Or "arm" length.)



iv) Exact calculations for nL_r are prohibitively complicated, when we consider how difficult it is even to obtain the results of 3.3.2. for $2L_r$.

Theorem 3.4.4.

A structured automaton evolving under R_0 in medium \mathcal{M}_n has nL_1 as a stable limiting structure, if $n\delta < 1$, $2p < 1$, and $m < nL_1$ is unstable.

Proof.

We write out:-

$$\bar{\Delta}\sigma_{ij}^{s=0} = K\sigma_{ij}^s \tau_i \left(\sum_{k \in \mathbb{N}_i} \sigma_{ik}^s \theta_{jk}(\sigma_{ij}^{s=0}) (p_i^j - p_i^k) \right) \quad \text{for } j \in \mathbb{N}_i. \quad 10.)$$

$$\bar{\Delta}\sigma_{ij}^{s=0} = K\sigma_{ij}^s \tau_i \left(\sum_{k \in \mathbb{N}_i} \sigma_{ik}^s \theta_{jk}(p_i^j - p_i^k) + \sum_{k \in \mathbb{N}_i} \sigma_{ik}^s \theta_{jk}(p_i^j - p_i^k) \right) \quad \text{for } j \notin \mathbb{N}_i. \quad 11.)$$

and $K = (1 - n\delta)(q-p)$ and $\sigma_{ij}^j = \sigma_{ij}^j(\{\sigma_{ij}^s\})$

and similarly for $s=1$, interchanging p and q ; and for M_i we just obtain our equations in 3.3.6.

The derivations of 10) and 11) are not of great interest, but we shall see in a later theorem why $\Delta\sigma_{ij}^s$ must have the above form.

The $\Delta\sigma_{ij}^s$ are taken with $\{\sigma_{ij}^s\}$ held fixed, as in 3.3.6., then the optimality property of $\theta_{jk}(\{\sigma_{ij}^s\})$ will give convergence to that action which maximizes expected payoff w.r.t. the equilibrium distribution. Let x_i be a +ve recurrent state.

a) Now for $s=0$, $j \in P_i$, $\Delta\sigma_{ij}^{s=0} > 0$ if $\sigma_{ij}^j / \sigma_{ij}^k > q/p \quad \forall k \notin P_i$

and so by boundary learning $\sigma_{ij}^s \rightarrow 1$ for some $j \in P_i$

iff $\sigma_{ij}^j / \sigma_{ij}^k > q/p \quad \forall k \notin P_i$

For if $\sigma_{ij}^j / \sigma_{ij}^k < q/p$ for some k then $\Delta\sigma_{ik}^s > 0$ with

$\sigma_{ij}^s = 1 - \epsilon_{ij}^s$ say in a sufficiently small nbd N_{ij}^s of the boundary v_{ij}^s

b) For $s=0$, $j \notin P_i$

$\Delta\sigma_{ij}^{s=0} > 0$ if $\sigma_{ij}^j / \sigma_{ij}^k > p/q$, $k \in P_i$ and $\sigma_{ij}^j / \sigma_{ij}^m > 1 \quad \forall m \notin P_i$

And by boundary learning $\sigma_{ij}^s \rightarrow 1$ for some $j \in P_i$

iff the above holds, else we get a contradiction as in a).

It is easy to see that this exhausts all possibilities for

if $\exists k$ s.t. $\sigma_{ij}^j / \sigma_{ij}^k < q/p$, $j \in P_i$, $k \notin P_i$ then take $\alpha: \sigma_{ij}^j > \sigma_{ij}^k$
 $\alpha, \beta \notin P_i$
 so that we have $\Delta\sigma_{i\alpha}^s > 0$ with $\beta \in P_i$.

Similarly for $s=1$ we get

c) $\Delta\sigma_{ij}^{s=1} > 0$ with $j \in P_i$ if $\sigma_{ij}^j / \sigma_{ij}^k > p/q$, $k \notin P_i$.
 so that $\sigma_{ij}^s \rightarrow 1$ for some $j \in P_i$ iff the above holds.

d) $\bar{\Delta}_{ij}^{s=1} > 0$ with $j \in P_i$ if $s_i^j/s_i^k > q/p$ $k \in P_i$ and $s_i^j/s_i^m > 1, \forall m \in P_i$
 and $\bar{\Delta}_{ij}^s \rightarrow 1$ for some $j \in P_i$ iff the above holds.

We must note, as in 3.3.6. that we have $\bar{\Delta}_{ij}^s > 0$ for all $j \in P_i$ where u_x maximizes \bar{g} payoff at the next trial. But we get convergence to precisely one $j \in P_i$ using the n-action theorem 1.6.4. with multiple optima, noting $q_i^{k,s} = q_j^{k,s} \forall i, j \in P_i, \forall k, s$.

It is now easy to see that nL_1 is stable and $m \subset nL_1$ is not stable by lemma 3.4.2.

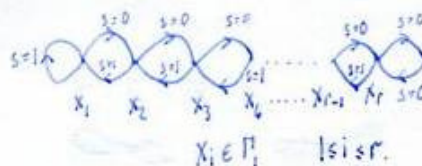
//

Corollary 3.4.5.

In M_n , the automaton nL_r is stable iff

$$s_{r-1}^1/s_{r-1}^2 > q/p$$

where:-



Proof.

By the symmetry of the inequality, if it holds for arm P_1 , then it is true for all arms P_i . Also $s_{r-1}^2 > s_{r-1}^{k \geq 2}$ since $\gamma_i^r \downarrow$

in r for nL_1 by 3.4.2. a).

Then apply 3.4.4. for the result.

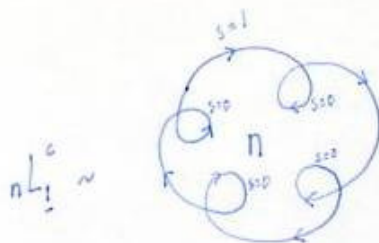
//

The extension of 3.3.7. b) requires the calculation of the equilibrium distribution for nL_2 which has not been done. However, it still remains clear that $\bar{R}(nL_1) = \bar{R}(nL_2)$ when $s_{r-1}^1/s_{r-1}^2 = q/p$, which we proved for $n=2$ in 3.3.2. d).

In general, the optimal v_{ij}^s automaton will not have the SOSA property although the Π_0 family may achieve close to optimal payoff. It is the discreteness of the formulation which causes the difficulties both here and in the general existence conjecture to be given later.

Corollary 3.4.6.

If $\bar{\Delta}_{ij}^s > 1$ then nL_1^C is a stable limiting structure.



with convention, reward links on exterior, and penalty links on interior.

Proof.

Just note $K = (1 - n\delta)(q - p)$ changes sign, and we also check SOSA property is satisfied by nL_1^C by modifying analysis of 3.4.4.

For $q < p$ we actually get no change in stable structures, since sign changes all eventually cancel out. Thus 3.4.4 could be stated for all $q \in [0, 1]$, with just the δ constraint. //

The nL_1^C is probably the unique limit, by symmetry requirements. However, in such enumeration problems of graph theory, there often seems to be the possibility of a subtle counter-example.

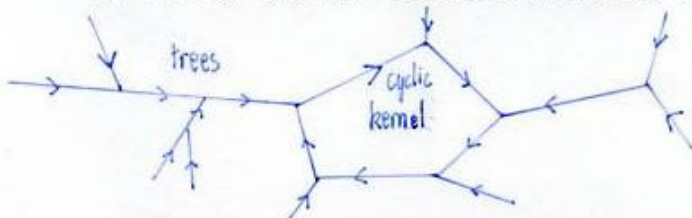
For fixed s , \mathcal{V}_s is a set of functional digraphs, and the following lemma gives the possible basic forms.

Lemma 3.4.7.

For each s , the graph \mathcal{V}_s consists of a set of disconnected subgraphs. Each has precisely one cyclic kernel with trees leading to it.

Proof.

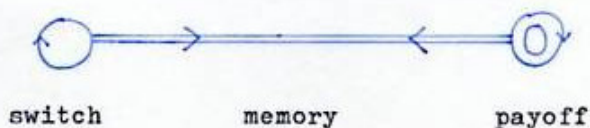
See Harary (1964). The result is just observation. //



For $s=0$ we call each cyclic kernel an action switch.

For $s=1$ " " " " " the payoff.

We represent an arbitrary deterministic \mathcal{V}_{ij}^s as:-



Such functional signed digraphs have not been investigated in the literature, but just from 3.4.7. , we see that we have a reasonably rich family of forms to adapt to a general environment \mathcal{M} .

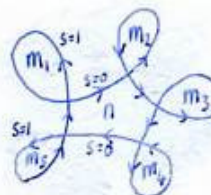
With a sufficiently large initial state space, we would expect to achieve near optimal memory depth, by combinatorial considerations of "cycle length" in each action class. However, such intuitive observations have not yet been rigorously proved. An analysis for a fixed finite state space of the asymptotic distribution over SOSA structures would appear extremely difficult to obtain, except for the proving of more general observations as above.

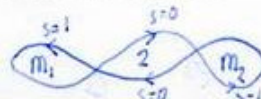
Definition 3.4.8.

i) An automaton is denoted linear if $n^L_{\underline{m}}$

m_i = memory depth of i^{th} action.

n = number of actions in switch.



In particular $2^L_{\underline{m}} \sim$ 



Previously we used $n^L_{\underline{r}} = n^L_{r1}$ so all memory depths were equal, in \mathcal{M}_n .

For arbitrary \mathcal{M} we cannot assert the existence of a limiting SOSA automaton but we can always achieve convergence using the rules $R_{\epsilon\delta}$. We then search for a SOSA automaton, yet we eventually converge to an unstable boundary if none exists. However, even if a SOSA automaton exists, under $R_{\epsilon\delta}$ rules we only have ϵ -optimality, so we might attain an unstable boundary. See 1.8. , rule 5).

ii) We use the notation $A_o(\mathcal{M}) \triangleq$ the set of SOSA automata existing in \mathcal{M} .

Theorem 3.4.9.

If $A_o(\mathcal{M}) \neq \emptyset$ then $\sigma_{ij}^s \rightarrow \nu_{ij}^s$ for some SOSA automaton,
under R_{ϵ} .

Proof.

We have seen that $A_0(M) \neq \emptyset$ for all symmetric media and we shall indicate later why it is plausible that this is true for an arbitrary M .

$$\text{Now } \Delta \sigma_{ij}^s(t) = \sigma_{ij}^s(t) \sum_{\alpha} w_i^{\alpha}(t) q_i^{\alpha,s} \sum_{h,\beta} \sigma_{ih}^s(t) \Delta_{\alpha\beta} \theta_{jh}(s) (\sigma_{ij}^{\beta,1} - q_h^{\beta,1}). \quad 12.$$

We take the equilibrium drift by fixing σ_{ij}^s in some small neighbourhood N_{ij}^s of the boundary V_{ij}^s , and then apply 1.12.9. to each $\theta_{(i,s)}$ associated with each σ_{ij}^s as discussed after 3.3.7. Thus we eliminate the transient effects, and consider the σ_{ij}^s -increment process w.r.t. the environmental equilibrium distribution in each state x_i which is +ve recurrent.

$$\text{Thus } \bar{\Delta} \sigma_{ij}^s = \sigma_{ij}^s \sum_{\alpha,h,\beta} \gamma_i^{\alpha} q_i^{\alpha,s} \sigma_{ih}^s \Delta_{\alpha\beta} \theta_{jh}(s) (q_j^{\beta,1} - q_h^{\beta,1}). \quad 13.$$

$$\text{now } \gamma_i^{\alpha} \sum_{\alpha} \sigma_{ih}^s q_i^{\alpha,s} \Delta_{\alpha\beta} = c_i \gamma_i^{\alpha} w_{\beta}^s \quad \text{where } w_{\beta}^s = \Pr(\text{in } E_{\beta} | s). \\ c_i = \sum_{\alpha} \sigma_{ih}^s q_i^{\alpha,s}.$$

$$\text{Thus } \bar{\Delta} \sigma_{ij}^s = \sigma_{ij}^s c_i \gamma_i^{\alpha} \sum_{\beta,h} w_{\beta}^s \sigma_{ih}^s \theta_{jh}(s) (q_j^{\beta,1} - q_h^{\beta,1}) = \sigma_{ij}^s c_i \gamma_i^{\alpha} \sum_{\beta,h} \sigma_{ih}^s w_{\beta}^s R_{jh}^{\beta}. \quad 14.$$

$$\text{now SOSA} \Rightarrow \sigma_{ij}^s \rightarrow 1 \quad \text{only if } u_j \text{ maximizes } \sum_{\beta} q_j^{\beta,1} w_{\beta}^s.$$

Then $\sum_{\beta} R_{jh}^{\beta} w_{\beta}^s > 0, \forall k \neq j \Rightarrow \bar{\Delta} \sigma_{ij}^s > 0$ so we have stability
and we can apply boundary learning theory
as in 3.3.6. and 3.4.4., using 1.12.9.

Clearly, if u_j does not maximize expected payoff at the next trial, yet $\sigma_{ij}^s \rightarrow 1$ then we get a contradiction, since reinforcement R_0 is optimal in a markov environment by 1.12.9.

//

Remarks 3.4.10.

$$\text{i) We can replace 13) by } \bar{\Delta} \sigma_{ij}^s = \sigma_{ij}^s c_i \sum_{\alpha,k} \mathcal{I}_{k,i}^{\alpha,s} R_{jk}^{\alpha}.$$

where $\mathcal{I}_{k,i}^{\alpha,s} = \Pr(\alpha, k \text{ at } t+1 \mid \gamma_i^{\beta}(t) \text{ and stimulus } s),$

and the stability of link σ_{ij}^s and hence $\theta_{(i,s)}$ depends only on:-

$$\text{sign}(\bar{\Delta}_{ij}^s) = \text{sign}\left(\sum_{k,j} R_{jk}^s \bar{I}_{ki}^{d,s}\right)$$

And boundary learning gives $\sigma_{ij}^s \rightarrow 1$ only if $\bar{\Delta}_{ij}^s \geq 0$ in arbitrarily small nbds of the boundary, under R_0 . Note that we need the whole of the apparatus of chapter 1 to assert this. We need now only test for deterministic stability to obtain stability of the stochastic process, if we use R_0 .

ii) We could by appropriate reinforcement achieve maximization over the next $r > 1$ trials, but it is our aim to keep the basic assumptions as natural as possible.

Conjecture 3.4.11.

a) For any $M(\Delta_{op}, q_{u_i}^{d,s})$ \exists a SOSA automaton.

b) The linear family $L_m^{(k_{ij})}$ covers M .

where m_i = memory depth of i^{th} action.

$k_{ij} = (k_{i1}, k_{i2})$ with k_{i1} = action we switch to at kernel under $s=0$ transition from u_i .

k_{i2} = action we switch to at kernel under $s=1$ transition from u_i .

so $k_{i2} = i$ under $m_i > 1$. (In many cases (k_{ij}) will be redundant.)

c) If $A_0(M) \neq \emptyset$ then $R(\mathcal{V}_{ij}^s) \geq \max_i e \cdot q_i$.

Remarks.

a) and b) go together since it is not possible to depart from $L_m^{(k_{ij})}$ and still obtain SOSA automata. If a) were false, then it would indicate that the graphs \mathcal{V}_{ij}^s are acting as a "strait-jacket" and that we should reformulate to give a larger family.

When $\Delta_{dd} \sim 1 \forall \alpha$ then we expect to be able to adjust $r_\alpha, \forall \alpha$ according to the equilibrium distribution e_α . Whilst if $\Delta_{dd} < 1/n$ with $n = \# \text{ media}$, then we adjust $k_{\alpha j}$ in an attempt to give a SOSA kernel. (k_{ij}) a kernel switching matrix, (m_i) a memory vector.

c) states that any limiting structure in a medium M will perform better than a singleton Π -cell. //

Examples 3.4.12.

a) For a cyclic 3-medium $\mathcal{M}_3^{\text{cyc}}(\Delta_{\alpha\beta}, q_{ij})$ with:-

$$\Delta_{\alpha\beta} = \begin{pmatrix} 1-\delta & \delta & 0 \\ 0 & 1-\delta & \delta \\ \delta & 0 & 1-\delta \end{pmatrix}$$

$$q_{i1}^{i1} = q, q_{j \neq i}^{i1} = p.$$

$$q_{i1}^{i0} = p, q_{j \neq i}^{i0} = q.$$

We have $\left(\begin{array}{c} \text{Diagram 1: A directed graph with 3 nodes and edges labeled } \delta, p, q. \\ \text{Diagram 2: A directed graph with 3 nodes and edges labeled } \delta, p, q. \end{array} \right) \cup \left(\begin{array}{c} \text{Diagram 3: A directed graph with 3 nodes and edges labeled } \delta, p, q. \\ \text{Diagram 4: A directed graph with 3 nodes and edges labeled } \delta, p, q. \end{array} \right)$ as a covering for all δ , and $p < \frac{1}{2}$.

and for $p=0$, ${}_3L_1$ is stable for $\delta < \frac{1}{2}$

$p=\frac{1}{2}$, ${}_3L_1$ is stable for $\delta < (3-\sqrt{5})/2 \approx 0.382$.

It is interesting that the threshold is not at $\delta = \frac{1}{3}$.

${}_3L_1$ is stable if $(1-\delta)(1-2q)(\delta(1-3q)+3\delta q-q) < 0$

${}_3L_1^*$ is stable if $(q-p)((3q-1)\delta^3-3q\delta^2+(3+pq)\delta-1) > 0$.

For $p > \frac{1}{2}$ we find that $({}_3L_1^* \vee {}_3L_1^A)$ is a covering for all δ .

For ${}_3L_1^*$ we have $(*) \sim (k_{ij}) = \begin{pmatrix} 3, 2 \\ 1, 3 \\ 2, 1 \end{pmatrix}$ the kernel switching matrix.

For ${}_3L_1^A$ we have $(A) = \begin{pmatrix} 3, 1 \\ 1, 2 \\ 2, 3 \end{pmatrix}$ so we cycle in the reverse direction to $\mathcal{M}_3^{\text{cyc}}$.

To obtain coverings we need to i) Calculate $\varphi_i^\alpha = \gamma_i^\alpha / \gamma_i$

$$\gamma_i^\alpha = \sum_{j, \beta} \sigma_{ji}^\beta q_{ij}^{\beta, \alpha} \Delta_{\beta\alpha} \gamma_j^\beta$$

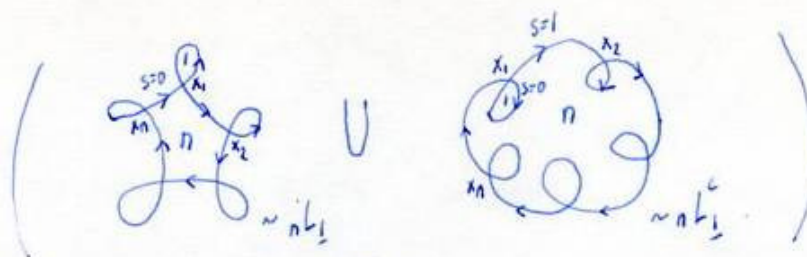
find ii) $\omega_\alpha^s = \sum_\beta \varphi_i^\beta \Delta_{\beta\alpha} q_{i\beta}^{\beta, s} / (\sum_\beta \varphi_i^\beta q_{i\beta}^{\beta, s})$

and then ensure $\sigma_{ij}^s = 1$ only if u_j maximizes $\sum_\beta \omega_\beta^s q_{ij}^{\beta, s}$.

We are just using 3.3.5. with φ_i^α as starting point.

The calculations are lengthy, with difficulties arising in the cyclic case since ω_β^s terms all interact, whilst in the symmetric \mathcal{M}_n we have ratios ω_i/ω_j determining structures.

b) For symmetric n-media \mathcal{M}_n , we have ${}_nL_1 \vee {}_nL_1^C$ as a covering.



with nL_1 stable for $n\delta \leq 1$.
 and nL_1^c stable for $n\delta \geq 1$.

//

Remarks. 3.4.13.

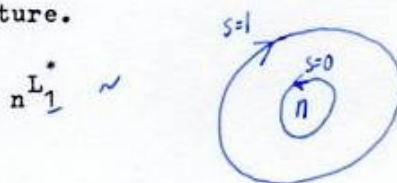
i) If we attempt to give a covering for $\mathcal{M}_n^{\text{cyc}}$, we get n^{th} order polynomial inequalities in δ which are most difficult to handle. Thus it has not been possible to show whether

$(nL_1^A \cup nL_1 \cup nL_1^*)$ is a covering for $\mathcal{M}_n^{\text{cyc}}$.

ii) In section 3.9. we shall see how hierarchies of structures can use \mathcal{M} coverings.

iii) In the case of $\mathcal{M}_3^{\text{cyc}}$ we actually obtain $3L_1^A \cup 3L_1^*$ as a covering for all δ and p . If the calculations are correct, then it is of interest that no $3L_1$ is required. However, for $\mathcal{M}_{n>3}^{\text{cyc}}$ we would still expect to require the full covering, as stated in i).

iv) The strategy of nL_1^* is of interest, since it has not appeared in the literature.



a) With $q > \frac{1}{2}$, then if we are in phase with $\mathcal{M}_n^{\text{cyc}}$ we shall receive $s=1$ with highest probability, so we switch through the actions x_i in the same order as the environment.

b) With $q > \frac{1}{2}$, then if we are out of phase we receive $s=0$ with highest probability so we switch through the actions in reverse order, since this is the quickest way to find the optimal action.

c) With $q < \frac{1}{2}$, we just have to avoid the one "bad" action so we attempt to be out of phase with $\mathcal{M}_n^{\text{cyc}}$, yet the nL_1^* is still SOSA.

v) Note that ${}_2L_1^A = {}_2L_1$, and that ${}_2L_1^*$ just cycles $x_1 \rightarrow x_2 \rightarrow x_1 \dots$ for all time, independent of stimulus. Thus only for $n > 2$ media do such automata exhibit useful adaptive behaviour.

3.5. Relationship between A_n and Likelihood Axis for a 2-Medium.

We shall consider in more detail the algorithm 3.3.5. and its relationship to our discrete structured automaton.

Lemma 3.5.1.

For 2-actions in M_2 with arbitrary δ and q ,

$$\exists \omega_{\max} \text{ and } \omega_{\min} \text{ s.t. } \omega_{\min} < \omega_i(t) < \omega_{\max} \text{ for } t \gg 0.$$

where $\omega_i(t) = \Pr(\text{in } E_1 \text{ at time } t \mid \omega_i(t-1), s(t-1), u(t-1))$.

Proof.

We call the set $Z_2 = \{\omega_i : \omega_{\min} < \omega_i < \omega_{\max}\}$ the operating zone for 2-actions.

Suppose $p, \delta < \frac{1}{2}$, then using the notation of 3.3.5.

$$\omega_i^*(t+1) = (\omega_i(t)(1-\delta)q + p\delta\omega_2(t)) / (\omega_i(t)q + \omega_2(t)q) = \omega_i(t) \text{ if } -$$

$$(\omega_i(t)(2q-1) + \omega_2(t)(1-2q+\delta) - p\delta = 0 \quad 15)$$

$$\text{and } \omega_{\max} = ((2q-1+\delta) + (2q-1)(1-2\delta + (\delta/(1-2q))^{1/2})) / 2(2q-1).$$

$$\text{and for } \delta \ll \frac{1}{2}, \quad \omega_{\max} \sim 1 - \delta q / (2q-1).$$

and similarly for $\omega_{\min} = 1 - \omega_{\max}$ by symmetry.

For $\delta = \frac{1}{2}$ we have $\omega_{\min} = \omega_{\max} = \frac{1}{2}$ so that Z_2 is just a singleton point.

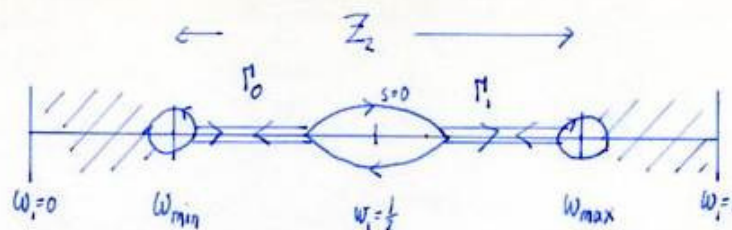
Whilst for $\delta > \frac{1}{2}, p < \frac{1}{2}$ we use $\omega_i^{s=0} = \omega_i$ to achieve bounds, so that

if $\delta = 1$ then $\omega_{\max} = 1 / (1 + (p/q)^{1/2})$ and in general for $\delta \in [\frac{1}{2}, 1]$ just interchange q, p in 15).

Finally if $p > \frac{1}{2}$ we interchange q and p in all the bounds so

$$\omega_{\max} \sim 1 - \delta p / (2p-1) \text{ if } \delta \ll \frac{1}{2}.$$

//



The Automaton $2^L_{\mathbb{F}}$ embedded in Z_2 .

Definitions 3.5.2.

- i) The set of $\{\underline{w}\}$ s.t. $\sum_{\alpha} w_{\alpha} (q_i^{\alpha} - q_j^{\alpha}) = 0$ for some $j \neq i$ and with $\sum_{\alpha} w_{\alpha} (q_i^{\alpha} - q_k^{\alpha}) > 0$ $k \neq j, i$, is called the threshold set.
- ii) The point $\underline{w} = \underline{\lambda}$ s.t. $\sum_{\alpha} \lambda_{\alpha} q_i^{\alpha} = \sum_{\alpha} \lambda_{\alpha} q_j^{\alpha} \forall i, j$ is called the action equilibrium point.
- iii) For an n -medium the operating zone Z_n is defined to be the \underline{w} absorbing set, for the optimal bayesian algorithm 3.3.5.

In the case of M_2 , $\underline{\lambda} = \underline{e} = \frac{1}{2} \underline{1}$ and for M_n and M_n^{cyc} we have $\underline{\lambda} = \underline{e} = \frac{1}{n} \underline{1}$. It is for this reason that the symmetric and cyclic environments were treated in 3.4.12., since it simplifies the analysis. Thus for M_2 it is optimal to take u_1 if and only if $w_1 > \frac{1}{2}$.

Lemma 3.5.3.

For 2-actions Z_2 spans \underline{e} , where $\exists w'_2 < e_2 < w''_2$ with $w'_2, w''_2 \in Z_2$.

Proof.

Suppose that \underline{e} is not spanned by Z_2 . Now by bayes' rule we must have $\overline{w}_2 = e_2$. So if $w_2 < e_2 \forall w_2 \in Z_2$ we have a contradiction, so that Z_2 spans \underline{e} as required, with $\underline{e} = \underline{e} \Delta$. //

It is unclear how this could be extended to give Z_n spans \underline{e} or whether we can obtain the stronger result $\underline{e} \in Z_n$.

Lemma 3.5.4.

For 2-actions:-

- a) A structured automaton may evolve to a singleton action class P_{α} iff $e_{\alpha}^s > \lambda_{\alpha}$ for $s=0,1$.
- b) If $\underline{\lambda}$ is not spanned by Z_2 then the maximum payoff

is achieved by using a single action for all time.

Proof.

a) If $e_\alpha^s > \lambda_\alpha$ then define structure:-

$$V_{ii}^s = 1, \quad s=0,1.$$



and this is clearly SOSA.

Suppose that $e_\alpha^s < \lambda_\alpha$ for $s=0$ say, but that \exists SOSA structure with equilibrium probabilities $p_i^s = \Pr(\text{in } E_\alpha \mid \text{in state } i)$.

$$\sum_i p_i^s = e_\alpha \quad \text{so} \quad \exists \quad p_{i_0}^s \leq e_\alpha \leq p_{i_1}^s.$$

However, by continuity we have either $\{p_{i_0}^s\}^{s=0} < \lambda_\alpha$ or $\{p_{i_1}^s\}^{s=0} < \lambda_\alpha$, 16)

and hence we have a contradiction as this gives

$u_{\beta \neq \alpha}$ is optimal at the next trial, so the SOSA structure has multiple action classes.

Note that we get 16) since if $\omega_i' > \omega_i''$ then $(\omega_i')^s > (\omega_i'')^s$ iff $\Delta_{ii} > \Delta_{2i}$.

b) This is immediate from the definition of λ .

//

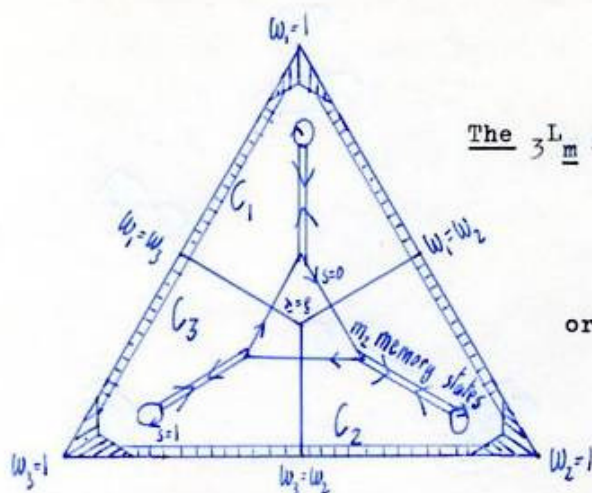
For a singleton action class P_α , we have $\bar{R}(v_{ij}^s) = e.g_\alpha$, as for a Π -cell. Then a) gives us conditions under which we have a stable limiting automaton in certain dynamic \mathcal{M} . The above lemmas have analogous forms for the n-action likelihood simplex but as yet, it is not possible to give such a precise characterization of their partitioning for automata learning.

3.6. The n-Medium Likelihood Simplex.

Definition 3.6.1.

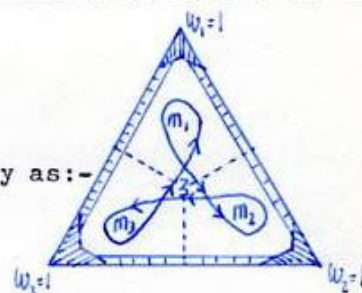
Let C_α be that part of the ω simplex in which u_α is the optimal action.

We first consider the symmetric medium \mathcal{M}_n and the relation between the SOSA property and the ω likelihood simplex.



The 3^L_m Automaton embedded in its w -Simplex.

or symbolically as:-



The boundary barriers in these diagrams are not intended to be exact representations.

Observations 3.6.2.

- The action switch forms around $\lambda = e$ for n^L_m , in M_n .
 - The Π_α payoff is in C_α .
 - The memory states of u_α are in C_α and they extend from the action switch towards $w_\alpha = 1$.
 - The w -operating zone Z_n is around $\lambda = e$ and it extends towards barriers around each $w_\alpha = 1$ and $w_\alpha = 0$.
- If $w_\alpha = 1$, $w_\beta^s = q_i^{s, \alpha} \Delta_{\alpha\beta} / q_i^{s, \beta} = \Delta_{\alpha\beta}$ under u_i .
- Thus $w_\alpha^s = \Delta_{\alpha\alpha}$ under any u_i , and this prevents Z_n from actually reaching $w_\alpha = 1$. Similarly for $w_\alpha = 0$.
- If $n\delta = 1$, $Z_n = \{w = (\frac{1}{n}, 1)\}$ whilst if $n\delta > 1$ we have $n^L_1^C$ stable around λ within $Z_n(M_n)$. d) and e) are the generalization of 3.5.1. from Z_2 to Z_n .

Remarks. 3.6.3.

- The λ point corresponds to the equilibrium strategy for "nature" if the process is considered as a game, since $\sum \lambda_\alpha (q_i^\alpha - q_j^\alpha) = 0 \forall j$. Then with n -actions, each of which is optimal in one of n -media, we have a "completely mixed" strategy. For \bar{w} -cell games we had $\sum \lambda'_\alpha (q_{\alpha i} - q_{\alpha j}) = 0 \forall j$ giving equilibrium values in 2.3.1.
- If $w = \lambda$ then all strategies u_i have the same expected payoff.
- If $(\bar{u}_i) = (\lambda'_i)$ then all strategies u_j for player 2 have the same expected payoff.

b) If $(q_i^\alpha - q_j^\beta) < \epsilon \forall i, \alpha, \beta$ then Z becomes an arbitrarily small region spanning \underline{e} . For M_2 we actually proved in 3.5.1 that $(w_{\max} - w_{\min}) \rightarrow 0$ as $(q - p) \rightarrow 0$, and the spanning of \underline{e} just follows from $\underline{w} = \underline{e}$.

When $q_i^\alpha = q_j^\beta \forall i, j, \alpha, \beta$ then $\underline{w} = \underline{e}$ and the statement above concerning the smallness of Z follows by the continuity implicit in Bayes' rule.

If $\underline{e} \neq \lambda$ then using 3.5.4, we can achieve singleton action class Γ_α iff $e_s^\alpha > \lambda_\alpha$ $s=0,1$ so that $\underline{e} \in \mathcal{C}_\alpha$. In this sense, the structure emerges from $\underline{w} = \underline{e}$ and as q_i^α separate we reach λ and an action switch forms, with memory extending out towards the vertices of the \underline{w} -simplex, always keeping within Z .

c) Just as the λ point gives nature an equilibrium strategy against the automaton, so with $q_i^\alpha = q_j^\beta \forall i, j, \alpha, \beta$ we find that "nature" has no influence on the automaton. Here we are representing the environmental actions by the states of M .

d) We have assumed that there is a unique optimal action for each environment.

i) If there are $n+r$ actions in an n -medium then take

u_i s.t. $q_{u_i}^\alpha \geq q_{u_j}^\alpha \forall j$ for each α .

ii) If there are $n-r$ actions in an n -medium, then we again

just take u_i s.t. $q_{u_i}^\alpha \geq q_{u_j}^\alpha \forall j$ for each α .

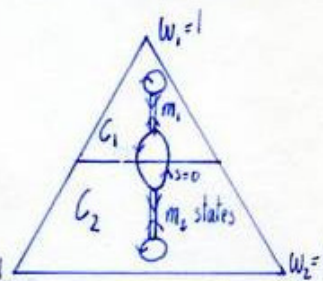
So with 2-actions in a 3-medium, we embed

2^L_M in the \underline{w} -simplex as:-

The threshold is the line $\sum_\alpha w_\alpha q_1^\alpha = \sum_\alpha w_\alpha q_2^\alpha$, and

if $q_{u_i}^2 = q_{u_i}^3$ for each u_i , then we effectively just have a 2-medium likelihood axis.

In general we obtain a hypersurface of dimension $n-2$ in an n -medium likelihood simplex, as the threshold set.



Summary 3.6.4.

We characterize the ω likelihood simplex by:-

i) The threshold set $\{\omega : \omega_{ij} = \omega_{ji} \geq \omega_{jk} \forall k \neq j, i\}$ for each u_i, u_j , and the λ point, where equality holds in the above in the case of a completely mixed equilibrium strategy.

ii) The environment equilibrium point e . $e_\alpha = \sum_{\beta} e_{\beta} \Delta_{\beta\alpha}$.

iii) The operating zone Z which is the ω absorbing set for the optimal bayesian algorithm 3.3.5.

We have now informally analysed the A_0 family, and so we turn now to networks of Π -cells $\mathcal{G}(\Theta)$.

3.7. Networks of Π -cells.

Instead of a unique action at each state we now have a Π -cell \otimes_k for each set of states x_i with $i \in \Pi_k$. There is now an alternative way of defining an underlying graph.

Let $\beta_{i,j}^{s_1, s_2}(t) = \text{Prob}(x_i(t) \rightarrow x_j(t+1) | s_1(t) \text{ and } u_k(t))$. Now since we asymptotically obtain the same family of structures, although with a different distribution, the greatly increased number of graphs ($2 \rightarrow 2^n$) does not seem justified. Indeed we shall see that we obtain all the desired results using $\sigma_{ij}^{s_1}$ as defined in 3.1.1.

The aim of this section is to show that even though we initially have all states possessing the same distribution over all actions, the actions automatically partition themselves through the tool of uniform learning. The actions interact through the equilibrium distribution γ_i^* of the structure. We first characterize the \mathcal{B}_0 family of structures in 3.7.1. and then relate this to $A_0(\mathcal{M}_n)$ in 3.7.2. and 3.7.3., whilst 3.7.4. considers the general $A_0(\mathcal{M})$.

As motivation for this abstract treatment we may consider the singleton Π -cells as entities which cluster together for their mutual benefit and then differentiate in their actions, in order to give structural adaptation in the markovian medium \mathcal{M} .

Theorem 3.7.1.

Let i be a +ve recurrent state in the limit $t \rightarrow \infty$ as defined in 3.2.1., then under R_0 and in environment M :-

$$\begin{aligned} & \text{a) } \sigma_{ij}^s \rightarrow 1 \\ & \text{and b) } \bar{u}_m^k \rightarrow 1 \end{aligned} \quad \text{with } i \in \Pi_k$$

only if $\max_1 \xi R$ w.r.t. the equilibrium distribution $\sum_i (v_{ij}^s)$, for $\sigma_{ij}^s \rightarrow 1$
and $\max_0 \xi R$ w.r.t. " " " $\sum_{\theta_k} \bar{u}_m^k$, " $\bar{u}_m^k \rightarrow 1$.

We denote $\{v_{ij}^s, y_m^k\}$ family as $B_0(M)$, where $\sigma_{ij}^s \rightarrow v_{ij}^s$ and $\bar{u}_m^k \rightarrow y_m^k$.

And where above we have used:-

$$\text{i) } \sum_{\theta_k} \tau_i^s = \sum_{i \in \Pi_k} \tau_i^s / \sum_{i \in \Pi_k} \tau_i^s$$

ii) $\max_r \xi R$ holds if we maximize the expected reward over the next r trials w.r.t the given environment equilibrium distribution.

Proof.

This is the basic convergence theorem for $B_0(M)$, and it is possible that a stable limit does not exist, as for $A_0(M)$. We shall assume here that $B_0(M) \neq \emptyset$.

Note that $\max_1 \xi R$ is precisely the SOSA property of 3.3.8. and that $\max_0 \xi R$ corresponds to the result 1.12.9. in which a \bar{u} -cell asymptotically takes that action which maximizes its average payoff.

We essentially mimic 3.4.9. and replace $q_i^{s,1}$ by $\sum_{j \in \Pi_k} \bar{u}_m^k q_j^{s,1}$ throughout. We could write out $\Delta \sigma_{ij}^s, \Delta \bar{u}_m^k$ for the actual process as in 3.3. but the equations are very similar, with $\bar{u}_m^k q_j^{s,1}$ replacing $q_i^{s,1}$.

Consider the process in an arbitrarily small neighbourhood of the boundary absorbing set, in the usual way. Let the equilibrium distribution for the fixed process be $\tau_i^s(\sigma_{ij}^s, \bar{u}_m^k)$.

$$\begin{aligned} \text{Now } \Delta \bar{u}_m^k &= \bar{u}_m^k \sum_{n \in \Pi_k} \tau_n^s \sum_i \pi_i^k \theta_{in}(\bar{u}_m^k) (q_m^{s,1} - q_i^{s,1}) \quad n \in \Pi_k \\ &= \bar{u}_m^k \left(\sum_{n \in \Pi_k} \tau_n^s \right) \sum_{s,i} \sum_{\theta_k} \pi_i^k \theta_{in}(\bar{u}_m^k) (q_m^{s,1} - q_i^{s,1}). \end{aligned}$$

17).

and $\Delta \sigma_{ij}^s = \sigma_{ij}^s \sum_{\alpha, k, \beta} \gamma_i^\alpha \left(\sum_m \bar{\pi}_m^k q_m^{\alpha s} \right) \sigma_{ik}^s \Delta_{\alpha\beta} \theta_{jk}(\sigma_{ij}^s) \left(\sum_m \bar{\pi}_m^a q_m^{\beta s} - \sum_m \bar{\pi}_m^b q_m^{\beta s} \right)$ 18).

where $i \in P_k, j \in P_a, k \in P_b$.

From 17) we get $\bar{\pi}_m^k \rightarrow 1$ only if $\sum_{\alpha} \gamma_i^\alpha R_{mi}^\alpha \geq 0 \quad \forall i$

and hence $\max_0 \mathcal{E} R$ w.r.t. γ_i^α by 1.12.9.

For σ_{ij}^s we consider $\{\bar{\pi}_m^k\}$ arbitrarily close to the boundary, in some nbd N_m^k of y_m^k say, so that $|\bar{\pi}_m^k - y_m^k| < \epsilon_m^k$.

Then $\Delta \sigma_{ij}^s = \zeta_i \sigma_{ij}^s \sum_{\alpha, k, \beta} \omega_\beta^s \sigma_{ik}^s \theta_{jk}(\bar{\pi}_m^a - \bar{\pi}_m^b) q_m^{\beta s}$ 19)

where $\zeta_i \omega_\beta^s = \sum_{\alpha} \gamma_i^\alpha \sum_m (\bar{\pi}_m^k q_m^{\alpha s}) \Delta_{\alpha\beta}$.

and $\zeta_i = \sum_{\alpha} \gamma_i^\alpha \left(\sum_m \bar{\pi}_m^k q_m^{\alpha s} \right)$.

and further $\text{sign}(\Delta \sigma_{ij}^s) > 0$ iff $\text{sign}(\sum_{\beta, k} \mathcal{I}_k^\beta R_{ab}^\beta) > 0$

with $k \in P_b, j \in P_a$.

and $R_{ab}^\beta = \sum_{k, m} \theta_{jk} q_m^{\beta s} (\bar{\pi}_m^a - \bar{\pi}_m^b)$

and from 19) we get $\max_1 \mathcal{E} R$ w.r.t. γ_i^α as we require

$$\sum_{\beta} R_{ab}^\beta \omega_\beta^s \geq 0 \quad \forall \theta_b.$$

Convergence follows as in 3.4.9., using a β -staircase construction and boundary learning theory in a markov medium, to give a form of 1.12.9. with the required equilibrium distribution for $y'(\theta)$.

//

Here v_{ij}^s has the SOSA property, but we have the additional constraint given on y_m^k , so it is unclear whether $A_0(m) = \theta_0(m)$. I shall indicate why I believe $\theta_0 \subset A_0$, but it is conceivable that $A_0 \not\subset \theta_0$ due to difficulties with discreteness.

We may find that u_j maximizes the expected reward at the next trial w.r.t. γ_i^α , yet the π -cell θ_k which executes this

action at the next trial is unstable with respect to θ_k^a .

However, no case of this has yet been found so the problem remains unresolved, and it seems clear that most structures will be well behaved.

The partitioning of actions is a stability result of a higher order than the boundary learning results we have so far considered. Consequently I shall present the analysis as a series of conjectures with an outlined "proof", since a truly rigorous proof would seem to be a deep excursion into probability theory.

Conjectures 3.7.2. and 3.7.3. are the results for $g'(\theta)$ which correspond to those for $g'(\omega)$ in 3.4.4. The corresponding extension of 3.4.9. is then 3.7.4. The \bar{u} -cells themselves effectively become "mixed actions" consisting of a distribution over the "pure" actions.

Conjecture 3.7.2.

In M_n under R_0 :-

$$\pi_j^k \rightarrow 1 \Rightarrow \pi_j^m \rightarrow 1 \quad \forall m \neq k.$$

Sketch.

This is a form of exclusion principle, so that only one \bar{u} -cell is allowed in each memory.

Suppose that u_j is the action which maximizes the expected reward at the next trial so that $\sum \omega_a^s (q_{u_j}^a - q_{u_k}^a) \geq 0 \quad \forall k$. Then it is optimal to choose θ_k s.t. $\pi_j^k \geq \pi_j^m \quad \forall m$. Hence the process is unstable if both θ_k and θ_m have u_j as the limiting action. Using R.W. theory it should be possible to prove that # axis crossings of $\pi_j^k = \pi_j^m$ is a.s. finite. Then after the last crossing θ_k say will always be the optimal \bar{u} -cell to pick if u_j is required.

Now use boundary learning to give $\Delta \sigma_j^s > 0$ for all $j \in P_k$, if θ_k gives $\max_j q_j^R$, when at the present trial we are in x_i .

For a rigorous formulation we need to consider $\gamma_i^s(\sigma_j^s, \pi_m^k)$,

so that for the instability in θ_k and $\theta_{m \neq k}$, when u_j is optimal we would expect $\gamma_{\theta_k}^{\alpha} \uparrow$ and $\gamma_{\theta_m}^{\alpha} \downarrow$ if $i_j^k \geq i_j^m \forall m$, $q_{u_j}^{a,1} \geq q_{u_k}^{a,1} \forall u_k$. //

Conjecture 3.7.3.

In M_n under R_0 :-

- nL_1 is in $B_0(M_n)$.
- $m < n L_1 \notin B_0(M_n)$.

Sketch.

Using 3.7.2. we have one i -cell on each memory, so that if $m < n L_1$ is stable we have θ_r say omitted. However, we know $m < n L_1$ is not SOSA by 3.4.2., and hence given $s=0$, it will be optimal to make a transition to some θ_r not "included in the structure". Only nL_1 is in $B_0(M_n)$, as it is SOSA and θ_r is stable by inspection of θ_r . //

This conjecture 3.7.3. is certainly true in itself, but it cannot be stated as a lemma since it relies on 3.7.2. for access to all θ_r which are required.

Now let nK_r^S be the complete digraph on rn states with n i_n -cells and each link $\sigma_{ij}^{s,b}$ having probability $1/rn$. Then we have 3.7.2. and 3.7.3. giving $nK_r^S \xrightarrow{M_n} nL_m$ as a mathematical

formalization of:- chaos $\xrightarrow{\text{environment}}$ structure adapted to its environment.

The next conjecture considers the problem of i -cell differentiation for general M .

Conjecture 3.7.4.

In M under R_0 :-

- $i_j^k \rightarrow 1 \Rightarrow i_j^m \rightarrow 1 \forall m \neq k$
- $B_0(M) \subset A_0(M)$.

Sketch.

In 3.7.2. we only had to consider the optimal actions, but for

general M , we must consider all u_i . It is possible that $\bar{u}_i^k > \bar{u}_i^m$ with u_i optimal at the next trial, yet θ_k need not be the optimal \bar{u} -cell to pick. θ_k may allocate probability $(1 - \bar{u}_i^k)$ to the action that $\min_1 \mathcal{E} R$, whilst θ_m allocates $(1 - \bar{u}_i^m)$ to the action which $\max_1 \mathcal{E} R$ if u_i is excluded.

However, we can still consider axis crossings of \mathcal{E} payoff (θ_k) against \mathcal{E} payoff (θ_m) at the next trial, when we wish to use u_i and obtain the instability as in 3.7.2. which prevents $\bar{u}_i^k, \bar{u}_i^m \rightarrow 1$.

Thus using boundary learning we can consider the equilibrium $\gamma_i^k(\theta_j^k, \bar{u}_m^k)$ and still obtain the exclusion principle, after proving that the # axis crossings ≥ 0 for \mathcal{E} payoff $(\theta_k) - \mathcal{E}$ payoff (θ_m) for each pair of \bar{u} -cells and each action u_i . (Construct a s/mg each side of the unstable axis.)

If M is static then for 3.2.2. we have that there exists a unique θ_k at $t = \infty$.

b) It is assumed that we have a structured automaton with every action u_i replaced by a \bar{u} -cell which incorporates all actions. If $B_0 \neq A_0$, let $v_{ij} \in B_0 \setminus A_0$ now v_{ij}^s is SOSA, so that the only possibility is that $\exists v_v$ s.t. this is not included in v_{ij}^s , yet would be in A_0 . However, just as in 3.7.3., the remaining θ_k have $\bar{u}_v^k > 0$, and these are used if u_v is optimal at the next trial, just as if we had a structured automaton with fixed actions evolving to A_0 .

Hence we have a contradiction unless $v_{ij}^s \in A_0$. Again it is a) which is difficult to prove, whilst b) is fairly immediate once we have a \bar{u} -cell available for each action in case it is required.

//

Networks of \bar{u} -cells would seem to act as a reasonable model for cellular differentiation. Waddington (1967) gives a useful basis for the biological framework.

It is perhaps useful to imagine the \bar{u} -cells of $\mathcal{E}'(\theta)$ as initially representing white light, with a spectrum of colours (actions).

Then as $q(\otimes)$ evolves, it adopts the colour most expedient in each environment, just as some animals, like the chameleon, change their colour according to their surroundings. Thus the white light within the e-m spectrum will partition itself into its component colours (frequencies) according to their uses within the environment.

The interaction between q and m enables the environment to "pull out" actions with certain properties from a "pool of actions", as they are required.

We can prove an analogue of $U\gamma = \gamma$ for structured automata, although the result has not been of use so far in generating bounds on absorption probabilities.

Theorem 3.7.5.

$$\begin{aligned} \text{a) } \lim_{n \rightarrow \infty} U^n \gamma_{v_{ij}}(\sigma_{ij}^s, \omega_i^s) &= \gamma_{v_{ij}}(\sigma_{ij}^s, \omega_i^s) \\ \text{b) } \Delta \gamma_{v_{ij}} &\equiv 0. \end{aligned}$$

where $\gamma_{v_{ij}}(\sigma_{ij}^s) = \begin{cases} 1 & \text{if } \sigma_{ij}^s = v_{ij}^s \\ 0 & \text{if } \sigma_{ij}^s \neq v_{ij}^s \end{cases}$ with γ_{ij}^s a deterministic SOSA automaton.
 $(\gamma_{ij}^s)^*$ " " " "

$$\text{and } U \gamma_{v_{ij}}(\sigma_{ij}^s(t), \omega_i^s(t)) = \sum (\gamma_{v_{ij}}(\sigma_{ij}^s(t+1), \omega_i^s(t+1)) | \sigma_{ij}^s(t), \omega_i^s(t)).$$

$$\Delta \gamma_{v_{ij}} = U \gamma_{v_{ij}} - \gamma_{v_{ij}}$$

$$\gamma_{v_{ij}}(\sigma_{ij}^s(0), \omega_i^s(0)) = \Pr(\sigma_{ij}^s(t) \rightarrow v_{ij}^s | \sigma_{ij}^s(0), \omega_i^s(t)).$$

$$\omega_i^s(t) = \Pr(\text{in } E_2(t) \text{ and } x_i(t) \text{ at time } t).$$

Proof.

$$\text{a) } \lim_{n \rightarrow \infty} U^n \gamma_{v_{ij}} = 1 \cdot \gamma_{v_{ij}} + 0(1 - \gamma_{v_{ij}}) = \sum (\gamma_{v_{ij}}(\sigma_{ij}^s(\infty), \omega_i^s(\infty)) | \mathcal{F}_0)$$

where it is assumed that $A_0(m) \neq \emptyset$ or else \mathbb{R}_{eg} is used.

$$\text{b) Clearly } \lim_{n \rightarrow \infty} U^n \gamma_{v_{ij}} = \gamma_{v_{ij}} \text{ and if } \Delta \gamma_{v_{ij}} \neq 0 \text{ we must get a}$$

$$\text{contradiction, for } \lim_{n \rightarrow \infty} U^n \gamma_{v_{ij}} = \lim_{n \rightarrow \infty} U^{n-1} \gamma_{v_{ij}} = U \lim_{n \rightarrow \infty} U^{n-1} \gamma_{v_{ij}} = U \gamma_{v_{ij}} = \gamma_{v_{ij}}.$$

We can interchange the U operator and the limit $n \rightarrow \infty$ since as in 1.12.1. $U^n \gamma$ converges uniformly to γ , even though

γ itself will have discontinuities, under R_0 reinforcement.

So we could verify that for each $\epsilon > 0 \exists N$ s.t. $n > N$ gives $|U^n \gamma - \gamma| < \epsilon$. //

Remarks 3.7.6.

i) If v_{ij}^s is the unique SOSA automaton then clearly $v_{ij}^s = 1$ under R_0 , if $\sigma_{ij}^s(0) = (0, i)$, but generally we may obtain the "same" structure in many different ways by relabelling the states in each V_{ij} . If we wish $\Pr(\sigma_{ij}^s \rightarrow {}_2L_1)$ we have to enumerate the graphs, with labelled states, that give ${}_2L_1$ from the initial state space X .

ii) Thus v_{ij}^s is the fixed functional solution of the operator U . Intuitively, this follows since in the limit $t \rightarrow \infty$, γ is an indicator function, assuming we have boundary convergence, with a probabilistic mass of $v_{ij}^s(\sigma_{ij}^s(0), \omega^s(0))$ on the set which has $v_{ij}^s(\sigma_{ij}^s(\infty), \omega^s(\infty))$ as unity. This is the fundamental equation of reinforcement learning, and if γ is defined appropriately, we shall achieve $\Delta \gamma = 0$, in all models using U.L. π -cells, when boundary convergence is assured.

Thus there would be a generalisation of 3.7.5. to the hierarchies $\mathcal{G}^n(\theta)$ of 3.9. and the games $\prod_{i=1}^n \mathcal{G}_i(\theta)$ of 3.10. So given any evolving stochastic automaton \mathcal{G} with absorbing structures defined in some environment \mathcal{M} , we can immediately form the functional equation which generates the absorption probabilities over the structures.

iii) This is essentially a discrete time result of probabilistic potential theory, where instead of a diffusion and Laplace's equation

$\nabla^2 \phi = 0$, we have stochastic difference equations and $\Delta \gamma = 0$. Thus any evolving stochastic automaton can be viewed as possessing a

certain potential at any time during its environmental adaptation. Such equations are treated in an abstract form by Meyer (1966). We can then consider a reward (penalty) stimulus as the reception

of a +ve(-ve) charge, say, with sign depending on boundary conditions.

iv) It is of interest to briefly note the relation between the

complementary pairs (reward, penalty), (survival, extinction), (+ve, -ve), (structure, chaos), (fixed, fluid) and the Ancient Chinese philosophy of Yin and Yang. Here all such fundamental duals are united through the Yin-Yang symbol,



where

the spots indicate that each side of any dual pair contains a little of the other. Gardener (1964) relates asymmetry in molecules to the possible origin of life, and considers the above symbol in relation to contemporary science.

Corollory 3.7.7.

$$\Delta \gamma_g(\omega) (\sigma_{ij}^1, \pi_j^1, \omega_i^1) = 0.$$

with $\gamma_g(y^*) = \delta_{gy^*}.$

Proof.

Again $\lim_{n \rightarrow \infty} \nu^n \gamma_g \rightarrow \gamma_g$ with $\gamma_g(y^*) = \delta_{gy^*}$. y, y^* are deterministic limiting structures

so $\lim_{n \rightarrow \infty} \nu^n \delta_g \rightarrow \gamma_g = \nu \delta_g$ by the same reasoning as in 3.7.5. //

We have reduced the analysis of structured automata, $\theta(\omega)$, and networks of π -cells, $\theta(\omega)$, evolving under θ in \mathcal{M} , to the investigation of:-

i) sign $(\Delta \sigma_{ij}^1)$. ii) sign $(\Delta \pi_m^k)$. iii) $\Delta \delta_g \equiv 0$.

in the usual notation.

3.8. π -cell Controllers and "Blueprint" Learning.

We may have several evolving automata and we wish to pick the one with the highest average payoff. We take a π_m -cell θ_c operating under θ which samples the $s_2(t)$ stimulus of the r^{th} automaton A_r , as its r^{th} action u_r at time t , $1 \leq r \leq m$.

θ_c is then called a π -cell controller.

Theorem 3.8.1.

For controller θ_c :-

a) $\exists r$ s.t. $\pi_r \rightarrow 1$.

$$b) \quad \bar{u}_k \uparrow 1 \quad \text{only if} \quad \overline{R}(\lim_{t \rightarrow \infty} A_k(t)) \geq \overline{R}(\lim_{t \rightarrow \infty} A_s(t)), \quad \forall s$$

where we assume $A_0(M) \neq \emptyset$ so that $\lim_{t \rightarrow \infty} A_r(t)$ has the SOSA property.

Proof.

We let $A_k(t)$ have transitions $\sigma_{ij}^s(k, t)$ and $\lim_{t \rightarrow \infty} A_k(t)$ is some SOSA automaton with transitions $\sigma_{ij}^s(k)$ and equilibrium distribution $\gamma_i^s(k)$. We could always ensure that $\lim_{t \rightarrow \infty} A_k(t)$ exists by using R_{csg} instead of R_0 as in 3.4.9.

We have Θ_0 in a markovian environment generated by $\prod_{r=1}^m A_r$. Now apply 1.12.9. to give $\bar{u}_i \uparrow 1$ iff $e_{\cdot i} > e_{\cdot j} \forall j \neq i$ where the equilibrium distribution e now becomes that for the A_r . Now combining this with 1.9.1. so that we need only consider the limiting probabilities as $t \rightarrow \infty$, we obtain:-

$$\bar{u}_r \uparrow 1 \quad \text{only if} \quad \sum_i \gamma_i^r(r) q_i^r \geq \sum_i \gamma_i^s(s) q_i^s \quad \forall s$$

If $\overline{R}(A_k(\infty)) = \overline{R}(A_h(\infty))$ for A_k and A_h say, then we still obtain convergence since we can relate the process to that of multiple optima in a static medium (1.6.4.) and use a ρ -staircase at each boundary to give $\# \text{upcrossings} < \infty$.

Clearly the $A_r(t)$ could be n -cell networks $g(\theta)$ and the result still follows if $\Theta_0(M) \neq \emptyset$, or else we use R_{csg} . //

This result is not elegant but a typical example of forced centralised learning, whilst we should prefer interactive behaviour between automata as the mechanism which increases their mutual adaptation.

If we are given a "blueprint" for any finite automaton then it is easy to formulate a learning process which converges to this structure. I am putting this result here because Suppes(1969) obtains a similar result with a different stimulus-response model. Using

R_0 the theorem is easily proved. Kieras (1976) considers the theorem of Suppes in great detail and extends the limiting family of automata.

Theorem 3.8.2.

Given $\{v_{ij}^{s_a}\}$ as a finite deterministic automaton where:-

$s_a(t)$ = input stimulus which gives transition $i \rightarrow j$ if $x_i(t)$
and $v_{ij}^{s_a} = 1$

Then if we put $q_{ij}^{s_a} > q_{ik}^{s_a}$ whenever $v_{ij}^{s_a} = 1$ then $q_{ij}^{s_a} \rightarrow v_{ij}^{s_a}$ under R .

Here $q_{ij}^{s_a} = \Pr(\sigma_{ij}^{s_a})$ is positively reinforced (reward), if we receive $s_a(t)$ when in $x_i(t)$ and $i \rightarrow j$.

Proof.

We apply the n-choice optimality theorem 1.6.4. for unstructured automata in a static medium, for each link $\sigma_{ij}^{s_a}$, and the result is immediate. Note that the $q_{ij}^{s_a}$ above differs from the usual $q_{ij}^{a,s}$ since here we must differentiate between the "blueprint" $\{v_{ij}^{s_a}\}$ and the learning mechanism. //

Remarks 3.8.3.

i). Essentially, we ensure that there is a unique absorbing state so that the evolving automaton is "always at the same potential" in that $v_{ij} = 1$ for the given "blueprint" (B.P.) if $\sigma_{ij}^{s_a}(0,1)$.

ii). The essence of \bar{n} -cell learning is that $M(A_{a\beta}, q_{a,i}^{a,s})$ is unknown so that learning is "blind", just as in our own evolution. If we allow the automaton to have access to $q_{a,i}^{a,s}$ we get a similar rather artificial "forced" result to 3.8.2. for it easy to obtain $q_{a,i}^{a,1} \rightarrow 1 \forall a, i$ through reinforcing $q_{a,i}^{a,1}$ whenever $u_i(t), E_a(t)$ hold at time t .

iii). We thus see that if we deviate from the \bar{n} -cell formulation through allowing the automaton to possess fuller knowledge of M then the results degenerate almost to remarks rather than leading to a fuller insight. So in representing structural adaptation through an evolving automaton G , we must ensure that we have no access to "global" information of M .

3.9 Hierarchical Automata.

We have considered static and markovian environments \mathcal{M} . In general Δ_{β} may consist of sub-environments, so that it would be natural to define a hierarchical automaton to adapt to such an environment. Tsetlin (1965) and Narendra and Viswanathan (1972) both used forms of 2-level systems for periodic media. The 1st level determines the period whilst the 2nd level operates in the selected environment. The $g(\theta)$ considered here is at present of conceptual rather than practical value and indicates how the theory of structural adaptation may develop in the future. This section contains work which is still rather speculative with many properties remaining to be investigated. We first motivate the theory through two simple examples.

Examples 3.9.1.

Here we take the structured automaton and substitute automata themselves as the actions executed in the states of the higher order automaton. Thus instead of a \bar{n} -cell being allocated to a set of states, we may replace it by a network of \bar{n} -cells. In these examples the automata are fixed and so we do not require a reinforcement stimulus; just a stimulus to determine state switching.

$$a) \quad \Delta_{\beta} = \begin{pmatrix} \Delta^1 & \Delta^e \\ \Delta^e & \Delta^2 \end{pmatrix} \quad \text{with} \quad \Delta^i = \begin{pmatrix} 1-\delta_i - \epsilon & \delta_i \\ \delta_i & 1-\delta_i - \epsilon \end{pmatrix} \quad \Delta^e = \begin{pmatrix} \epsilon & 0 \\ 0 & \epsilon \end{pmatrix}$$

and where $\epsilon \ll \delta_i \quad i=1,2$

and $\delta_1 < \frac{1}{2} \quad \delta_2 > \frac{1}{2}$

We denote the environment states as:- $E_{11}, E_{12}, E_{21}, E_{22}$.

Then for Δ^1 we use ${}_2L_F(1)$.

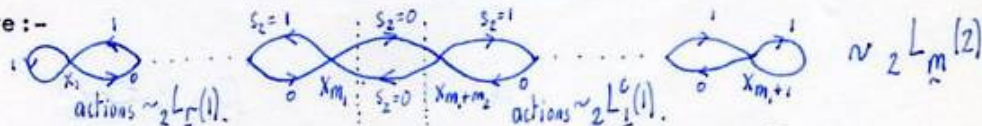
" " Δ^2 " " ${}_2L_1^C(1)$.

and for Δ use ${}_2L_m(2)$ with m chosen according to ϵ .

We define stimulus $s(t) = (s_1(t), s_2(t))$, where $s_j(t)$ is applied to the j^{th} -level in the hierarchy.

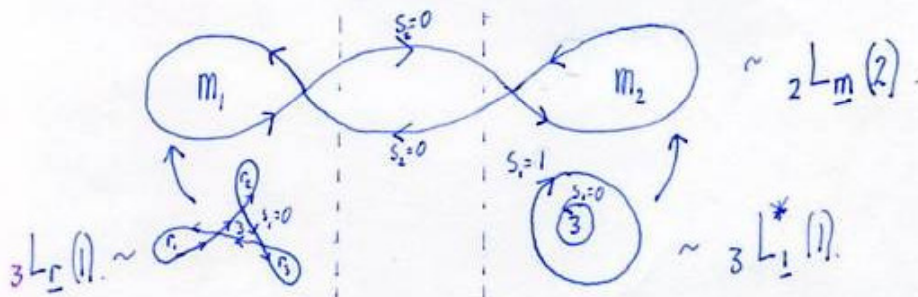
For level 1 automata in E_{ab} we put $q_{i=b}^{b,1} = q_{j \neq b}^{b,0}$ else $q_i^{a,s} = p = 1 - q$. And similarly for level 2 we put $q_{i=a}^{a,1} = q_{j \neq a}^{a,0}$ else $q_i^{a,s} = p$.

Thus we have:-



The actions on $2L_m(2)$ are replaced by automata $2L_r(1)$ and $2L_1^c(1)$ of the first level.

b) We can form a hierarchical automaton to operate in a cyclic 3-medium in a similar manner, where the δ parameter switches between say $\delta = \frac{1}{10}$ and $\delta = \frac{2}{3}$. Thus we have:-



The technique of 3.9.1. allows us to adapt to environments which split naturally into sub-environments. If we have $m(\Delta_{ap}, q^{a,s})$ and the parameters change, then we switch between the covering automata of $A_0(m)$ as in 3.9.1., where the coverings were investigated in 3.4.12.

Definitions 3.9.2.

- Let $g^n(\omega)$ be an n -level hierarchical automaton of actions, so that the r^{th} level has $g_i^{r-1}(\omega)$ as actions.
- Let $\sigma_{ij}^{s,r}(t)$ be the transition function for the r^{th} level at time t , with reinforcement under U.L. rules.
- Let $\Delta_{\alpha\beta}^m(r) = \Pr(E_\alpha \rightarrow E_\beta \text{ in } r^{\text{th}} \text{ environmental level when } n^{\text{th}} \text{ level is in } m_1, \text{ and } r^{\text{th}} \text{ level is in } m_{n-r+1}).$

iv) Let $\alpha_{q_i^m, s_{j,r}} = \Pr(\text{stimulus } s_{j,r} \text{ in level } r \text{ on using action } i \text{ in } E$
with $m_r = \alpha$, and $(r+s)$ th environment level in $m_{n-r+s+1}$.)

v) We have matrix stimulus $s_{.j} = (s_{1j}, s_{2j})$ with $s_{ij} = 1$
representing reward, as usual.

Here s_{1j} is used for r^{th} level reinforcement and s_{2j} is
used for r^{th} level transitions, each being independent with the
same stimulus probability. Thus we have a vector stimulus for each level.

We have defined a stimulus for each level rather than a single
global stimulus since the latter would require large memory in many
cases and would fail to operate if we required n^L , say, as actions.

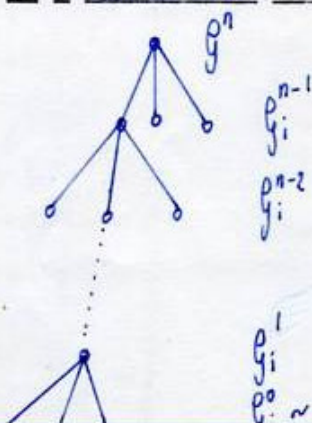
Using s_{ij} we get in phase with ∞ by successive levels,
rather than requiring the whole automaton to be in phase before
we benefit from the increased rewards, as in the case of global $s = (s_1, s_2)$.

vi) We define Γ_α^r as before so that g_α^{r-1} operates at time t
if and only if $x_i^r(t) \in \Gamma_\alpha^r$ with x_i^r an r -level state. Thus Γ_i^r is
the set of states in the r^{th} level which use action i . If each
level is held in Π -cells then Γ_k^r denotes the r -level $\Theta_k(r)$.

If we change $x_i^r(t) \rightarrow x_j^r(t+1)$ with $i \in \Gamma_\alpha^r$ and $j \in \Gamma_{\beta \neq \alpha}^r$
then the automaton action g_β^{r-1} will initially require time to
"acclimatize" after its period of inactivity.

We consider each level r , with its associated memory held in each Γ_α^r ,
reaching out to "tap" different environmental facets.

a) g^n as a Hierarchical Controller.



b) Fully Extended g^n .



action switch n

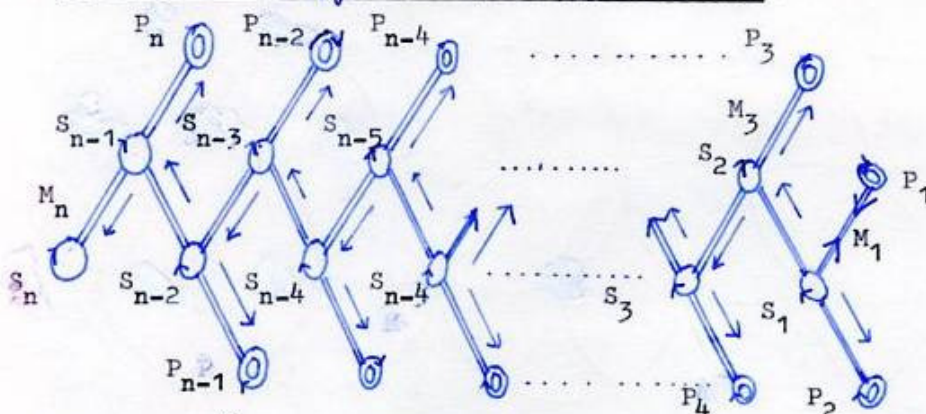
memory n

payoff n and action
switch $n-1$.

levels r s.t. $1 < r < n$.

payoff level 1.

c) Representation of Hierarchy in General State.



where $S_i = i^{\text{th}}$ action switch.

$M_i = i^{\text{th}}$ memory.

$P_i = i^{\text{th}}$ payoff.

} and $g^r = \bigcup_{i=1}^r (S_i, M_i, P_i)$ in terms of these "elements".

We have an abstract "expanding tongs" exploring the environment, which can be likened to a robot under hierarchical control, as considered by Albus and Evans (1976), in a practical context.

On the memory of each g_i^r we have automaton actions g_i^{r-1} .

Remarks 3.9.3.

- i) The environment transition matrix $\Delta_{\alpha\beta}^m$ as defined in 3.9.2. iii) differs from that used in the motivating examples 3.9.1.

In our new notation we write:-

$$\Delta_{\alpha\beta}^1 = \begin{pmatrix} 1-\delta_1 & \delta_1 \\ \delta_1 & 1-\delta_1 \end{pmatrix}, \quad \Delta_{\alpha\beta}^2 = \begin{pmatrix} 1-\delta_2 & \delta_2 \\ \delta_2 & 1-\delta_2 \end{pmatrix}$$

and $\Delta_{\alpha\beta} = \begin{pmatrix} 1-\epsilon & \epsilon \\ \epsilon & 1-\epsilon \end{pmatrix}$

with ϵ small: $0 < \epsilon \ll \delta_i$, $i=1,2$.

- ii) When we hold the actions at each level in \bar{n} -cells we can write:-

$$g^n(\otimes) = g(g(g(\dots g(\otimes)\dots))_{n\text{-times}})$$

So for g^n we have \bar{n} -cells of g^{n-1} . This is called \bar{n} -cellising the automaton actions at each level.

Then we obtain the results for hierarchies of \bar{n} -cells by the theory previously developed for $g(\otimes)$.

We define $\pi_i^k(r, t) = \text{Prob}(\text{use action } g_i^{r-1} \text{ at } r^{\text{th}} \text{ level on using } \bar{u}\text{-cell } \bar{u}_k(r), \text{ at time } t.)$

With the stimulus defined for each level we can immediately modify the action partitioning conjectures of 3.7. to the case of $g^n(\bar{u})$. Thus we can \bar{u} -cellise the actions at each level to obtain the full $g^n(\bar{u})$, or define a fixed action for each level state, to give $g^n(\bar{u})$. These are then the natural generalisations of \bar{u} -cell networks and structured automata, respectively, which are 1-level hierarchies g^1 . Clearly we could also construct mixed hierarchies so that we \bar{u} -cellise some levels and leave the remaining levels as fixed actions.

iii) The actions at the r^{th} level are denoted $u_i(r)$ and are held in the $\pi_i^k(r, t)$ if we \bar{u} -cellise. The $u_i(r)$ are g_i^{r-1} automata which receive stimuli s_{ir} with probability $q_{u_i}^{s_{ir}}$

$$\text{We put } \bar{R} = \sum_{r, i} \gamma_i^{\alpha} q_{u_i(r)}^{\alpha, s_{ir}} = \lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \sum_{r=1}^n s_{ir}(t) \quad \text{where } \gamma_i^{\alpha} = \text{eqm dsn for } g^n \text{ over } E_{\alpha} \text{ and states } i.$$

So the average reward is obtained by summing over each level r .

iv) For hierarchical automata, we reinforce transition links $\sigma_{ij}^{s_{ir}}$ only if we use the same action automaton g^r on ≥ 2 successive trials.

v) We note that biological systems themselves have a hierarchically stable organisation, so that each level is stable within the appropriate environmental level.

We can obtain a generalised SOSA limit theorem, with the stable limiting family defined as A_0^n for an n -level hierarchy $g^n(\bar{u})$.

Theorem 3.9.4.

If $A_0^n(n) \neq \emptyset$ then under A_0 , $\sigma_{ij}^{s_{ir}} \rightarrow \nu_{ij}^{s_{ir}}$ for some hierarchy that is SOSA at each level, for $g^n(\bar{u})$.

Proof.

As in 3.7.1. we can write the SOSA property as $\max_1 \bar{R}$ at each level with respect to the equilibrium distribution over the complete hierarchy.

We essentially need only apply 3.4.9. to each level r .

Thus $\sigma_{ij}^{s,r} \rightarrow 1$ only if $\max_1 \mathcal{E} R$ w.r.t. $\gamma_i^s(r)$ where:-

$\gamma_i^s(r) = \Pr$ (in E_β at r^{th} level and $m_{n-r+1+s}$ at $r+s$ level, and in automaton state i_{r+s} at level $r+s$, with $0 \leq s \leq n-r$.)

So the full equilibrium distribution is $\gamma_i^s(n)$, yet for level r we are only dependent on environment and automaton levels $s \geq r$ above us.

Now, as in 3.4.9., we take +ve recurrent state x_i^r at level r and consider $\sigma_{ij}^{s,r}$ in an arbitrarily small neighbourhood of the absorbing boundary. Then we fix $\{\sigma_{ij}^{s,r}\}$ and find $\Delta \sigma_{ij}^{s,r}$ in equilibrium, on considering the hierarchical automaton as a markov process. Note that each level r must have at least one +ve recurrent state, since automaton has finite state space.

Now we have $\Delta \sigma_{ij}^{s,r} > 0$ with $j \in P_j^r$ if and only if $u_j(r)$ gives $\max_1 \mathcal{E} R$ w.r.t. the equilibrium distribution $\gamma_i^s(r)$.

Then under \mathcal{R}_0 , we use 1.12.9., with e being replaced by $\gamma_i^s(r)$, and hence the limiting structure at each level is SOSA as required. //

We can prove a similar theorem to 3.7.1. for $g^n(\theta)$ which will give us the limiting family \mathcal{B}_0^n for an evolving n -level hierarchy.

Theorem 3.9.5.

Let i_r be a +ve recurrent state as $t \rightarrow \infty$ at level r , for each r .

Then under \mathcal{R}_0 and in \mathcal{M} :-

- i) $\sigma_{ij}^{s,r} \rightarrow 1$ with $i \in P_k^r$ for $g^n(\theta)$.
- ii) $\pi_m^k(r) \rightarrow 1$

only if $\max_1 \mathcal{E} R$ w.r.t the equilibrium distribution $\gamma_i^s(r)$ for $\sigma_{ij}^{s,r} \rightarrow 1$
and $\max_0 \mathcal{E} R$ " " " " $\gamma_i^s(r)$ " $\pi_m^k(r) \rightarrow 1$

where $\gamma_{\theta_n}^s(r) = \Pr$ (environment level $s \geq r$ is in m_{n-s+1} | we use $\theta_k(r)$, in eqm)

We denote $\{\sigma_{ij}^{s,r}, \pi_m^k(r)\}$ family as $\mathcal{B}_0^n(m)$, where $\sigma_{ij}^{s,r} \rightarrow \gamma_i^s(r)$ and $\pi_m^k(r) \rightarrow \gamma_m^k(r)$.

Proof.

This is just a restatement of 3.7.1 for n-levels, and the same method of proof applies in each level. Thus we extend B_r to B_r^n just as 3.9.4. gives us $A_r \rightarrow A_r^n$. //

Now at each level we conjecture that the \bar{u} -cellised actions partition themselves, so that 3.7.2. \rightarrow 3.7.4. have natural extensions. Thus on each r-level memory p_r^n we conjecture that only one g_r^{r-1} is allowed, and then that $B_r^n(m) \subset A_r^n(m)$, for each m .

An evolving $g^n(\otimes)$ will first adapt to the higher level environments before proceeding to "explore" the more minor environmental facets of the lower levels with the lower $g^r(\otimes)$ hierarchies. Thus, intuitively, we can imagine a "wave" of adaptation and action partitioning travelling down the "branches" of the hierarchical controller g^n .

3.10 Games between Structured Automata.

In the Russian Literature there are many papers which have been published on games between fixed deterministic automata, which were first formulated by Tsetlin (1963). Subsequent work has been published by:- Butrimenko (1967), Gersht (1967), Ginzburg and Stefanyuk (1970), Gurvich (1975), Kalinin (1965, 1966), Krinskii (1963, 1966), Takeuchi (1974), Tsertsvadze (1970), Tsetlin (1964, 1965, 1974), Vaisbord (1968), Varshavskii (1972) and Volkonskii (1965). Some related papers, including applications, are also to be found in the bibliography at the end of the thesis.

We define g_k^k as before, in 2.5.1., and then we obtain 3.10.1. quite easily. The game between \bar{u} -cell networks is denoted $\prod_{i=1}^n g_i^i(\otimes)$.

Theorem 3.10.1.

For n, structured automata, $g_k^i(\otimes)$ $1 \leq k \leq n$, evolving under R_r with game matrix g_k^k then:-

a) A pure strategy p^* is stable in the limit ~~if~~ and only if it is a Nash Point of g_k^k .

b) $\sigma_{ij}^s(r) \rightarrow 1$ only if the limiting structure has the SOSA property w.r.t τ_r , where τ_r is the equilibrium distribution for $\prod_{i=1}^n g_i(\omega)$.

Proof.

a) This does not assert uniqueness; only that any Nash Point is stable. We follow 2.5.1. and 3.2.2. in considering $\sigma_{ij}^s(r, t)$ near an absorbing barrier of pure strategies β^* for each automaton r at time t . Then $\Delta \sigma_{i, \beta^*}^s(r) > 0 \forall r, i$, in a sufficiently small nbd of the boundary $\{\sigma_{i, \beta^*}^s(r) = 1, \forall r, i\}$ iff β^* is a Nash strategy. Now we use the boundary learning theory of 1.7.5. under θ_1 to give convergence to β^* only if it is a Nash Point, so that each G_r^t just contains β_r^* in the limiting structure. It is unresolved whether other limiting stable structures exist if $\prod_{i=1}^n g_i^k$ has Nash Points. For a singleton \bar{u} -cell this is solved in 2.5.1.

b) This does not give the existence of limiting deterministic structures for it is possible that such limits do not always exist, as in the periodic trajectories of 2.3.

The result is obtained in the same way as the theorem 3.4.9., by writing out $\bar{\Delta \sigma_{ij}^s(r)}$, where the average is taken w.r.t. the equilibrium distribution τ_r . Then apply 1.12.9., with e replaced by τ_r to obtain the SOSA convergence constraint. //

Remarks 3.10.2.

i) Again, we could extend 3.10.1. to $\prod_{r=1}^n g_r(\omega)$ on defining β_k^r for the product gaming environment, with strategy vector β .
 ii) We would like the limiting structures to resemble the automata of Krinskii (1963), which have a deterministic "circular" behaviour, but it is unclear whether they have the SOSA property in a gaming environment. Thus although the linear automata are ideally suited to learning in markovian media, we need cyclic behaviour for automata games, with the digraph resembling the $\bar{Q}_k = \text{const}$ trajectories of section 2.3.

iii) A structural analysis of $\prod_{i=1}^n g_i(\omega)$ still remains to be done, but

it seems most unlikely that they will achieve worse results than the gaming singleton \mathcal{H} -cells, just as we proved 3.2.2. for games against "nature".

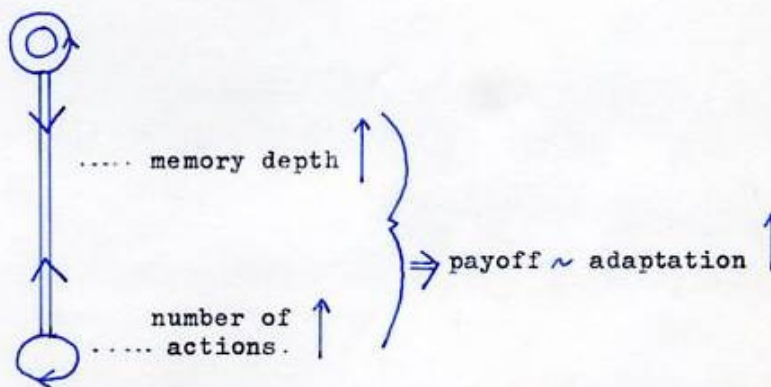
3.11. Concluding Remarks.

3.11.1. There are essentially two means by which automata can increase their average payoff in the formulation of this chapter.

a) Memory extension. ("Tree" \uparrow)

b) Action switch enlargement. ("Kernel" \uparrow)

We represent this diagrammatically as:-



We saw in 3.8. that efforts to improve on suboptimality using controllers, were rather artificial. It would be preferable to postulate community interaction so that automata learn by their mutual errors and successes.

For a) we need to "fluidise the structure" and increase the adaptation through observing other automata with larger memories. In addition, we could allow automata to possess individual curiosity so that some may act as pioneer automata in trying out new structures. A preliminary model for community behaviour is presented as 3.11.2.

For b) new actions may arise which are expedient in some environment and which automata can now incorporate in their structure. The cities of Paris and Rome have new "actions" in their "kernel of operation"; La Défense and the E.U.R. respectively. Bacon (1974) adopts such an evolutionary approach in his qualitative study of cities.

Model 3.11.2. ("The Sheep Effect".)

Here we are not concerned explicitly with the automaton structure but just in their relative expediency in an environment. In this very simple model we assume perfect mutual observation between automata, rather than the usual observation with uncertainty as in n -cell theory.

- i) We define states $x_i = 1, 2, \dots, m$ with state 1 "best" and then the utility monotonically decreases to give state m "worst".
- ii) We define automata a_1, \dots, a_n .
- iii) If an automata observes another automata doing "better" than itself, then it jumps with probability $1/m$ to some x_1 .
- iv) Let $A_{j_1, \dots, j_n}(s) = \Pr(\text{all automata are eventually absorbed in state } s \text{ where initially } a_i \text{ is in state } j_i.)$
- v) Clearly $A_{j_1}(s)$ is only a function of $\min_i j_i$ and the number which attain this.

So define $A_{jr} = \Pr(\text{all automata eventually absorbed in } x=1, \text{ given } r \text{ start in } x=j > 1.)$

$$\begin{aligned} \text{Thus } A_{jr} &= A_{jr} \left(\frac{m-j}{m}\right)^{n-r} + \sum_{s=1}^{n-r} A_{j,r+s} \binom{n-r}{s} \left(\frac{1}{m}\right)^s \left(\frac{m-j}{m}\right)^{n-r-s} \\ &\quad + \sum_{k < j} \sum_{s=1}^{n-r} A_{ks} \binom{n-r}{s} \left(\frac{1}{m}\right)^s \left(\frac{m-k}{m}\right)^{n-r-s}. \\ \text{or } A_{jr} (m^{n-r} - (m-j)^{n-r}) &= \sum_{s=1}^{n-r} A_{j,r+s} \binom{n-r}{s} (m-j)^{n-r-s} \\ &\quad + \sum_{k < j} \sum_{s=1}^{n-r} A_{ks} \binom{n-r}{s} (m-k)^{n-r-s}. \end{aligned}$$

And if $A_{1n} = 1$, $A_{jn} = 0$, $j \neq 1$, then:-

$$A_{jr} \sim 1 - \left(\frac{1}{j}\right)^{n-r} \text{ for } j > 1 \text{ and } n \text{ large. So that the probability}$$

of non-optimal absorption becomes geometrically small as n increases.

- vi) Thus we have a bootstrapping effect in that the "pioneer" automata "pull up" the trailing "sheep" automata, yet which themselves may

change rôles and act as the new pioneers. Since the automata are only guided by each other, there is always a finite probability that they will all become absorbed together in the same sub-optimal state.

//

3.11.3. As motivation for extending the π -cell theory to hierarchies of automata g^n in 3.9. , we provide the following brief intuitive outline.

Initially we suppose there are singleton π -cells which proceed to cluster for their mutual benefit and subsequently differentiate in their rôles as in 3.7. Now to form g^{r+1} we require a mechanism of collective behaviour between the g_i^r which eventually gives rise to a critical point and the gelation of structural links, followed by rôle differentiation. (Whittle 1971). In this way automata can attain adaptation in higher environmental levels of \mathcal{M} , within a single entity g^n . The automaton g^n acts as a hierarchical controller with autonomous behaviour within each lower level, and which we could now embed in the appropriate ω -likelihood simplex, as in 3.6., for g^i .

Such a process of clustering and subsequent rôle differentiation appears frequently in nature. In particular, we may proceed up the evolutionary tree from the single-celled amoeba, eventually obtaining the multicellular mammals and the cultural society of man. The actual mechanism of reinforcement here is natural selection over many generations with the genetic encoding of this information giving embryonic cell differentiation and subsequent learning over the lifetime of each individual.

3.11.4. At present, most models in psychological learning theory are based on the unstructured automaton, yet semantic ideas express the brain as a structure evolving through the reception of environmental stimuli, as in the simple qualitative models of De Bono (1971).

So as a simple conceptual model, we have investigated the π -cell as a basis for structural learning, with its family of learning functions.

3.11.5. Probability expresses our environmental uncertainty so that through learning, our ideas become rigid, and we act according to that structure.

Research requires us to "fluidise" this structure to enable us to "extend the trees" and "enlarge the kernels", in our "blind" attempts at environmental understanding and adaptation.



Bibliography

4 Bibliography.

..... even as the sand that is upon the sea shore
in multitude.

Joshua 11 v 4.

4 Bibliography.

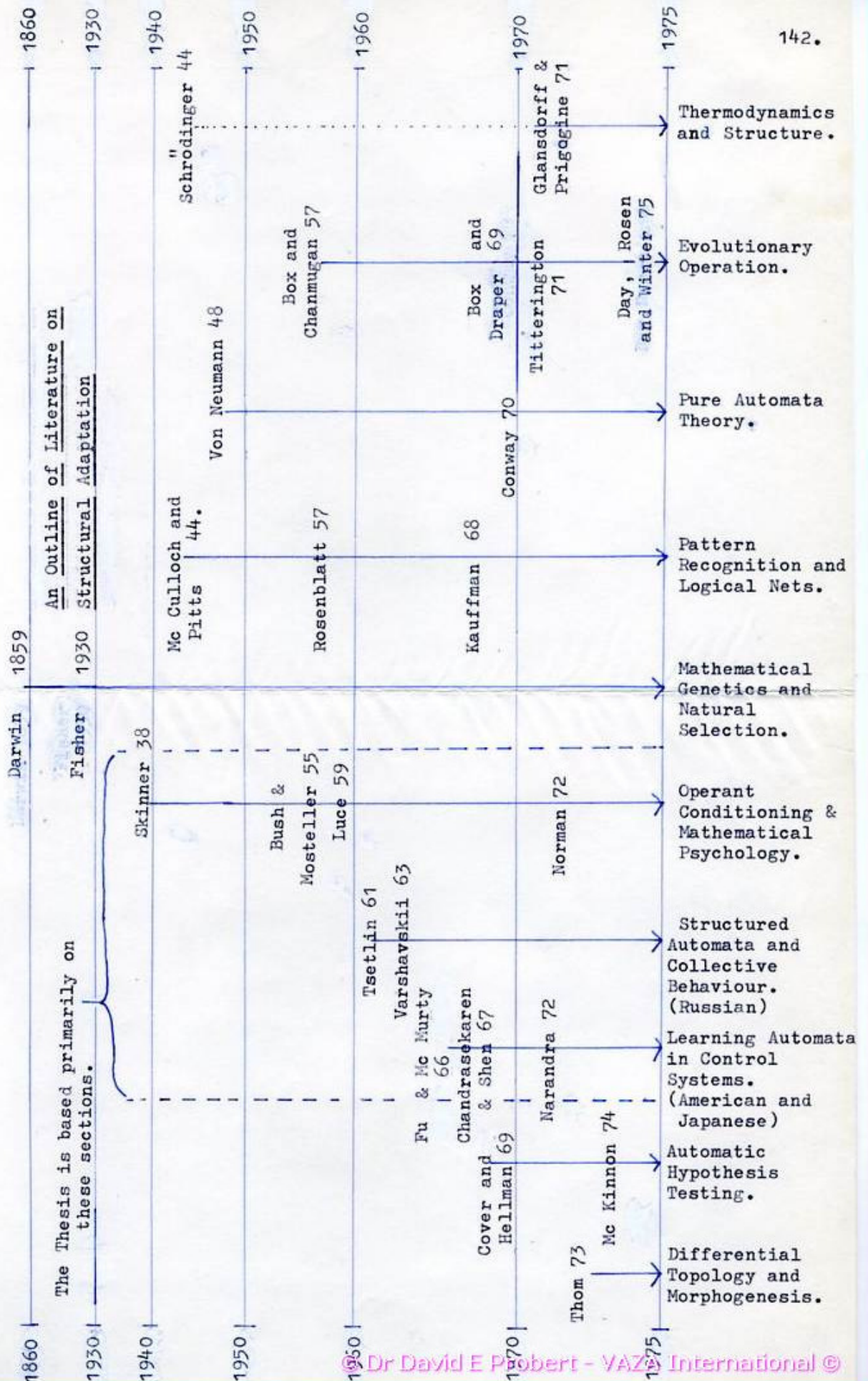
A comprehensive bibliography is given here including some papers not actually cited in the text, although they all helped in the development of the underlying theme of the thesis.

Much of this literature was gathered before the survey paper of Narendra and Thathachar (1974), yet the field covered is remarkably similar. This includes the work of mathematical psychologists, American and Japanese control engineers and Russian automaton theorists.

Prior to 1973, the literature was virtually partitioned into the three categories above, but it is now clear that all approaches are merging to form a common framework of learning automata and structural adaptation. More abstractly, we could view the field within Thom's (1975) theory of morphogenesis.

Abbreviations:-

- a) A.R.C. — Automation and Remote Control.
(Translation of Avtomatika i Telemekhanika.)
- b) P.I.T. — Problems of Information Transmission.
(Translation of Problemy Peredachi Informatsii.)
- c) S.S.C. — I.E.E.E. Transactions on Systems, Science and Cybernetics.
(1965-1970)
- d) S.M.C. — I.E.E.E. Transactions on Systems, Man and Cybernetics.
(1971-)
- e) A.C. — I.E.E.E. Transactions on Automatic Control.
- f) I.T. — I.E.E.E. Transactions on Information Theory.
- g) J. of Math Psych — Journal of Mathematical Psychology.
- h) J. of Theor Biol — Journal of Theoretical Biology.



- 1) Accardi L. (1974). Description of Probabilistic Evolution
Dependent on the Past. A.R.C. 4 50-61.
- 2) Asai K and Kilajima S (1971). A Method of Optimising Control of
Multimodal Systems using Fuzzy Automata. Information Sciences 3 ~~343~~-353.
- 3) Albus J and Evans J.M. (1976). Robot Systems. Scientific American 234 2.
- 4) Aleksander I (1971) Electronics for Intelligent Machines.
New Scientist 11th March.
- 5) Aiserman M.A et Al (1971). Logic, Automata and Algorithms.
Academic Press.
- 6) Aiserman M.A. (1975). A Man and a Collective as Elements of a
Control System. A.R.C. 36 No 5 Pt 1.
- 7) Agasandyan G.A. (1969) Some Aspects of the Theory of Automata
with Variable Structure. P.I.T. 5 No 1 71-78.
- 8) Arbib M.A. (1969) Theories of Abstract Automata. Prentice-Hall Inc.
- 9) Atkinson R.C. (Ed 1964) Studies in Mathematical Psychology.
Stanford University Press.
- 10) Atkinson R.C. et Al (1974). Contemporary Developments in Mathematical
Psychology. Freeman San Francisco.
- 11) Bacon E.N. (1974). The Design of Cities. Thames and Hudson.
- 12) Barbour A. (1973). The Principle of the Diffusion of Arbitrary
Constants. J. of Applied Probability. 9 519-541.
- 13) Bellman R (1956). A Problem in the Sequential Design of Experiments.
Sankhya. 16 221-229.
- 14) Bellman R (1975). Communication and Mind. Biosciences. 27 p347-
- 15) De Bono E (1971). The Mechanism of Mind. Penguin Books.
- 16) Bernal J.D. (1967). The Origin of Life. W.Clowes and Son.
- 17) Borovikov V.A. and Bryzgalov V.I. Simplest Symmetrical Game of
Many Automata. A.R.C. 26 No 4.
- 18) Breiman L (1968). Probability. Addison Wesley.
- 19) Bush R.R. and Mosteller F. (1955) Stochastic Models for Learning. Wiley.
- 20) Butrimenko A.V. (1967). Games of Automata Possessing Different
Activities. P.I.T. 3 No 4 81-88.

- 21) Buyakas V.I. (1974) Problem of Controlling Collective Behaviour.
A.R.C. 35
- 22) Box G. and Jenkins G (1962). Some Statistical Aspects of Adaptive
Optimisation and Control. J.R.S.S. B. 24 297-343.
- 23) Box G and Draper N. (1969). Evolutionary Operation. Wiley.
- 24) Cairns-Smith (1971). The Life Puzzle. Oliver and Boyd.
- 25) Calvin (1969) Chemical Evolution. Oxford University Press.
- 26) Chaikovskii Yu.V. (1970) A Group of Gain Comparing Automata.
P.I.T. 6 No 4 56-65.
- 27) Chandrasekaren B and Shen D.W.C. (1967) Adaptation of Stochastic
Automata in Nonstationary Environments. Proc of National Electronics
Conference.
- 28) Chandrasekaren B and Shen D.W.C. (1968) On Expediency and Convergence
in Variable Structure Automata. S.S.C. 4 No 1.
- 29) Chandrasekaren B and Shen D.W.C. (1969) Stochastic Automata Games.
S.S.C. 5 No 2.
- 30) Chandrasekaren B. (1970). Finite-Memory Testing. A Critique. I.T. 16.
- 31) Chandrasekaren B. et Al. (1974). Artificial Intelligence.
S.M.C. 4 p88-103.
- 32) Chien Y.T. and Fu K.S. (1967a) Learning in Nonstationary Environments
using Dynamic Stochastic Approximation. Proc 5th Annual Allerton
Conference on Circuit and System Theory.
- 33) Chien Y.T. and Fu K.S. (1967b) On Bayesian Learning and Stochastic
Approximation S.S.C. 3 28-38.
- 34) Chien Y.T. and Fu K.S. (1969) Stochastic Learning of Time-Varying
Parameters in Random Environments. S.S.C. 5 No 3 237-246.
- 35) Codd E.F. (1968) Cellular Automata. Academic Press.
- 36) Coxeter et Al (1953) Uniform Polyhedra. Phil Trans A. 401-450.
- 37) Day R and Groves T. Eds (1975) Adaptive Economic Models. Academic Press.
- 38) Dodson M.M. (1976) Darwin's Law of Natural Selection and Thom's
Theory of Catasrophes. Math Biosciences. 28
- 39) Doob (1953). Stochastic Processes. Wiley.

- 40) Dubins L.E. and Savage L.J. (1965). How to Gamble if you Must.
Mc Graw Hill.
- 41) Estes W. (1970). Learning Theory and Mental Development. Academic Press.
- 42) Feichtinger G. (1970). Lernprozesse in Stochastischen Automaten.
Springer-Verlag.
- 43) Fisher R.A. (1930). The Genetical Theory of Natural Selection. Oxford.
- 44) Fogel L. et Al (1966) Artificial Intelligence through Simulated
Evolution. Wiley.
- 45) Freedman D. (1973). Another Note on the Borel-Cantelli Lemma and
Strong Law. Annals of Probability 1 No 6 910-925.
- 46) Fu K.S. and Wee W.G. (1967). A Formulation of Fuzzy Automata and
its Application to Learning Systems. Proc 5th Annual Allerton
Conference on Circuit and Systems Theory.
- 47) Fu K.S. and Li T.J. (1969) Formulation of Learning Automata and
Automata Games. Information Science 1 237-256.
- 48) Fu K.S. (1970). Learning Control Systems : Review and Outlook. A.C. 15 2.
- 49) Gardener M. (1970) The Ambedextrous Universe. Pelican Books.
- 50) Gardener M. (1971). Mathematical Games-"Life". Sci American 224 2 112-117.
- 51) Gersht A.M. (1967). Games of Continuous Automata. P.I.T. 3 No1 50-56.
- 52) Ginzburg S.L. and Stefanyuk V.L. (1970). n-Automata Games with
Single Nash Point. A.R.C. 31 Pt 2 1264-1272.
- 53) Glorioso R.M. and Grueneich G.R. (1971) A Training Algorithm for
Systems described by Stochastic Transition Matrices. S.M.C. 1 86-87.
- 54) Glushkov V.M. (1966). Introduction to Cybernetics. Academic Press.
- 55) Goel N.S. et Al (1971). Non-Linear Models of Interacting Populations.
- 56) Golovchenko V.B. (1974). Self-Organisation of a Collective of
Probabilistic Automata with the Two Simplest Motives of Behaviour.
A.R.C. 35 No 4 151-156.
- 57) Gurvich E.T. (1975). Method for Asymtotic Investigation of Games
between Automata. A.R.C. 36 No 2 80-94.
- 58) Harary F. et Al (1965). Structural Models. Wiley.
- 59) Hellman M.E. and Cover T.M. (1970a). The Two-Armed Bandit Problem
with Time Invariant Finite Memory. I.T. 16 No 2.

- 60) Glansdorff and Prigogine. (1971). Thermodynamic Theory of Structure and Fluctuations. Wiley Interscience.
- 61) Hellman M.E. and Cover T.M. (1970b) Learning with Finite Memory. Annals of Math Stat 41 No 3 765-782.
- 61) " " " (1970c) Comments on Automata in Random Media. P.I.T. 6 No 2 21-30.
- 62) " " " (1971) On Memory Saved by Randomisation. Annals of Math Stat 42 1075-1078.
- 63) Hellman M.E. (1972) The Effects of Randomisation on Finite Memory Decision Schemes. I.T. 18 No 4.
- 64) Hirschler D. and Cover T (1975) A Finite Memory Test of the Irrationality of the Parameter of a Coin. Annals of Statistics 13 No 4.
- 65) Holman E.W. (1969) Asymptotic Properties of Monotonic Learning Models. J. of Math Psych. 6 456-469.
- 66) Iosifescu M. and Theodorescu (1969) Random Processes. Springer-Verlag.
- 67) Jarvis R.A. (1975) Adaptive Global Search by the Process of Competitive Evolution. S.M.C. 5 297-311.
- 68) Jones P.W. (1975) Two-Armed Bandit. Miscellanea, Biometrika 62 No 2
- 69) Kalinin D.I. and Epshtein I.M. (1965) A Note on an Automaton Game with Partner using Correlated Mixed Strategy. A.R.C. 26 No 11.
- 70) Kalinin D.I. and Rotvain I.M. (1966) Some Asymptotic Estimates for Games of Automata in Distribution. 27 No 4 642-644.
- 71) Kanal L. (1962) A Functional Equation Analysis of Two Learning Models. Psychometrika. 27 89-109.
- 72) Kandelaki N.P. and Tsertsvadze G.N. (1966) Behaviour of Certain Classes of Stochastic Automata in Random Media. A.R.C. 27 No 6.
- 73) Karlin S. (1959) Matrix Games, Programming and Mathematical Economics. Vol 1. Addison Wesley.
- 74) Kauffman S.A. (1969). Metabolic Stability and Epigenesis in Randomly Constructed Genetic Nets. J. of Theor Biol. 22 437-467.
- 75) Keilson J and Wishart D. (1964). A Central Limit Theorem for Processes defined on a Finite Markov Chain. Proc Cam Phil Soc 60 p547-567.

- 76) Keilson J and Wishart D. (1965). Boundary Problems for Additive Processes defined on a Finite Markov Chain. Proc Cam Phil Soc 61 173-190.
- 77) Keilson J and Wishart D. (1967). Addenda to Processes defined on a Finite Markov Chain. Proc Cam Phil Soc. 63 187-193.
- 78) Kieras D.E. (1976). Finite Automata and S-R Models. J. of Math Psych. 13 127-147.
- 79) Koganov A.V. (1973). Automata that Distinguish Random Media. P.I.T. 9 No 2 68-80.
- 80) Krinskii V.I. (1963). An Asymptotically Optimum Automaton. Biofizika 6 No 6 484-487. (In Russian)
- 81) Krinskii V.I. (1966). Zero-Sum Games for Two Asymptotically Optimal Sequences of Automata. P.I.T. 1 No 2 43-53.
- 82) Krylov V.Y. (1963) On One Automaton that is asymptotically Optimal in a Random Media. A.R.C. 24 1226-1228.
- 83) Kushner H.J. (1967) Stochastic Stability and Control. Academic Press.
- 84) Lakshmivarahan S and Thathachar M.L. (1972) Optimal Non-Linear Reinforcement Schemes for Stochastic Automata. Information Sciences 4 121-128.
- 85) Lakshmivarahan S. and Thathachar M.L. (1973). Absolutely Expedient Learning Algorithms for Stochastic Automata. S.M.C. 3
- 86) Linkin V. (1972) An Adaptive Algorithm for Determining Variations of the Characteristics of an Observed Random Process. P.I.T. 8 40-45.
- 87) Luce R.D. (1959) Individual Choice Behaviour.
- 88) Lyubihik L.M. and Poznyak A.S. (1974) Learning Automata in Stochastic Plant Control Problems. A.R.C. 35 95-109.
- 89) Mc Culloch W and Pitts W. (1943) A Logical Calculus of the Ideas Immanent in Nervous Activity. Bull Math Biophysics. 5 113-133.
- 90) Mc Kinnon K.I.M. (1974) Optimal Finite-State Discrimination Procedures. Phd Thesis Cambridge.
- 91) Mc Laren R.W. (1966) A Stochastic Automaton Model for the Synthesis of Learning Systems. S.S.C. 2 No 2.
- 92) Mc Murty C.J. and Fu K.S. (1966) A Variable Structure Automaton used as a Multimodal Searching Technique. A.C. 11 No 3.

- 93) Malishevski A.V. and Tennisberg P.D. (1969) One Class of Games connected with Models of Collective Behaviour. A.R.C. 30 1828-1837.
- 94) Marley A.A.J. (1967) Abstract One-Parameter Families of Commuting Learning Operators. J. of Math Psych. 4 414-429.
- 95) Mason I.G. (1973) An Optimal Learning Algorithm for S-Model Environments. A.C. 18 493-496.
- 96) Maynard-Smith J. (1974) The Theory of Games and the Evolution of Animal Conflicts. J. of Theor Biol. 47 p207-
- 97) Meleshina M.V. (1969) Automaton Model of Client Interaction Organisation in Queuing Systems with Waiting. A.R.C. 30.1 782-783.
- 98) Mendal J.M. and Fu K.S. (1970) Adaptive, Learning and Pattern Recognition Systems Theory and Applications. Academic Press.
- 99) Mendal J.M. (1973) Reinforcement Learning Models with their Applications to Control Problems. Joint Automatic Control Conference 1973 p 3-18.
- 100) Meyer P.A. (1966) Probability and Potentials. Ginn (Blaisdell).
- 101) Miller H.D. (1962b). Absorption Probabilities for Sums of Random Variables defined on a Finite Markov Chain. Proc Cam Phil Soc. 58 286-316.
- 102) Miller H.D. (1962a) A Matrix Factorization Problem in the Theory of Random Variables defined on a Finite Markov Chain. Proc Cam Phil Soc 58 268-285.
- 103) Milutin A.A. (1965) On the Automaton with Optimal Expedient Behaviour in a Random Environment. A.R.C. 26 No 1.
- 104) Monod J. (1971). Chance and Necessity. A.A. Knopf Inc.
- 105) Narendra K.S. and Viswanathan (1972) A Two-Level System of Stochastic Automaton for Periodic Random Environments. S.M.C. 2.
- 106) Narendra K.S. and Thathachar M (1974) Learning Automata- A Survey. S.M.C. 4 No 4
- 107) Neimark E and Estes W.K. (1967) Stimulus Sampling Theory. Holden Day.
- 108) Neumann Von J (1948) The General and Logical Theory of Automata. Collected Works Vol 5 Pergamon Press 1963.
- 109) Nilsson N.J. (1965) Learning Machines. Mc Graw Hill.
- 110) Norman M.F. and Yellott J.I. (1966) Probability Matching. Psychometrika 31 43-60.

- 111) Norman M.F. (1968a) Some Convergence Theorems for Stochastic Learning Models with Distance Diminishing Operators.
J. of Math Psych. 5 61-101.
- 112) Norman M.F. (1968b) On the Linear Model with Two Absorbing Barriers.
J. of Math Psych. 5 224-241.
- 113) " " (1970) Limit Theorems for Additive Learning Models.
J. of Math Psych 7 1-11.
- 114) " " (1971) Slow Learning with Small Drift in Two Absorbing Barrier Model. J. of Math Psych. 8 1-21.
- 115) Norman M.F. (1972) Markov Processes and Learning Models.
Academic Press.
- 116) Norman M.F. (1974) Markovian Learning Processes. SIAM Review 16 No 2.
- 117) " " (1975) An Ergodic Theorem for Evolution in a Random Environment. J. Applied Probability. 12 661-672.
- 118) Owen G. (1968) Game Theory. Saunders Company.
- 119) Page C.U. (1965) Equivalences between Probabilistic and Deterministic Machines. Information and Control. 9 469-520.
- 120) Paz A. (1971) Introduction to Probabilistic Automata. Academic Press.
- 121) Pittel B.G. (1965) The Asymptotic Properties of One Form of Gur Game.
P.I.T. 1 No 3 76-89.
- 122) Ponomarev V.A. (1964) On a Finite Automaton Asymptotically Optimal in a Stationary Random Media. Biofizika 9 No 1 104-110.
- 123) Pyatetskii-Shapiro I.I. and Vasil'ev N.B. (1967). The Time for an Automaton to Adapt to the External Media. A.R.C. 28.
- 124) Pyatetskii-Shapiro I.I. (1970) Mathematical Problems associated with Morphogenesis. P.I.T. 4 94-107.
- 125) Rabin M.O. (1963) Probabilistic Automata. Information and Control. 6 No 3 230-245.
- 126) Richardson D (1976) Self-Replication by Template.
Math Biosciences. 28 1-24.
- 127) Riordon J.S. (1969) An Adaptive Automaton Controller for Discrete Time Markov Processes. IFAC Journal 5 No 6 721-730.
- 128) Sawaragi Y. and Baba N. (1973) A Consideration on the Learning Behaviour of Variable Structure Stochastic Automata. S.M.C. 3 644-647.

- 129) Sawaragi Y. and Baba N. (1974) Two ϵ -Optimal Non-Linear Reinforcement Schemes for Stochastic Automata. S.M.C. 4 126-131.
- 130) Scarfe H. (1974) The Computation of Economic Equilibria. Yale Press.
- 131) Schmukler Y. (1970) Unidimensional and Bidimensional Gur Games. A.R.C. 31 Pt 2 1634-1642.
- 132) Schrödinger (1945) What is Life? C.U.P.
- 133) Shapiro I.J. and Narendra K.S. (1969) Use of Stochastic Automata for Parameter Self-Optimisation with Multimodal Performance Criteria. S.S.C. 5 No 4.
- 134) Skilling (1975) The Complete Enumeration of Uniform Polyhedra. Phil Trans A 111-137.
- 135) Skinner B.F. (1938) The Behaviour of Organisms. Appleton. New York.
- 136) Stefanyuk V.L. (1963). Example of a Problem in Joint Behaviour of Two Automata. A.R.C. 24 No 6 716-719.
- 137) Stefanyuk V.L. (1971) Description of Games of Two ϵ -Optimal Automata. A.R.C. 32 Pt 1 No 4.
- 138) Stewart D.J. et Al (1969) Automaton Theory and Learning Systems. Academic Press.
- 139) Stratonovich R.L. and Dobrovidov A.V. (1964) On the Synthesis of Optimal Automata to Operate in Random Media. A.R.C. 25 1289-1296.
- 140) Suppes P. and Lamperti J. (1960) Some Asymptotic Properties of Luce's β -Learning Model. Psychometrika 25 233-241.
- 141) Suppes P. (1969) Stimulus - Response Theory of Finite Automata. J. of Math Psych. 6 327-355.
- 142) Sushkov B.G. (1973). Symmetric Games for Large Aggregates of Probabilistic Automata. P.I.T. 9 95-99.
- 143) Takeuchi et Al (1974) Two Automata Games. Information Sciences. 7 81-93.
- 144) Tenisberg Yu. D. (1969) Some Models of Collective Behaviour in Dynamic Processes of Market Price Formation. A.R.C. 30 2 1140-1149.
- 145) Thom R. (1975) Structural Stability and Morphogenesis. Benjamin Inc.
- 146) Trakhtenbrot B.A. and Barzdin Ya. (1973) The Behaviour and Synthesis of Finite Automata. North Holland Press.
- 147) Tsertsvadze G.N. (1968) On the Asymptotic Proportion of Purposive Automata in a Stationary Environment. A.R.C. 29 Pt 2 No 8.

- 148) Tsertsvadze G.N. (1970) Rate of Establishment of the Final Distribution in a Game of Many Identical Automata. A.R.C. 31 1 583-586.
- 149) Tsypkin Yu.Z. (1966) Adaptation, Training and Self-Organisation in Automatic Systems. A.R.C. 27 No 1 23-61.
- 150) Tou J.T. et Al (1968) Applied Automata Theory. Academic Press.
- 151) Tsetlin M.L. (1961) On the Behaviour of Finite Automata in a Random Medium. A.R.C. 22 No 10.
- 152) Tsetlin M.L. and Krylov V.Yu. (1963) On Games of Automata. A.R.C. 24 No 7.
- 153) Tsetlin M.L. (1963) Finite Automata and Simple Models of Behaviour. Russian Mathematical Surveys. 18 No 4 1-27.
- 154) Tsetlin M.L. et Al (1964a). One Example of a Game for Many Identical Automata. A.R.C. 25 No 5 608-611.
- 155) Tsetlin M.L. et Al (1964b). Homogenous Games for Automata and their Computer Simulation. A.R.C. 25 No 11.
- 156) Tsetlin M.L. and Ginzburg S.L. (1965) Some Examples of Simulation of the Collective Behaviour of Automata. P.I.T. 1 No 2 54-62.
- 157) Tsetlin M.L. et Al (1965). The Behaviour of Automata in Periodic Random Media and the Problem of Synchronisation in the Presence of Noise. P.I.T. 1 No 1 65-71.
- 158) Tsetlin M.L. and Stefanyuk V.L. (1967) Power Regulation in a Group of Radio Stations. P.I.T. 3 No 4 49-57.
- 159) Tsetlin M.L. (1974). Automaton Theory and Modelling of Biological Systems. Academic Press.
- 160) Titterington (1971). Adaptive Optimisation of Yield. Phd Thesis Cambridge.
- 161) Tsuji H. et Al. (1973) An Automaton in the Non-Stationary Random Environment. Information Science. 6 123-142.
- 163) Vaisbord E.M. (1968a). Game of Two Automata with Differing Memory Depths. A.R.C. 29 Pt 1 440-451.
- 164) " " (1968b). Game of Many Automata with Differing Memory Depths. A.R.C. 29 Pt 2 1938-1943.
- 165) Varshavskii V.I. and Vorontsova I.P. (1963) On the Behaviour of Stochastic Automata with Variable Structure. A.R.C. 24 353-360.

- 166) Varshavskii V.I. and Gersht A.M. (1966) Behaviour of Continuous Automata in Stochastic Media. P.I.T. 2 No 3 68-75.
- 167) Varshavskii V.I. (1968a) Collective Behaviour and Control Problems. Machine Intelligence 3 No 14.
- 168) Varshavskii V.I. et Al (1968) Priority Organisation in Queuing Systems using a Model of the Collective Behaviour of Automata. P.I.T. 4 No 1 73-76.
- 169) Varshavskii V.I. et Al (1968b). Some Variants on the Problem of Synchronisation of Automata. P.I.T. 4 No 3 73-83.
- 170) Varshavskii V.I. (1969a) Synchronisation of a Collection of Automata with Random Pairwise Interaction. A.R.C. 30 Pt 1 224-228.
- 171) Varshavskii V.I. (1969b) The Organisation of Interaction in Collectives of Automata. Machine Intelligence. 4 No 16 285-311
- 172) Varshavskii V.I. et Al (1969c) Use of a Model of Collective Behaviour in the Problem of Resource Allocation. A.R.C. 30 Pt 1 1107-1114.
- 173) Varshavskii V.I. (1972a) Some Effects in the Collective Behaviour of Automata. Machine Intelligence. 5 No 22.
- 174) Varshavskii V.I. (1972b) Automata Games and Control Problems. Proceedings of the 5th World IFAC Congress Pt 3 37.4.
- 175) Vasershtein (1969) Markov Processes over Denumerable Products of Spaces, Describing Large Systems of Automata. P.I.T. 5 No 3 64-72.
- 176) Vasil'ev N.B. (1969). The Limiting Behaviour of a Random Media. P.I.T. 5 No 4 68-74.
- 177) Vasil'ev N.B. and Koganov A.V. (1973). A Model of Optimal Behaviour in an Unknown Media. P.I.T. 9 No 4 58-65.
- 178) Viswanathan R. and Narendra K.S. (1972) A Note on the Linear Reinforcement Scheme for Variable Structure Automata. S.M.C. 2 292-294.
- 179) Viswanathan R. and Narendra K.S. (1973). Stochastic Automata Models with Applications to Learning Systems. S.M.C. 3
- 180) Viswanathan R. and Narendra K.S. (1974). Games of Stochastic Automata. S.M.C. 4 131-135.
- 181) Volkonskii V.A. (1965). Asymptotic Properties of the Behaviour of Elementary Automata in a Game. P.I.T. 1 No 2 36-53.

- 182) Vorontsova I.P. (1965). Algorithms for Changing Stochastic Automata Transition Probabilities. P.I.T. 1 No 3 122-126.
- 183) Waddington C.H. (1966). Principles of Development and Differentiation. Macmillan.
- 184) Watanabe S. (1975). Creative Learning and Propensity Automaton. S.M.C. 5 No 6.
- 185) Whittle P. (1971). Stochastic Automata and Co-operative Effects. Technical Report No 69. Stanford University.
- 186) Wiener N. (1948). Cybernetics. Wiley. New York.
- 187) Witten I.H. (1973). Finite-Time Performance of Some Two-Armed Bandit Controllers. S.M.C. 3 194-197.
- 188) Witten I.H. (1974). On the Asymptotic Performances of Finite-State Two-Armed Bandit Controllers. S.M.C. 4.
- 189) Yasii T. and Yajima S. (1970). Two-State, Two-Symbol Probabilistic Automata. Information and Control 16 203-224.

List of Frequently Occurring Symbols. (with brief statement of meaning)

θ_k	k^{th} \bar{n} -cell, 11	$V^k(\bar{n})$	average reward received by θ_k , 76
$\bar{\pi}_i$	action distribution, 11	λ_j	Von Neumann saddle strategy, 76
u_i	action i, 11	Φ^{const}	approximate trajectories for mixed strategy saddles, 77
E_i	environment i, 11	$g_{\bar{n}}^k$	game matrix for θ_k in n-automata game, 85
$s=0$	penalty stimulus, 11	$\underline{s}(t)$	vector stimulus, 89
$s=1$	reward stimulus, 11	$\sigma_{ij}^s(t)$	structure transition matrix, 90
$q_i^{d,s}$	stimulus probability, 11	Γ_k	set of states which use θ_k , 90
$\Delta_{\alpha\beta}$	environment transition matrix, 11	x_i	state i, 90
$\mathcal{M}(A_{\alpha\beta}, q_i^{d,s})$	environment, 11	$g^1(\bar{n})$	network of \bar{n} -cells, 91
$R(\bar{n})$	average reward, 12	\bar{O}	\bar{n}_1 -cell or pure action u, 92
$\Delta X(t)$	expected increment in $X(t)$, 12	E_i	environmental equilibrium probability when in x_i , 95
U.L.	uniform learning, 12	γ_i^x	joint probability of E_{α} and x_i , 95
$\theta_{ij}(\bar{n})$	n-action learning function, 13	R_{jk}^{α}	reward comparison matrix for E_{α} , 98
s/mg	semi-martingale, 15	$\Delta\sigma_{ij}^{\alpha}$	equilibrium increment, 99
$\delta(\bar{n})$	absorption function, 19	$\omega_{\alpha}(t)$	probability in E_{α} at time t, 98
$U_{\theta_{ij}(\bar{n})}$	learning operator for $\theta_{ij}(\bar{n})$ rules, 18	SOSA	self-one-step-ahead, 101
R_0	optimal U.L. family, 23	\mathcal{M}_n	symmetric n-medium, 102
$i_1(\pi_i)$	probability take same action i for <u>all</u> time, 28	L_{R-P}	linear reward-penalty rule, 104
#alts	number of response alternations, 27	\mathcal{A}_0	limiting family of $g'(\bar{n})$ structures under R_0 , 109
$X(\bar{n})$	expected #alts, 30	$n \frac{L_{ij}^{(k_{ij})}}{m}$	linear family of automata with kernel (k_{ij}) , memory m_i and n-actions, 111
R_{ϵ}	ϵ -optimal U.L. family, 32	\bar{Z}_n	operating zone for n-medium, 115
$\text{Pr}(aBb c)$	skeleton absorption function, 56	\mathcal{B}_0	limiting family of $g'(\bar{n})$ structures under R_0 , 120
G_{ij}	stochastic game matrix, 72	$E_{\theta_k}^{\alpha}$	environmental equilibrium probability if in θ_k , 120
$\bar{\pi}_i^k$	action distribution for θ_k , 72	$g^n(\bar{n})$	n-level hierarchy of \bar{n} -cells, 131
e_i	equilibrium vector for $\Delta_{\alpha\beta}$, 62		

The fundamental idea underlying this research is that an initially randomly structured stochastic automaton placed in an environment can change its structure to increase its adaptation. The automaton gains information from the environment by executing actions and consequently receives reward or penalty stimuli with some probabilistic distribution.

The first part of the thesis concerns optimal reinforcement rules for markovian learning and a comprehensive theory is developed which also gives new insight into the operation of the many existing non-optimal rules. The behaviour of optimal rules in both static and dynamic environments is considered. An automaton evolving through "uniformly learning" rules which embrace both the new non-linear optimal rules and existing linear rules, is defined and called a \bar{U} -cell.

Next, a model for games between \bar{U} -cells is formulated and it is proved that they converge to pure saddle strategies when they exist. Deterministic approximations are found in the case of games with equilibrium mixed strategies and the trajectories are shown to be closely related to those of the Volterra model for predator-prey behaviour.

The thesis then develops a theory for the operation of networks of \bar{U} -cells. It is proved that starting from a random structure, the automaton evolves to an expedient structure which represents a discrete approximation to bayesian updating. The structure characteristically consists of a central action switch, surrounded by "arms" of memory states which reach out to the vertices of the likelihood simplex in which it is embedded. It is shown how such evolving automata can serve as a simple model of cellular differentiation. Finally it is shown how the theorems have natural extensions to \bar{U} -cells arranged hierarchically, which give greater adaptation in hierarchical environments.

