

# MPRA

Munich Personal RePEc Archive

## **A dynamic model of interactions between conscious and unconscious**

Lotz, Aileen and Gosselin, Pierre

Cerca Trova, BP 114, 38001 Grenoble Cedex 1, France,

Institut Fourier, UMR 5582 CNRS-UJF, Université Grenoble

I, BP 74, 38402 St Martin d'Hères, France

15. February 2012

Online at <http://mpra.ub.uni-muenchen.de/36697/>

MPRA Paper No. 36697, posted 16. February 2012 / 05:58

# A dynamic model of interactions between conscious and unconscious

Aileen Lotz\*      Pierre Gosselin†

February 2012

## Abstract

This paper advocates that some limits of the rational agent hypothesis result from the improper assumption that one individual should be modeled as a single rational agent. We model an individual composed of two autonomous and interacting structures, conscious and unconscious. Each agent utility form depends both on external signals and other structures' actions. The perception of the signal depends on its recipient and its grid of interpretation. We study both the static and dynamic version of this interaction mechanism. We show that the dynamics may display instability, depending on the structures interactions' strength. However, if unconscious has a strategic advantage, greater stability is reached. By manipulating other structures' goals, the strategic agent can lead the whole system to an equilibrium closer to its own optimum. This result shows that some switch in the conscious' objective can appear. Behaviors that can't be explained with a single utility can thus be rational if we add a rational unconscious agent. Our results justify our hypothesis of a rational interacting unconscious. It supports the widening of the notion of rationality to multi-rationality in interaction.

Key words: dual agent; conscious and unconscious, rationality; multi-rationality; emotions; choices and preferences; multi-agent model; consistency;

JEL Classification: B41,D01, D81, D82.

---

\*Aileen Lotz: Cerca Trova, BP 114, 38001 Grenoble Cedex 1, France. E-mail: a.lotz@ercercatrova.eu

†Pierre Gosselin : Institut Fourier, UMR 5582 CNRS-UJF, Université Grenoble I, BP 74, 38402 St Martin d'Hères, France. E-Mail: gosselin@ujf-grenoble.fr

*"[...] in this you should let yourself be guided not by any fixed purpose but mainly by intellectual curiosity and a spirit of exploration."*  
Hayek, Friedrich A. 1944. 'On Being an Economist'

## 1 Introduction

For about a hundred years, Slutsky's results<sup>1</sup> have been consistently rejected. Only recently did Browning and Chiappori (1998) show that individual consumers do indeed solve Slutsky's equations, even if this optimization is not done consciously. This result has generally been taken as a confirmation of the rational agent hypothesis : individuals are infinitely rational, and their decisions are those of utility-maximizing agents.

To be fair, what Browning and Chiappori results mostly show is the existence of an unconscious rationality. The shift is however easy to understand : conscious and unconscious being part of the same individual and receiving the same signals, it can be assumed that they are endowed with the same information and utility, i.e. preferences, goals, actions, etc... Nothing distinguishing them, it is logical to combine them in one rational action, i.e. the agent's action. In this respect, to demonstrate the unconscious rationality is equivalent to proving the whole agent's rationality.

Yet countless arguments and experimental facts seem to demonstrate the agent's irrationality. Allais<sup>2</sup> and Kahneman-Tverski<sup>3</sup>, for example, have shown that under simple situations, agents display systematic psychological biases. And indeed, agents' irrationality is usually attributed to psychological factors. It is assumed that individuals are, most of the time, driven by their emotions, and, ultimately, their unconscious. The action being irrational, so is the unconscious, and so is the agent.

If the unconscious is indeed irrational, nothing can be said about it, and we must reduce ourselves to sum the list of its behaviors. If we are to understand anything about the unconscious, we have to suppose him to be, at least partly, rational : we must endow him with all the attributes of the rational agent, in the economic sense of the term. Besides, this is what Browning and Chiappori results tend to prove. The question therefore is not to know whether the unconscious is rational or not, but rather in what respect his rationality differs from the conscious' rationality.

Note that utility optimization, possibly including some forecatings, has already successfully modeled (unconscious) automatic behavior, for example in the motor or visual system<sup>4</sup>. These results confirm that economics is relevant to explain unconscious neural processes. Yet similar models for conscious decision making, and, more generally, non automatic processes, have been highly criticized.

---

<sup>1</sup>See [38].

<sup>2</sup>See [2].

<sup>3</sup>See [26].

<sup>4</sup>See [35] for an overview.

The main criticism is that they should include some unconscious phenomena. It advocates the use of partial rationality to describe seemingly unrational, or sudden switches in, choices.

For example, a well know anomaly is described by the following situation : In a restaurant, a consumer can choose chicken or beef. He orders chicken. But when the waiter suggests a third dish, fish, the consumer orders beef. This exemplify the independance of irrelevant alternatives. Some choices, even when they may be rational, lack intelligibility.

In this paper we confront this problem by extending the range of the rational explanation. We advocate that some of the limits of the rational agent hypothesis result from an improper assumption: the fact that one individual should be modeled as a single rational agent. What can suffice to explain automatic processes should be extended for more complex tasks. It seems natural to postulate decision making not only involves the conscious individual but also some rational unconscious processes.

This question has already been adressed in [31]. The author considers the individual as a "dual agent", i.e. the result of two distinct agents, conscious and unconscious, endowed with their own utility, goals and preferences. Importantly, although they receive the same signals from the outer world, both the conscious and unconscious interpret these signals in different way. Their actions are rational, but the differential of information between them creates a bias in the action of the dual agent.

However, this simple model does not explain the interactions between conscious and unconscious. Besides, it lacks of dynamics. One of our goals is to remedy these weaknesses.

In order to implement this program, we will show that a cognitive unconscious can rationally modify conscious actions and representations, by inducing changes in choices, goals, as well as systematic biases in actions. This does not solve the above anomaly, but gives some intelligibility to seemingly irrational an fluctuating preferences and choices.

Our results justify our hypothesis of a rational interacting unconscious, and supports the widening of the notion of rationality to multi rationality in interactions.

More precisely, we model an individual, the dual agent, as composed of several autonomous and interacting rational agents, or structures. These structures differ in their goals, information, and action. For the sake of simplicity, we will first assume two structures, the "conscious" and the "unconscious" The first structure schematically describes the agregate of the whole set of conscious processes. For the ease of exposition, the second structure merely gathers the remaining processes.

Fairly enough, this assumption, although practical in our demonstration, could be taxed as dubious. The unconscious could be more acuretely modeled by a set of multiple autonomous processes. Indeed, our simplifying assumption will later be relaxed to include an arbitrary number of structures.

Besides, one could argue that the distinction between conscious and unconscious is not clear cut. Unconscious contents can become conscious and vice-

versa in a continuous flow. In this circumstances, disentangling the two could seem artificial. However, we argue in this paper that this complexity is the mere consequence of two interacting agents, exchanging information, and impacting one another. The dual agent is precisely the result of this interaction and entanglement.

The key point of our model is that each agent utility varies in time. Actually, its form depends both on external signals and on internal ones, i.e. the other structure actions. Moreover, each structure has its own grid of interpretation. Therefore, each structure will interpret the same signal in its own, specific, way. Some signals may be considered as irrelevant, and as such be disregarded by a structure<sup>5</sup>. Typically, we can assume that the conscious structure interprets external signals accurately, whereas the unconscious interpretation may refer to past or even symbolic events. The archetypal example is the psychoanalysis view of the unconscious that may, more often than not, use condensation and displacement.

These assumptions allow us to model reactive structures : the scene they build is adapted to the outer world, yet is also influenced by the other structure reactions. In other words, each structure tailors the system of representation of the other. It is this mutual influence that creates the entanglement. It can also explain the variation in goals without referring to external modifications.

We study both the static and dynamic version of this entanglement mechanism. Considering in parallel the case of, first, two non strategic structures, and second, one structure having strategic advantage, we show that the dynamics may display instability, depending on the structures interactions' strength. However, under strategic advantage, greater stability is reached. By manipulating the other structure's goals, the strategic agent can lead the whole system to an equilibrium closer to its own optimum. He can do so without inducing an over reaction from the non strategic agent.

This result shows that some switch in the conscious' objective can appear, even in a constant environment, through this interaction mechanism. Behaviors that can't be explained with a single utility can thus be rational if we add a rational unconscious agent.

The paper is organized as follows : in section 2, we describe a general pattern of interaction between various autonomous structures. Those structures interact through their exchange of information. Both agents receive outer signals, as well as inner signals that are the other structures' actions. They process these signals to produce some information. The signals may, when relevant, activate the structure. Once activated, the structure builds a utility function from its information, and optimizes it through its action. In section 3 we present a static model of interaction between two structures, "conscious" and "unconscious". Its simple pattern is general enough to convey the main points of section 2. The model's equilibrium is studied and discussed for two cases : two non strategic agents, and one strategic agent (unconscious) facing a non strategic one (conscious).Section 4 generalizes the pattern of section 3 and develops a

---

<sup>5</sup>In that respect, we follow the view presented in [33].

general dynamic model of interaction between  $n$  structures. At each moment  $t$ , the structure's utility is shaped by some parameters. These parameters depend on the other structure signals, i.e. actions, as well as on the external signals<sup>6</sup>. Section 5 applies the section 4 results to a dynamical version of section 3 model, and draws the implications of its dynamics. It shows, in particular, how the strategic agent can shift the goal of the non strategic agent to reach his own objective more easily. The last section concludes and present some perspectives for further research.

## 2 Setup

We will first precise our assumptions about the conscious and unconscious, as well as their interactions as rational agents.

### 2.1 Hypotheses

#### 2.1.1 Conscious and unconscious

We suppose that all human activity originates in the unconscious, in the neuroscientific acceptance of the term. More specifically, we suppose the existence of autonomous and interacting unconscious structures, endowed with a certain degree of planification and reflexion. It is the notion of "cognitive unconscious", introduced by Kihlstrom (1987)<sup>7</sup>. According to this notion, some unconscious processes can show, at least to some degree, awareness. Besides, the survey [32] for example, describes how these autonomous units are capable of coordination, and present a pattern of metastability.

The notion of unconscious studied here is inspired by the combination of these two approaches : the unconscious is a set of structures more or less autonomous and linked one to the other. Under some conditions, part of these structures become conscious. We call "unconscious" the structures non emerging to the conscious at a given point in time. By opposition, we call "conscious" the set of unconscious structures having emerged to the conscious. The conscious can therefore be seen as a workspace directly fed by the unconscious.

As described, conscious and unconscious are, technically, merely aggregates of sub-structures loosely linked to one another. In the following, we will set aside this fact, and will indifferently qualify them as "structures". This is a useful approximation that will not modify our results.

These structures are rational agents receiving and sending signals.

---

<sup>6</sup>This fit at least partly with Edelman's presentation of consciousness in [10]: each structure builds a scene at a certain time through the entrance and reentrance of external as well as internal correlated signals.

<sup>7</sup>An account of this theory can be found in [23].

### 2.1.2 Signal and structures

Both conscious and unconscious receive a continuous flow of signals originating from the outside environment, or equivalently, from the other structure. Each structure receives all the signals fully, and screens them according to its own grid.

Each structure operates according to its own pattern. This pattern is a grid of a priori characteristics, according to which the structure analyzes a signal. Exactly in the same way as a photoelectric cell fail to react to a sound wave, and a phone ignores a ray of light, each structure reacts to the signals it can screen through its own pattern, and extract from these signals the inputs it can perceive. We will call "informations" these inputs. A structure that has managed to extract an information from a signal is said to be activated.

The information is a signal to which the structure has given a special meaning, corresponding to its own operating system. The signals that will not have been successfully processed by the structure will be fully disregarded.

Since the information is the interpretation of a signal, and not the signal itself, the information extracted by the structure can be radically different from the original signal.

Note that our approach is similar to the one developped in [33] where a structure has its own grid of lecture, through which it attempts to construct an object based on the signals it receives. Through a feed back process, it may qualify this signal as a noise, and not activate itself.

### 2.1.3 Structure and utility

Conscious and unconscious are modeled as rational economic agents. This is not unusual in neurosciences, where, for instance, the motor system of the body is modeled as an agent minimizing a loss function according to an incoming information<sup>8</sup>.

Each structure has its own utility, which is a function of parameters and a set of possible actions. Recall that the structure's environment (i.e. the outer world and other structure) send some signals that describe a situation. The structure attempts to extract information from these signals. If it succeeds, it will then be activated, meaning that it will adapt some parameters of its utility according to this new information. Alternatively, if the structure does not succeed in extracting information from the signals, it will not be activated, and its parameters will remain set to zero.

For a structure to be activated, only few parameters are sufficient, provided that they are significant enough. The parameters of an unactivated structure are set to 0. When a structure is activated, its relevant parameters take a specific value that quantifies the information extracted from the signal. The utility thus build will induce an action directly linked to the incoming stimulus.

---

<sup>8</sup>See [40] and [35] for example.

## 2.2 Actions

With its significant parameters now set to a certain value, the structure can react to the context. The action it will take can be manifold : physical, physiological or neuronal. It will optimize the structure utility, taking as given the value of the parameters, as would any rational agent do.

### 2.2.1 The agents grid of lectures

Let us come back to the main point of this section. The information is merely the structure's interpretation of a signal. A single signal can therefore produce different informations, depending on its recipient. This information is quantified into parameters, the utilities true inputs. Each structure therefore perceives, through one single signal, two different situations.

A minima, the unconscious is that part of the individual that is not conscious, and continuously manages one's vital functions. Breathing, hunger, thirst, for instance, .are all vital functions and that are not performed consciously by the individual and denote the action of an unconscious agent. Although minimalist, this conception highlights the fact that, within each and every one of us, a hidden agent operates according to its own objectives.

That the signals of the environment should have the same meaning for him than for the conscious may seem dubious. The unconscious defined here is characterized for a large part by its intemporality. Each second, it must act and react to maintain an equilibrium aquired in the past, and still at stake today. On the contrary the conscious is marked by temporality. By its intentionality, the conscious can handle multiple dynamic and variable representations, to play with memories, to project itself. We can therefore suppose that its grid of lecture is much more flexible and mobile than the one of the unconscious.

Consequently, the crucial distinction between conscious and unconscious lies in the interpretation of the signals, and in the construction of the information. Whereas the information of the conscious is marked by its temporality, the information of the unconscious will be characterized by its intemporality. Given a determined set of external signals, the unconscious has a static lecture of the events: it continuously reacts to signals from which it systematically extracts the same type of information. The conscious, on the contrary, extracts a wider range of information from the signals. This gives him a more dynamic vision of the situation, but also a more short-sighted one. The unconscious will interpret the signals under a more permanent light.

### 2.2.2 Notations

To make things easier,  $S_j$  will be any structure, aggregated or not, conscious or unconscious. Only two (agregated) structures will be considered in what follows, the conscious and the unconscious. By convention  $S_1$  will be the conscious, and  $S_2$  the unconscious.

Each structure, once activated, produces its own action,  $a_1$  and  $a_2$  respectively. Unconscious actions emerge into consciousness: a reflex, for instance,



stemming directly from the unconscious, is perceived by the conscious, and used as information and parameter in its utility.

These two actions are not independent :  $a_2$  is taken before  $a_1$  and, being itself a signal, will impact  $S_1$ 's information.  $a_2$  will therefore be both an action and a signal to  $S_1$ .

### 2.2.3 Interaction between structures

Much of the structures' interactions dynamics rests on the mutual perception each structure has of the other. Whereas we may well, as individuals, be fully aware that unconscious processes exist, this does not necessarily mean that our conscious can actively use this information to his own advantage. To do so, he should, first recognize the impact of the unconscious on its behavior, and then be able to manipulate it.

Therefore, the analysis must take into account the number of strategic agents it involves We can distinguish three situations :

1. No strategic agent : Conscious and unconscious are unaware of each other, and consider the information they get as an outside data.
2. Two strategic agents.: Conscious and unconscious are aware of the other, but, having no control over it, treat its signals as any other external signal. This will result in a Nash equilibrium, with an outcome similar to the above case.
3. One strategic agent: One of the agents is aware of the other and can manipulate it through its signals.

## 3 Static model

### 3.1 Setup

#### 3.1.1 Preliminary remarks

To illustrate our point, we will apply our setup to a specific and very basic case. In what follows, we will suppose an agent  $\lambda$  whose utility depends positively on meeting people. We will measure this outcome as a variable depending on the amount of social life  $\lambda$  can get. By assumption, each time  $\lambda$  can go out (signal), he feels a sudden strain that effectively limits his ability or willingness to do so.

The economic theory does not need the unconscious to explain this outcome. It would model this sequence of facts by setting a cost to  $\lambda$ 's outing. And indeed, that's what it boils down to, eventually. Yet one may wonder how such a cost could arise, and be rationalized. Should it be fixed, or should it be reset at each new situation in an ad hoc manner?

We will depart from this solution and assign to the unconscious agent a utility distinct from the conscious. In  $\lambda$ 's case, the unconscious has indeed perceived the possibility to go out. However, unlike the conscious, going out

will diminish its utility. When  $\lambda$  goes out, an additional adverse signal is sent to the unconscious. He will need to counterbalance it with the means at its disposal.  $\lambda$ 's strain can now be read as the optimal reaction of the unconscious to optimize his utility, given the (bold) action of the conscious agent.

These hypotheses model how and why the cost of an action can arise. They further give a rational explanation to a fact that would otherwise seem irrational. More generally, they show how a utility can arise to the conscious. Under the impulse of external signals, the structures exchange signals, that in turn model each other's utility. Here, the conscious utility has been modified by the cost arising from the unconscious action.

This approach is in line with the neurosciences and their study of the formation and variations of conscious contents. We are all aware of what is conscious, we are aware of our goals, yet we seldom know the origin of these goals. Our model rationalise these facts by supposing an unconscious agent continually shaping both our utility and reality.

### 3.1.2 Description of the model

Having underlined the dynamic nature of the structures interaction, let us now detail their modelisation.

The two structures are activated by one (common) signal. Each one reacts by sending signals perceived by the other structure. This could go on indefinitely, but to simplify the matter we will first consider a basic, one period interaction. With a one period horizon and no temporal dependance, solving the model conveniently boils down to finding its long-term equilibrium.

Over this period,  $S_1$ , the conscious, extracts and plans first. Only then does  $S_2$ , the unconscious structure, react to  $S_1$ 's planned action. This setup, open to criticism, allows two concomitant actions, which satisfies our one period hypothesis. We will later relax this assumption by introducing a delayed effect in the dynamic model.

**The conscious grid of lecture and action** We mentioned in section 2.1.2. that  $S_1$ , the conscious, has a "temporal" grid of lecture. It is adaptable, and relatively flexible to the present context.  $S_1$  is therefore able to read, interpret and extract information from the signals in a relatively efficient way. So that, were it not for  $S_2$ ,  $a_1$  would be set to  $a_0$ , up to a random error term, as is usual for the rational agent.

$S_1$ , receives a set of signals  $\sigma$ , from which it extracts information and deduce an action,  $a_1$ , associated to a specific reward. This reward is a utility  $U_1$  that is a function of  $a_1$  and  $a_0$ , the parameter describing the optimal action in the present context.

In quantifying the information through its parameters, the structure has inferred an external context. It has associated to it a utility. The shape of this utility depends therefore on the information  $a_0$  extracted from the signals  $\sigma$ .

For our agent  $\lambda$ , the set of signals  $\sigma$  is an invitation to go out.  $a_0$  would be  $\lambda$ 's optimal action, i.e. a time out, given his constraint.  $\lambda$  knows that  $a_0$  will

maximize his utility, and if he is rational, he will choose  $a_1$  such that  $a_1=a_0$ .

**The unconscious grid of lecture and action** Having an intemporal grid of lecture,  $S_2$  interpretation of the signals can very well be poorly related to the present context. Its reactions, however optimal in view of its own grid of lecture, will, more often than not be inadapted to the situation.

Whereas  $S_1$  has received a signal  $\sigma$  and considers  $a_0$ ,  $S_2$  receives both the signal  $\sigma$  and the planned optimal action  $a_0$ . It screens them through its own grid of lecture, extracts from them an information, and parting, a utility radically different from  $S_1$ . We will call  $a'_0$   $S_2$ 's interpretation of  $a_0$ .  $S_2$  sees  $a'_0$  as inducing a loss in its utility. In the same vein, it interprets  $a_1$  as a sub-optimal action  $a'_1$ , and will, as such, try to counter it.

We have defined the unconscious as chiefly ensuring the agent's vital functions. The planned outing  $a_0$  could, in our example, be seen as detrimental for the metabolism. Other, deeper reasons could come into play. A past, long forgotten experience could trigger a general fear of outings, seen as sources of danger. Eventhough the information extracted by  $S_2$  has nothing to do with  $a_0$ ,  $a'_0$  will then activate the unconscious utility. This will in turn induce a reaction. We will call  $a_2$  the unconscious reaction to the (combined) conscious action and external signals. This action  $a_2$  will be meant to set a cost to limit the conscious action, prevent him to perform  $a_0$ , and therefore limit its action.

To go back to our example,  $a_2$  could be a sudden strain or anxiety due to the phobia, rendering the action of  $\lambda$  difficult or impossible. Each structure's action being specified, we can now specify their respective utilities, as specified by their interactions.

### The agents utilities

**The conscious utility**  $S_1$  is described by the following quadratic utility :

$$U_1 = -\frac{1}{2}(a_1 - a_0)^2 - a_2 a_1 \quad (1)$$

where the first term describes  $S_1$ 's "own" utility, and the second is the cost imposed by  $S_2$  to counter  $S_1$  action. Were there only  $S_1$ , this second term would disappear.

$a_2$  is  $S_2$ 's reaction to  $S_1$ 's action. Insofar as  $S_1$  ignores the existence of  $S_2$ ,  $S_1$  faces  $a_2$  without being able to precisely determine its origin. Consequently,  $S_1$  will optimize its time out taking into account both its will to go out, and the strain induced by actually going out.

**The unconscious utility**  $S_2$  remains unconscious, i.e. it does not appear in the conscious framework. But being activated by the signal, it will send to the conscious its own action,  $a_2$ . Its utility will also be described as quadratic, and will be written :

$$\begin{aligned}
U_2 &= E_{s_2} \left[ -\gamma \frac{a_1'^2}{2} + \alpha a_2 a_1' - \frac{a_2^2}{2} \right] \\
&= -\gamma \frac{a_1^2}{2} + \alpha a_2 a_1 - \frac{a_2^2}{2}
\end{aligned} \tag{2}$$

$E_{s_2}$  is the expectation that characterize  $S_2$ 's interpretation of the signals.  $S_2$  intends to react to  $a_1'$ , and will use  $a_1$  as a proxy for  $a_1'$ .

$S_2$  suffers  $a_1'$ , his interpretation of  $a_1$ , and his own action to counter  $a_1'$ ,  $a_2$  : every action bears its own cost. Here,  $-\gamma \frac{a_1'^2}{2} - \frac{a_2^2}{2}$  is  $S_2$ 's cost, both to suffer  $a_1$  and to react to it, i.e. to produce a strain.

$S_2$  will nonetheless gain something in the process :  $\alpha a_2 a_1'$ . This gain in  $S_2$ 's utility is directly proportional to the cost it imposes on  $S_1$ ,  $a_2 a_1$ . One structure's cost is the gain in utility of the other.  $S_2$  somehow discovered that the strain diminishes the planned action. Setting  $\gamma - \alpha^2 > 0$  ensures a stable maximum for the utility in the two variables  $a_2$  and  $a_1$ .

Given this setup, and depending on the nature of the agents, two cases can arise : either both agents are not strategic, or one of them is strategic and manipulates the other. This would imply that one agent ignores the existence of the other, or at least cannot counter it effectively.

Both agent could be defined as strategic. However dreams, somatic disorders, phobias, amply demonstrate that we are more affected by our unconscious than we would wish to be, and are condemned to take his actions as given constraints. We therefore suppose the unconscious to be strategic, and the conscious to be non strategic. We will successively study the case for a non strategic unconscious, and for a strategic one.

## 3.2 Interactions between agents

### 3.2.1 The non strategic unconscious agent

Here, structures can be seen as being non strategic, or both strategic and neutralizing each other. Alternatively, they could even be unaware of each other's existence. Conscious and unconscious playing simultaneously, the situation will result in a Nash equilibrium.

Each agent will choose its action according to its grid of lecture and given the other action : whereas the conscious read  $a_0$  and infers  $a_1$  for what they truly are, the unconscious, will read  $a_0$  and  $a_1$  as  $a_0'$  and  $a_1'$ .

The equilibrium of the system is found by considering simultaneously the two agents. Optimizing  $U_1$  yields

$$a_1 = a_0 - a_2$$

and the equilibrium is found by replacing  $S_2$  optimal response :

$$a_2 = \alpha E_{s_2} a_1' = \alpha a_1$$

This action could as well be the answer of an automate mechanically reacting to  $a'_1$ . The rational agent hypothesis does not necessarily refers to an individual, but rather to a structure optimizing a utility. The motor system, for instance, is well modeled by a learning rational agent<sup>9</sup>.

Resolving the system yields the following actions:

$$\begin{aligned} a_1 &= a_0 - \alpha a_1 \\ a_1 &= \frac{a_0}{1 + \alpha} \\ a_2 &= \frac{\alpha a_0}{1 + \alpha} \end{aligned}$$

where  $a_1^{(NS)} = \frac{a_0}{1+\alpha}$  et  $a_2^{(NS)} = \frac{\alpha a_0}{1+\alpha}$  are the actions of the non strategic agents.

Without the unconscious, the conscious optimal action would have been  $a_0$ . As soon as  $a_0$  is systematically missed, we can infer from this bias the presence from the unconscious. This bias results from the differential of information between conscious and unconscious. This phenomenon has been explained in [31]. Actually  $a_2$  results from a misreading of  $a_1$  by  $S_2$ . If the unconscious had not set  $E_{s_2} a'_1 = a_1$ , it could have given to  $a'_1$  its true value, 0, and would have chosen  $a_2 = 0$ . As a consequence, the conscious could have set  $a_1 = a_0$ . With this in mind, we rewrite :

$$\begin{aligned} a_1 &= a_0 - a_2 \\ &= a_0 - \alpha E_{s_2} a'_1 \\ &= a_0 - \alpha (E_{s_2} a'_1 - E_{s_1} a'_1) \end{aligned}$$

The bias between  $a_1$  and its optimum  $a_0$  is  $-\alpha (E_{s_2} a'_1 - E_{s_1} a'_1)$ , and is the exact expression of a difference of information between conscious and unconscious.

[31] shows this action is in fact the optimal action of a single agent, the dual agent.

This dual agent is the individual whose utility explicitly encompasses both his conscious and unconscious utilities. As such his action is a combination of two actions, the conscious and unconscious ones respectively, thus including the previous bias in its action. Note that the dual agent, seen as a single individual, has the remarkable property to produce two different forecasts for one unique signal, revealing the combination of two interacting autonomous agents. This notion is equivalent to this paper approach.

How could the bias be qualified? It originates in a difference between perceptions, that counters or shift an action. It can be psychological or physical. Besides, it does not correspond to a seemingly rational reality, and can result in a strain or a well-being. We will sum up these elements by stating that this shift is an emotion.

---

<sup>9</sup>See [40].

This emotion is, by definition, the result of a differential of information between conscious and unconscious. As such, it could potentially be reduced to 0, provided two identical grids of lectures.

Here, this differential of information results in a loss in welfare : had the unconscious recognized that  $a_1 \neq E_{s_2} a'_1$ , he would not have acted this way, and the conscious would have reached  $a_0$ . The two agents' utilities would be equal to their reference value 0. Here, we rather have:

$$\begin{aligned} U_1^{NS} &= -\frac{1}{2} \frac{\alpha(\alpha+2)}{(\alpha+1)^2} a_0^2 < 0 \\ U_2^{NS} &= \frac{1}{2} a_0^2 \frac{\alpha^2 - \gamma}{(\alpha+1)^2} < 0 \end{aligned}$$

We will compare these values to the case of the strategic agent below.

### 3.2.2 The strategic unconscious agent

The general setup is unchanged :  $S_1$  and  $S_2$  act simultaneously.  $S_1$  still optimizes its utility taking as given  $S_2$ 's action, and sets its optimal action  $a_1 = a_0 - a_2$ . However,  $S_2$  now observes  $S_1$  and knows its optimal reaction  $a_1$ . It will take it into account and use it as a parameter to set its own action  $a_2$ .

We find  $a_2$  by replacing  $a_1$  as a function of  $a_2$  in  $S_2$ 's utility, and by optimizing over  $a_2$ , which gives :

$$\begin{aligned} U_2 &= E_{s_2} \left[ -\gamma \frac{a_1'^2}{2} + \alpha a_2 a_1' - \frac{a_2^2}{2} \right] \\ &= -\gamma \frac{a_1^2}{2} + \alpha a_2 a_1 - \frac{a_2^2}{2} \\ &= -\gamma \frac{(a_0 - a_2)^2}{2} + \alpha a_2 (a_0 - a_2) - \frac{a_2^2}{2} \end{aligned}$$

The optimum is  $a_2 = \frac{(\alpha+\gamma)}{2\alpha+\gamma+1} a_0$ . We then deduce that  $a_1$  is :

$$\begin{aligned} a_1 &= a_0 - a_2 \\ &= \frac{1+\alpha}{1+2\alpha+\gamma} a_0 \end{aligned}$$

We note  $a_1^{(S)} = \frac{1+\alpha}{1+2\alpha+\gamma} a_0$  et  $a_2^{(S)} = \frac{(\alpha+\gamma)}{2\alpha+\gamma+1} a_0$  the agents' actions under the strategic case.

We can check that  $a_1^{(S)} - a_1^{(NS)} = \frac{\alpha^2 - \gamma}{(\alpha+1)(2\alpha+\gamma+1)} a_0 < 0$ . This shows that the unconscious, by manipulating the conscious, has reached a higher equilibrium from his perspective : compared to the previous case, the conscious' action is reduced. Concretely, our agent  $\lambda$  will go out even less than before.

But since  $a_2^{(S)} - a_2^{(NS)} = \frac{(\alpha+\gamma)}{2\alpha+\gamma+1} a_0 - \frac{\alpha a_0}{1+\alpha} = a_0 \frac{\gamma - \alpha^2}{(\alpha+1)(2\alpha+\gamma+1)} > 0$ , this results from a stronger reaction of the unconscious than in the non strategic case.

### 3.2.3 Utility comparison

We can compute the difference of utilities between the strategic and non strategic case for each agent  $\therefore$

$$\begin{aligned} U_2^S - U_2^{NS} &= \frac{1}{2} \frac{(\alpha^2 - \gamma)^2 a_0^2}{(\alpha + 1)^2 (2\alpha + \gamma + 1)} > 0 \\ U_1^S - U_1^{NS} &= \frac{1}{2} \frac{(\alpha^2 - \gamma) (4\alpha + \gamma + \alpha^2 + 2) a_0^2}{(\alpha + 1)^2 (2\alpha + \gamma + 1)^2} < 0 \end{aligned}$$

There is a clear gain for  $S_2$  to be strategic, as well as a loss for  $S_1$  to be manipulated. Nevertheless, the unconscious gain in utility refers, by definition, to past situation and fully-reconstructed situation, furthermore based on biased signals. From the social point of view, and over the long run, the conscious loss in utility should result in a general loss for the dual agent.

## 4 General case and Dynamics

### 4.1 General setup

Let us now introduce a general dynamic model of interactions that will describe  $n$  structures (practically,  $n = 2$ ) sending arbitrary signals and taking any possible actions. It will encompass a dynamic version of the previous model.

Recall the general pattern of interaction that arised in the static example  $\therefore$

- A structure  $S_i$  is described by a vector of possible actions  $X_i(t)$  and a set of parameters  $P_i(t)$ . Both live in some, possibly different and infinite-dimensional spaces.

- $P_i(t)$  is the information  $S_i$  extracts from the signals and actions sent by other structures and the outer world.

- $S_i$ 's actions  $X_i(t)$  can act as signals for an other structure  $S_i$ , and in turn influence  $P_j(t)$ .

- $S_i$  utility directly depends on its actions, as well as on  $P_i(t)$ .  $S_i$  utility at time  $t$  is given by, at the quadratic approximation<sup>10</sup>:

$$V_i(t) = -\frac{1}{2} X_i^t(t) A_i X_i(t) + X_i^t(t) B_i P_i(t) - \frac{1}{2} P_i^t(t) C_i P_i(t) \quad (3)$$

At each point in time, the utility is fully dependent on both the outer world and the other structures through  $P_i(t)$ . Besides, the utility exhibits its specific pattern through the matricial coefficients  $A_i$ ,  $B_i$  and  $C_i$ .

The form of the utility  $V_i(t)$  deserves to be emphasized : it explains why, and how, the conscious utility evolves through time. The structure's rational choices can vary at each moment in time and evolve with its environment. But they also, and more importantly, depend on other structures' perception of the

<sup>10</sup>In the following,  $Y^t(t)$  denotes the transpose of any vector  $Y(t)$ .

environment. This can explain why some facts may appear important at time  $t$ , and no longer at time  $t + 1$ . This does not necessarily occur in a seemingly continuous manner, since the conscious is unaware of the unconscious, and is subject to it.

Ultimately,  $S_i$  optimizes at time  $t$  an intertemporal utility :

$$U_i(t) = \sum_{n \geq 0} \beta_i^n V_i(t + n)$$

where  $\beta_i$  is the rate of time preference and . ranges between 0 and 1. The higher the  $\beta_i$ , the less the future is discounted by the structure.

## 4.2 Time schedule and information process

The interaction between structures is dynamic : there is a delay between the moment an action is taken and the moment it is processed as an information. Each period  $t$  is thus subdivided in two steps: first, all structures process the information, then all react to this information at the same time. The delay between the information process and the action depends on each structure time scale, and can be relatively short. For the sake of simplicity, this time scale will be considered uniform across structures. The structures will extract the parameters from the signals, through the information process, as mentioned earlier.

We model our hypotheses on the parameters  $P_i(t)$  in the following way.

$$P_i(t) = \sum_{i \neq j} F^{i,j}(t) X_j(t-1) + F^{i,ext}(t) P_{ext}(t)$$

where  $P_{ext}(t)$  describe the external situation.

$F^{i,j}(t)$  and  $F^{i,ext}(t)$  are filter matrices through which  $S_i$  interprets signals of  $S_j$  and the external signals, respectively. It models a linear information process where information  $P_i(t)$  is reconstructed from the signals.

## 4.3 Non strategic agents

Let us now describe the dynamics of a system of  $n$  non strategic structures. Here again, as in the static case, structures do not deliberately influence one another . Both take their parameters as given, and choose their action regardless of its effect on other structures at time  $t + 1$ . Since there is no intertemporal constraint, each period is independent. Each structure thus optimizes  $V_i(t)$  one period at a time, that yields the optimal action for  $S_i$  :

$$X_i(t) = D_i P_i(t) = D_i \left( \sum_{i \neq j} F^{i,j}(t) X_j(t-1) + F^{i,ext}(t) P_{ext}(t) \right)$$

with  $D_i = A_i^{-1} B_i$ . This equation represents the dynamics for the action vector  $X_i(t)$  for  $S_i$ . It is not self-consistent, since it involves the other structures



dynamics. To solve the dynamics for each structure, we need to stack all the structures dynamical equations into one unique system. We introduce a vector  $(X(t)) = \begin{pmatrix} X_1(t) \\ \dots \\ X_n(t) \end{pmatrix}$  that encompasses all actions at time  $t$ . The previous dynamical equations can be rewritten as:

$$(X(t)) = \left(\hat{F}(t)\right) (X(t-1)) + \left(\hat{F}^{ext}(t)\right) P_{ext}(t)$$

where  $\left(\hat{F}(t)\right)$  and  $\left(\hat{F}^{ext}(t)\right)$  are the notations for the following concatenated matrices.

$$\left(\hat{F}(t)\right) = i \begin{pmatrix} \dots & \dots & \dots \\ \dots & D_i F^{i,j}(t) & \dots \\ \dots & \dots & \dots \end{pmatrix}_{\substack{i=1\dots n \\ j=1\dots n}}, \left(\hat{F}^{ext}(t)\right) = \begin{pmatrix} D_1 F^{i,ext}(t) \\ \dots \\ D_n F^{i,ext}(t) \end{pmatrix}$$

To simplify the computation, we will assume that the filters  $F^{i,j}(t)$  and  $F^{i,ext}(t)$  are independent of time. The solution of the system is

$$(X(n)) = \left(\hat{F}\right)^n (X(0)) + (X_e) \quad (4)$$

Where  $(X_e)$  is the equilibrium solution given by a static situation where  $(X(t)) = (X(t-1))$ . We find that  $(X_e) = \left(1 - \left(\hat{F}\right)\right)^{-1} \left(\hat{F}^{ext}\right) P_{ext}(t)$ .

The stability of the dynamics and its specificities will depend on  $\hat{F}$  eigenvalues. Since little can be inferred from the general case, we will study in greater detail a practical example in section 5.

#### 4.4 The unconscious as a strategic agent

**Setup and information** As in the static case, one structure, say  $S_i$ , among  $S_j$  with  $j = 1, \dots, n$ , is strategic. For the sake of simplicity, filters  $F^{i,j}(t)$  and  $F^{j,ext}(t)$  are independent of time.

For all  $j \neq i$ , the optimization problem is unchanged and leads to the optimal response to the parameters:  $X_j(t) = D_j P_j(t)$  for  $j \neq i$ .

However,  $S_i$  now differs from others structures in its optimization.

Consider  $S_i$ 's intertemporal utility:

$$U_i(t) = \sum_{n \geq 0} \beta^n V_i(t+n)$$

with

$$V_i(t) = -\frac{1}{2} X_i^t(t) A_i X_i(t) + X_i^t(t) B_i P_i(t) - \frac{1}{2} P_i^t(t) C_i P_i(t)$$

The important point here is that  $S_i$  actions  $X_i(t)$  also influence  $V_i(t+n)$  for all  $n$ , through  $P_i(t+n)$ .

Actually, all periods are interdependent in the optimization problem :  $X_i(t)$  influences other structures' parameters at time  $t+1$ . Doing so, it impacts  $X_j(t+1)$ , that in turn influences  $P_k(t+2)$ , for  $k \neq j$  at time  $t+2$ . The process goes on indefinitely, spreading  $X_i(t)$ 's effects over all subsequent periods.

$S_i$ , being strategic, knows this. At time  $t$ , it will take all subsequent periods into account, and, to do so, will forecast the future parameters.  $E_{S_i,t}O(t+n)$  will denote  $S_i$ 's forecast at time  $t$  of an arbitrary quantity  $O(t+n)$ . We will assume these forecasts to be linearly dependent on  $S_i$ 's information at time  $t$ . We also assume that  $S_i$  knows the whole set of information at time  $t$ , i.e. all utility functions, external signals, past actions and filters : it knows the whole vector of parameters ( $P(t)$ ), the filters  $F^{i,j}(t)$  and  $F^{j,ext}(t)$  and  $S_j$ 's utilities<sup>11</sup>.

$S_i$  linear forecasts of the future external parameters, given his set of information, are  $E_{S_i,t}P_{ext}(t+n) = F_{t,t+n}^{ext}(P(t))$ .  $F_{t,t+n}^{ext}$  is thus the matrix that expresses  $E_{S_i,t}P_{ext}(t+n)$  as a linear function of the present information ( $P(t)$ ).

$E_{S_i,t}P_{ext}(t+n)$  is the only forecast  $S_i$  needs to build all its expectations about the future : actions at time  $t+n$  depend on the signals received ( $P(t+n)$ ), themselves depending on actions at  $t+n-1$ , up until time  $t$ , where optimization is performed.  $S_i$  being rational, it is able, through the dynamic equations, to reconstruct the whole set of future actions and parameters, provided that the future exogenous parameters are forecasted.

$\lambda$ 's unconscious knows how his conscious works, and, provided accurate forecasts of future external signals, will be able to infer  $\lambda$ 's conscious actions and parameters. He will therefore be able to manipulate him.

**The strategic agent's optimisation** To optimize  $U_i(t)$  and resolve  $\frac{\partial}{\partial X_i(t)}U_i(t) = 0$  requires deriving the dependence of  $E_{S_i,t}P_i(t+n)$  on  $X_i(t)$ .

Appendix A shows that  $E_{S_i,t}P_i(t+n)$  is given by:

$$\begin{aligned} E_{S_i,t}P_i(t+n) &= \Pi_i M^n(P(t)) \\ &+ \Pi_i \sum_{l=1}^n M^{l-1} E_{S_i,t}(VX_i(t+n-l) + WP_{ext}(t+n+1-l)) \end{aligned} \quad (5)$$

where  $\Pi_i$  is the projection operator on  $S_i$  space. Matrices and vectors  $(M_{i,j})_{i=1,\dots,n,j=1,\dots,n}$ , ( $P(t)$ ),  $V$ ,  $W$  are concatenations of the structures data<sup>12</sup>.

Each block of this concatenation is defined by :

$$\begin{aligned} M_{j,k} &= F^{j,k}(t)(1 - \delta_{i,k})D_k \text{ for } j \neq i \\ M_{i,k} &= F^{i,k}(t)D_k \\ V_j &= F^{j,i}(t) - \delta_{ij}F^{i,i}(t) \\ W_j &= F^{j,ext}(t) \end{aligned}$$

<sup>11</sup>A lack of information about the parameters could be modeled. It could be done through setting some parameters estimation to 0. However, assuming that  $S_i$  fully knows ( $P(t)$ ) does not impair the generality of our results.

<sup>12</sup>See section 4.3.

The formula 5 yields the dependence of  $P_i(t+n)$  in the action  $X_i(t)$  :

$$\frac{\partial}{\partial X_i(t)} E_{S_{i,t}}(P_i(t+n)) = \Pi_i M^{n-1} V$$

Appendix B shows that, inserting this result in  $S_i$ 's intertemporal optimization problem,  $\frac{\partial}{\partial X_i(t)} U_i(t) = 0$  yields :

$$0 = \beta_i \left( \sum_{n \geq 1} \Pi_i (\beta_i M)^{n-1} V (B_i G_i - C_i \Pi_i) [(M + V G_i)^n + M_n^{ext}] (P(t)) \right) - A_i X_i(t) + B_i P_i(t)$$

whose solution for  $S_i$ 's optimal action is

$$X_i(t) = G_i(P(t))$$

where the matrix  $G_i$  satisfies :

$$G_i = D_i \Pi_i + \beta_i \Pi_i \left( \sum_{n=0}^{\infty} (\beta_i M)^n V (B_i G_i - C_i \Pi_i) \left( (M + V G_i)^{n+1} + M_{n+1}^{ext} \right) \right)$$

$\Pi^{ext}$  is the matrix that projects  $(P(t))$  on  $P_{ext}(t)$ .  $M_n^{ext}$  is given by

$$M_n^{ext} = \sum_{l=1}^n (M + V G_i)^{l-1} W F_{t,t+n}^{ext} \Pi^{ext}$$

The defining equation for  $G_i$  being matricial, it has usually to be solved numerically. However, an explicit solution will be given in section 5 for a basic example.

**The dynamics of the system** We have found the strategic structure action :  $X_i(t) = G_i(P(t))$ , and we further know that the non strategic structure choose  $X_j(t) = D_j P_j(t)$  for  $j \neq i$ . We can gather all these dynamical equations in a unique concatenated system. This leads to:

$$(X(t)) = \left( \hat{F}(t) \right) (X(t-1)) + \left( \hat{F}^{ext}(t) \right) P_{ext}(t)$$

where  $(X(t)) = \begin{pmatrix} X_1(t) \\ \dots \\ X_n(t) \end{pmatrix}$  encompasses all actions at time  $t$  and

$$\left( \hat{F}(t) \right) = i \begin{pmatrix} \dots & \dots & \dots \\ \dots & \left( \hat{F}(t) \right)_{j,k} & \dots \\ \dots & \dots & \dots \end{pmatrix}_{\substack{i=1 \dots n \\ j=1 \dots n}}, \left( \hat{F}^{ext}(t) \right) = \begin{pmatrix} \left( \hat{F}^{ext}(t) \right)_1 \\ \dots \\ \left( \hat{F}^{ext}(t) \right)_n \end{pmatrix}$$

with the block matrices  $\left(\hat{F}(t)\right)_{j,k}$  and  $\left(\hat{F}^{ext}(t)\right)_1$  defined as:

$$\begin{aligned}\left(\hat{F}(t)\right)_{j,k} &= \left(D_i F^{j,k}(t)\right)_{\substack{j=1\dots n, j \neq i \\ k=1\dots n}} \\ \left(\hat{F}(t)\right)_{i,k} &= \left(G_i F^{j,k}(t)\right)_{k=1\dots n} \\ \left(\hat{F}^{ext}(t)\right)_j &= \left(D_j F^{j,ext}(t)\right)_{j=1\dots n, j \neq i} \\ \left(\hat{F}^{ext}(t)\right)_i &= \left(G_i F^{i,ext}(t)\right).\end{aligned}$$

The solution of this dynamical system is thus strictly similar to the previous non strategic case. The equilibrium is obtained by setting  $(X(t)) = (X(t-1))$ , and leads to :

$$(X_e) = \left(1 - \left(\hat{F}(t)\right)\right)^{-1} \left(\hat{F}^{ext}(t)\right) P_{ext}(t)$$

This allows to solve the system at each time  $t = n$ :

$$(X(n)) = \left(\hat{F}(n)\right)^n ((X(0)) - (X_e)) + (X_e)$$

This dynamics is similar to the non strategic case. However, one major difference arises here: whereas, in the non strategic case,  $S_i$ 's action was based on its own, sole parameters,  $S_i$ 's action is now based on the whole set of parameters at its disposal, including those pertaining to other structures. This will have a significant impact on the results, as shown now.

## 5 A dynamic two structures model

Our general set up can be straightforwardly applied to the simple example of section 3.

Let us consider a dynamic version of this model with  $n = 2$  agents. Agent  $i$  intertemporal utility is written:

$$U_i(t) = \sum_{n \geq 0} \beta_i^n V_i(t+n)$$

and utilities  $V_i(t)$  at time  $t$  are

$$V_1(t) = -\frac{1}{2}(a_1(t) - a_0)^2 - a_2(t) a_1(t) \quad (6)$$

$$V_2(t) = -\gamma \frac{a_1^2(t)}{2} + \alpha a_2(t) a_1(t) - \frac{a_2^2(t)}{2} \quad (7)$$

These are merely (1) and (2) with time dependent actions.

These formulas are encompassed in our general set up. Actually  $V_1(t)$  and  $V_2(t)$  can be cast in the form **3** when we identify:

$$\begin{aligned}
X_1(t) &= a_1(t), X_2(t) = a_2(t) \\
P_1(t) &= \begin{pmatrix} a_0 \\ a_2(t) \end{pmatrix}, P_2(t) = a_1(t), P_{ext}(t) = a_0 \\
A_1 &= 1, B_1 = \begin{pmatrix} 1 & -1 \end{pmatrix}, C_1 = 0, A_2 = 1, B_2 = \alpha, C_2 = \gamma \\
F^{1,2}(t) &= \begin{pmatrix} 0 \\ 1 \end{pmatrix}, F^{2,1}(t) = 1, F^{1,ext}(t) = \begin{pmatrix} 1 \\ 0 \end{pmatrix}, F^{2,ext}(t) = 0
\end{aligned}$$

## 5.1 Non strategic agents

Applying 4 yields the dynamical system:

$$\begin{pmatrix} a_1(t) \\ a_2(t) \end{pmatrix} = \begin{pmatrix} 0 & -1 \\ \alpha & 0 \end{pmatrix} \begin{pmatrix} a_1(t-1) \\ a_2(t-1) \end{pmatrix} + \begin{pmatrix} a_0 \\ 0 \end{pmatrix}$$

The equilibrium denoted by the upperscript  $NS$  is:  $\begin{pmatrix} a_1^{NS} \\ a_2^{NS} \end{pmatrix} = \begin{pmatrix} \frac{a_0}{\alpha+1} \\ \frac{\alpha a_0}{\alpha+1} \end{pmatrix}$ .

As a consequence, the dynamics can be solved as:

$$\begin{pmatrix} a_1(n) \\ a_2(n) \end{pmatrix} = \begin{pmatrix} 0 & -1 \\ \alpha & 0 \end{pmatrix}^n \left( \begin{pmatrix} a_1(0) \\ a_2(0) \end{pmatrix} - \begin{pmatrix} a_1^{NS} \\ a_2^{NS} \end{pmatrix} \right) + \begin{pmatrix} a_1^{NS} \\ a_2^{NS} \end{pmatrix}$$

The equilibrium replicates the static case, and consequently yields the same conclusions. However, the dynamics reveals an important additional result: the equilibrium stability depends on  $\alpha$ .

When  $0 < \alpha < 1$ , the unconscious reaction to the conscious action is mild, and the equilibrium is relatively quickly reached, although with some oscillations. Each agent react to others actions in a damped way, leading to a stable equilibrium.

When  $\alpha > 1$ , the unconscious overreacts to the conscious actions. It induces an increasing and explosive oscillatory movement around the equilibrium, that results in big losses for both agents : each attempt by  $\lambda$ 's conscious to go out will be met by a stronger unconscious reaction. The conscious will try to counter the perceived strain, further increasing the strain. The effect being multiplicative, the strain will increase exponentially. We can assume that this unstable dynamics will lead to a real disorder.

If, on the contrary,  $\lambda$ 's unconscious propensity to react to the conscious is smaller ( $0 < \alpha < 1$ ) the dynamics, will gradually fade away and settle down to the equilibrium.

## 5.2 $S_2$ as a strategic agent

Appendix C shows that the unconscious optimal action is:

$$a_2(t) = \alpha a_0 + c a_1(t-1)$$

with:

$$\begin{aligned}
a &= \frac{\sqrt{(\gamma\beta_2^2 + 1)^2 - 4\alpha^2\beta_2^2} + \gamma\beta_2^2 - 1}{\sqrt{(\gamma\beta_2^2 + 1)^2 - 4\alpha^2\beta_2^2} + 2\alpha\beta_2^2 + (\gamma\beta_2^2 + 1)} \\
c &= \frac{(\gamma\beta_2^2 + 1) - \sqrt{(\gamma\beta_2^2 + 1)^2 - 4\alpha(\alpha\beta_2^2)}}{2(\alpha\beta_2^2)} \tag{C}
\end{aligned}$$

Whereas, in the non strategic case,  $S_2$  was reacting to  $a_1(t-1)$  only,  $S_2$  optimal action now depends on  $a_0$ ,  $S_1$ 's objective.

Why is it so? Inserting  $a_2(t) = aa_0 + ca_1(t-1)$  in 6,  $S_1$ 's "effective" utility can now be rewritten as:

$$V_1(t) = -\frac{1}{2}(a_1(t) - a_0)^2 - aa_0a_1(t) - [ca_1(t-1)]a_1(t)$$

The last term is the cost  $S_2$  imposes on  $S_1$  described in section 3.

$-\frac{1}{2}(a_1(t) - a_0)^2 - aa_0a_1(t)$  is a utility whose optimum is  $a_1(t) = a_0 - aa_0$ .

$S_2$  has clearly manipulated  $S_1$  system of representations by reducing its optimal goal. It is now  $a_0 - aa_0$ .

Rewriting  $-\frac{1}{2}(a_1(t) - a_0)^2 - aa_0a_1(t)$  as  $-\frac{1}{2}(a_1(t) - (a_0 - aa_0))^2$  up to an irrelevant constant shows more clearly this downward shift imposed on  $a_0$ .  $S_2$  gains in this : in both utilities,  $a_2(t)$  was a cost. By shifting  $S_1$ 's goal,  $S_2$  reduces its own strain and in turn increases its utility.

$S_2$ 's ability, as a strategic agent, to manipulate  $S_1$  through its goal is of course inherent to the model, where actions are signals and signals modify the other structures parameters. However, this particular example clearly shows how a conscious rational agent representations and goals can shift over time, under the action of the unconscious.

On top of the strain to go out, our agent  $\lambda$  now experiences a decrease in its preference to go out : the unconscious has succeeded in shifting his tastes.

The dynamical system is now straightforward :

$$\begin{pmatrix} a_1(t) \\ a_2(t) \end{pmatrix} = \begin{pmatrix} 0 & -1 \\ c & 0 \end{pmatrix} \begin{pmatrix} a_1(t-1) \\ a_2(t-1) \end{pmatrix} + \begin{pmatrix} a_0 \\ aa_0 \end{pmatrix}$$

and its equilibrium is given by:  $\begin{pmatrix} a_1^S \\ a_2^S \end{pmatrix} = \begin{pmatrix} \frac{(1-a)a_0}{c+1} \\ \frac{(c+a)a_0}{c+1} \end{pmatrix}$ .

A direct computation yields  $a_1^S < a_1^{NS}$ ,  $a_2^S > a_2^{NS}$  : compared to the non strategic case,  $S_2$  imposes an equilibrium closer to  $a_1 = 0$ . Both structures face a higher cost  $a_2$ . It is however optimal for  $S_2$ , since it compensates this loss through the reduction of  $a_1$  and  $a_0$ . We will detail below the mechanism  $S_2$  uses to reach its goal.

Given the equilibrium, the dynamics is easily solved:

$$\begin{pmatrix} a_1(n) \\ a_2(n) \end{pmatrix} = \begin{pmatrix} 0 & -1 \\ c & 0 \end{pmatrix}^n \left( \begin{pmatrix} a_1(0) \\ a_2(0) \end{pmatrix} - \begin{pmatrix} a_1^S \\ a_2^S \end{pmatrix} \right) + \begin{pmatrix} a_1^S \\ a_2^S \end{pmatrix} \tag{8}$$

When  $c < 1$ , the system converges toward the equilibrium. Inspecting C shows that the dynamics is converging for all values of  $\alpha$  and  $\gamma$ , except in a intermediate zone:

$1 < \gamma < \frac{1}{\beta_2^2}$  and  $\sqrt{\gamma} > \alpha > \frac{(\beta_2^2 \gamma + 1)}{(1 + \beta_2^2)}$  : by manipulating  $S_1$ ,  $S_2$  stabilizes the dynamics.

When  $\gamma$  and thus  $\alpha$  are lower than 1, the non strategic mechanism remains. The unconscious underreacts to the conscious actions, quickly converging toward the equilibrium.

When  $\gamma > \frac{1}{\beta_2^2}$  the unconscious overreacts but, since  $a$  is large enough,  $S_2$  will successfully manipulate  $S_1$ 's goal. The larger  $\gamma$ , the more  $a_0 - aa_0$ ,  $S_1$ 's goal, will be reduced. This goal having now shifted close to 0,  $S_1$ 's incentive to react is reduced. The equilibrium is reached and stable.

Only in the intermediate zone does  $S_2$ 's overreaction fail to reduce  $S_1$ 's goal.  $S_1$  reacts in turn, which triggers an explosive dynamics. It may seem surprising that the strategic  $S_2$  could induce a suboptimal outcome in the long run. However, one should recall that at time  $t$ ,  $S_2$  optimizes a discounted sum of utilities. The discount rate is  $\beta_2$ . The dynamics is explosive because  $S_2$ , at least partly, disregard the future<sup>13</sup>.

### 5.3 Utility comparison

The utilities at time  $t$ , in the equilibrium for both situations can be compared. Denoting  $NS$  and  $S$  the non strategic and strategic agent respectively, in the non strategic case, we find that:

$$\begin{aligned} (1 - \beta_1) U_1^{NS} &= -\frac{1}{2} \left( \frac{a_0}{1 + \alpha} - a_0 \right)^2 - \alpha \left( \frac{a_0}{1 + \alpha} \right)^2 = -\frac{1}{2} \frac{\alpha(\alpha + 2)}{(\alpha + 1)^2} a_0^2 < 0 \\ (1 - \beta_2) U_2^{NS} &= -\gamma \frac{\left( \frac{a_0}{1 + \alpha} \right)^2}{2} + \frac{\alpha^2 \left( \frac{a_0}{1 + \alpha} \right)^2}{2} = \frac{1}{2} a_0^2 \frac{\alpha^2 - \gamma}{(\alpha + 1)^2} < 0 \end{aligned}$$

and in the strategic case.:

$$\begin{aligned} (1 - \beta_1) U_1^S &= -\frac{1}{2} \frac{(\alpha + 2\alpha\beta_2^2 + \gamma\beta_2^2 + 2)(\alpha + \gamma\beta_2^2)}{(\alpha + \alpha\beta_2^2 + \gamma\beta_2^2 + 1)^2} a_0^2 \\ (1 - \beta_2) U_2^S &= -\frac{(\gamma - \alpha^2)(\gamma\beta_2^4 + 2\alpha\beta_2^2 + 1)}{2\alpha^2\beta_2^4 + 4\alpha^2\beta_2^2 + 2\alpha^2 + 4\alpha\gamma\beta_2^4 + 4\alpha\gamma\beta_2^2 + 4\alpha\beta_2^2 + 4\alpha + 2\gamma^2\beta_2^4 + 4\gamma\beta_2^2 + 2} a_0^2 \end{aligned}$$

The difference in utility is therefore :

<sup>13</sup>The intertemporal utility is  $\sum_{n \geq 0} \beta_2^n V_2(t + n)$ . 8 shows that the variables and  $V_2(t + n)$  behave as  $c^{\frac{n}{2}}$  with  $c > 1$ . We can show that  $c < \frac{1}{\beta_2}$ . Therefore,  $\beta_2^n V_2(t + n)$  behaves as  $\beta_2^{\frac{n}{2}}$  : the intertemporal utility converges.

$$\begin{aligned}
(1 - \beta_1) (U_1^S - U_1^{NS}) &= \frac{1}{2} \beta_2^2 (\alpha^2 - \gamma) \frac{2\alpha + \alpha^2 \beta_2^2 + 2\alpha \beta_2^2 + \gamma \beta_2^2 + 2}{(\alpha + 1)^2 (\alpha + \alpha \beta_2^2 + \gamma \beta_2^2 + 1)^2} a_0^2 < 0 \\
(1 - \beta_2) (U_2^S - U_2^{NS}) &= \frac{1}{2} \frac{\beta_2^2 (\alpha^2 - \gamma)^2 (2\alpha + \gamma \beta_2^2 + 2 - \beta_2^2)}{(\alpha + 1)^2 (\alpha + \alpha \beta_2^2 + \gamma \beta_2^2 + 1)^2} a_0^2 > 0
\end{aligned}$$

As expected, there is a gain of utility for  $S_2$  to be strategic. This, of course, occurs at the expense of  $S_1$  whose actions and goals have been distorted.

For the conscious, this loss in utility may well remain unnoticed : an external witness would observe it, but as long as  $\lambda$ 's preferences have been shifted, there no reason why he should bother. Only by keeping in mind the past objective will  $\lambda$  realize that a shift has indeed occurred, as well as a loss.

## 6 Conclusion

We modeled the "irrational" behavior of a single conscious agent, such as systematic bias, sudden changes in the preferences, by introducing interactions between two rational agents, the "conscious" and "unconscious" .

This model has several distinctive features : First, each agent's utility function is a function of other agents' actions through some costs. Second, a strategic agent can, to some extent, modify the other's goal. The "conscious" can therefore experiment a switch in his representations, directly leading him to reduce his objective, and in turn reducing the cost of the "unconscious'" own action.

This result sheds a different light on the agent rationality debate. A model with a single agent can hardly explain why goals can endogeneously change through time or why some systematic bias appear in actions . In our setup, these change and bias are the consequence of an interaction : an external signal induces a conscious behavior, such as utility formation, information process, planning...The unconscious uses its own grid of interpretation to react both to the outer world and to the conscious action. Doing so, it induces a change in the conscious' perception of reality, such as a change in costs, or goals. This could explain why seemingly irrational or inconsistent reactions appear through time. It is the result of the manipulation of the conscious goals by a strategic unconscious. Moreover, We have shown that the systematic bias between the conscious' actions and objectives reveals the difference in the two agents' information processes. [31] advocates that emotions, associated to a welfare loss, are the sign of this differential of information. If the unconscious is invisible, and, like a blackhole, reveals itself through its manifestations, then emotions could provide a practical way to explore and interpret the unconscious' grid of interpretation.

We can now take a different look at the independance of irrelevant alternatives and give another explanation to the behavior of the consumer ordering his meal. We could admit that the conscious agent preferences are ranked in the following order : *Fish*  $\succ$  *Beef*  $\succ$  *Chicken*. When beef and chicken are



suggested, his rational choice would be to choose beef first, and fish when it is suggested. However the unconscious agent may have other views on the matter. He may consider that hiding his favorite choice is optimal : under this assumption, when two options exist, chicken is optimal. When fish is suggested, beef becomes optimal, being an equilibrium between the two agents preferences. We can of course not prove this assertion for a particular individual. Yet it is characteristic of our approach.

More generally, this set up provides a first step towards introducing multi-rationality to describe neural processes and choice formation. If automatic processes, such as motor control, can be described by a utility set up including some kind of forecasting, more complex processes could well be described by the interaction of autonomous and possibly strategic structures. This view fits relatively well with two recent approaches in neurosciences: the cognitive unconscious, and the complex system approach, namely the cooperative and "self-assembly" view of the mind.

Moreover, our set up of interactions raises the question of an effective unity of the individual. It also questions the aggregation of structures over time. We may wonder if interacting structures can learn from each other to reach a cooperative equilibrium. This question will be left for further research.

## References

- [1] Antonelli Giovanni. Sulla teoria matematica della economia politica. Pise, 1886 ; traduit en Preferences, utility and demand, J. Chipman, L. Hurwicz, M. Richter, & H. Sonnenschein (eds). New York : Harcourt Brace Jovanovich, 1971.
- [2] Allais, Maurice. « Le comportement de l'homme rationnel devant le risque: critique des postulats et axiomes de l'école américaine ». *Econometrica*, vol. 21, n° 4, 1953,503-546.
- [3] Baltussen, Guido, Behavioral Finance: An Introduction (January 13, 2009).
- [4] Becker, Gary (1968). "Crime and Punishment: An Economic Approach". *The Journal of Political Economy* 76, pp. 169-217.
- [5] Becker, G. S. 1962. Irrational behavior and economic theory, *Journal of Political Economy* 70:1-13.
- [6] Berg, Nathan & Gigerenzer, Gerd, 2010. "As-if behavioral economics: Neo-classical economics in disguise?," MPRA Paper 26586, University Library of Munich, Germany.
- [7] *The Methodology of Economics, or How Economists Explain*, revised edition, 1993. Cambridge.
- [8] Martin Browning and Pierre-André Chiappori: "Efficient Intra-Household Allocations: a General Characterization and Empirical Tests", *Econometrica* 66 (1998), p. 1241-1278
- [9] Camerer, Colin F. 1999. "Behavioral Economics: Reunifying Psychology and Economics." *Proceedings of the National Academy of Sciences*, 96, 10575-77.
- [10] Gerald Edelman (2004) *Wider than the Sky. The Phenomenal Gift of Consciousness*. Yale University Press.
- [11] Ivar Ekeland & Jon Elster, *Théorie économique et rationalité*, Vuibert - SPS 2011
- [12] Jon Elster, 1998. "Emotions and Economic Theory," *Journal of Economic Literature*, American Economic Association, vol. 36(1), pages 47-74, March.
- [13] Jon Elster, *Le désintéressement : Traité critique de l'homme économique*, tome 1, 2009
- [14] Jon Elster, *L'irrationalité : Traité critique de l'homme économique*, tome 2, 2010

- [15] Farmer J. Doyne, Geanakopulos John. «The virtues and vices of equilibrium and the future of financial economics ». arXiv:q-fin/0803.2996.
- [16] Freud, S. *The Interpretation of Dreams (Die Traumdeutung)*, 1899/1900
- [17] Friedman, M. (1953). *Essays in positive economics*. Chicago : University of Chicago Press.
- [18] Glimcher P.W. (2003) *Decisions, Uncertainty, and the Brain : The Science of Neuroeconomics*, Cambridge, MIT Press.
- [19] Gold J.I. and Schadlen M.N. (2000) "Representation of a perceptual decision in developing oculomotor commands", *Nature*, 404, p. 390-394.
- [20] Gold J.I. and Schadlen M.N. (2001) "Neural computations that underlie decisions about sensory stimuli", *Trends in Cognitive Science*, 5, p. 10-16.
- [21] Gosselin, P. Lotz,A. Wyplosz,C. When Central Banks Reveal Future Interest Rates: Alignment of Expectations Vs.Creative Opacity. *International Journal of Central Banking* Vol. 4 Nbr 3, 145 (2008)
- [22] Gosselin, P. Lotz,A. Wyplosz,C. How Much Information should Interest Rate-Setting Central Banks Reveal ? *NBER International Seminar on Macroeconomics*, 8 (2007)
- [23] *The New Unconscious*. Edited by Ran R. Hassin, James S. Uleman and John A. Bargh, *Oxford Series in Social Cognition and Social Neuroscience*). Oxford University Press.
- [24] Hayek, Friedrich A. 1944. 'On Being an Economist', an address given to economics students at the London School of Economics in 1944, first published in *The Trend of Economic Thinking: Essays on Political Economists and Economic History* (vol. III of *The Collected Works of F. A. Hayek*), edited by W. W. Bartley and Stephen Kresge (Chicago: University of Chicago Press, 1991), pp. 35–48.
- [25] James, W. (1890). *Principles of Psychology* (2 vols). New York. Holt
- [26] Kahneman, D. E. and Tversky, A. 1979. 'Prospect theory: An analysis of decision under risk', *Econometrica* 47: 263–91.
- [27] Katona, G. 1951. *Psychological Analysis of Economic Behavior*. New York : McGraw-Hill.
- [28] Keynes, J. M. 1936. *The General Theory of Employment, Interest and Money*. London : Macmillan
- [29] Lea, S.E.G., Webley, P. & Young, B.M. (1992). (Eds). *New Directions in Economic Psychology: theory, experiment and application*. Cheltenham: Edward Elgar.

- [30] Rick, Scott and Loewenstein, George F., The Role of Emotion in Economic Behavior (January 3, 2007).
- [31] Lotz, Aileen, An Economic Approach to the Self: The Dual Agent (April 11, 2011). Available at SSRN: <http://ssrn.com/abstract=1798999> or <http://dx.doi.org/10.2139/ssrn.1798999>
- [32] Neurodynamics of cognition and consciousness. Leonid I. Perlovsky , Robert Kozma. Editors Understanding complex systems Springer: complexity. Springer
- [33] Neural Networks and Intellect: Using Model-Based Concepts: Leonid I. Perlovsky, Oxford University Press
- [34] Rabin, M. 2002. 'A perspective on psychology and economics', European Economic Review 46: 657–85.
- [35] Neuroeconomie, Christian Schmidt, Odile Jacob économie.
- [36] Shiv, Baba, Loewenstein, George, Bechara, Antoine, Damasio, Hanna & Damasio, Antonio R. (2005). Investment behavior and the negative side of emotion, Psychological Science, 16(6), 435-439.
- [37] Simon, Herbert A, 1986. "Rationality in Psychology and Economics," Journal of Business, University of Chicago Press, vol. 59(4), pages S209-24, October.
- [38] Slutsky, E. E. (1915). "Sulla teoria del bilancio del consumatore". Giornale degli Economisti 51 (July): 1-26;
- [39] Smith, A. 1759/2002. The Theory of Moral Sentiments, 6th edn, K. Haakonssen (ed.), Cambridge: Cambridge University Press.
- [40] The Encyclopedia of Neuroscience. Editor: Larry R. Squire, Elsevier Science & Technology, 2006
- [41] Tarde, G. (1902). La psychologie économique (2 vols). Paris : Alcan.

## 7 Appendix A. Form of the parameters.

We compute  $E_{S_{i,t}}P_i(t+n)$ ,  $S_i$ 's expectation at time  $t$ , given the expression for  $P_i(t+n)$ :

$$P_i(t+n) = \sum_{j \neq i} F^{i,j}(t+n) X_j(t+n-1) + F^{i,ext}(t+n) P_{ext}(t+n)$$

ad thus

$$E_{S_{i,t}}P_i(t+n) = E_{S_{i,t}} \sum_{j \neq i} F^{i,j}(t+n) X_j(t+n-1) + E_{S_{i,t}} F^{i,ext}(t+n) P_{ext}(t+n)$$

Using  $S_j$ 's solution to the optimization problem, i.e.  $X_j(t) = D_j P_j(t)$ , for  $j \neq i$ , leads to:

$$\begin{aligned} E_{S_{i,t}}P_i(t+n) &= E_{S_{i,t}} \sum_{j \neq i} F^{i,j}(t+n) D_j P_j(t+n-1) + E_{S_{i,t}} F^{i,ext}(t+n) P_{ext}(t+n) \\ &= \sum_{j \neq i} F^{i,j}(t+n) D_j E_{S_{i,t}} P_j(t+n-1) + E_{S_{i,t}} F^{i,ext}(t+n) P_{ext}(t+n) \end{aligned}$$

On the other hand, we now need  $P_j(t+n-1)$  for  $j \neq i$ . It is given by:

$$E_{S_{i,t}}P_j(t+n) = E_{S_{i,t}} \sum_{k \neq j} F^{j,k}(t+n) X_k(t+n-1) + F^{j,ext}(t+n) P_{ext}(t+n)$$

Using again the non strategic agents optimization  $X_k(t) = D_k P_k(t)$  for  $k \neq i$

$$\begin{aligned} E_{S_{i,t}}P_j(t+n) &= \sum_{k \neq j, k \neq i} F^{j,k}(t) D_k E_{S_{i,t}} P_k(t+n-1) \\ &\quad + F^{j,i}(t) E_{S_{i,t}} X_i(t+n-1) + E_{S_{i,t}} F^{j,ext}(t+n) P_{ext}(t+n) \end{aligned}$$

leads to

$$E_{S_{i,t}}(P(t+n)) = E_{S_{i,t}} M(P(t+n-1)) + V E_{S_{i,t}} X_i(t+n-1) + E_{S_{j,t}} W E_{S_{j,t}} P_{ext}(t+n)$$

with  $M_{j,k} = F^{j,k}(t)(1 - \delta_{i,k}) D_k$  for  $j \neq i$  and  $M_{i,j} = F^{i,j}(t) D_j$  and  $V_j = F^{j,i}(t) - \delta_{ij} F^{i,i}(t)$ ,  $W_j = F^{j,ext}(t)$ . This allows to find ultimately  $E_{S_{i,t}}(P(t+n))$  and  $E_{S_{i,t}}P_i(t+n)$ . First:

$$E_{S_{i,t}}(P(t+n)) = M^n(P(t)) + \sum_{l=1}^n M^{l-1} E_{S_{i,t}}(V X_i(t+n-l) + W P_{ext}(t+n+1-l))$$

and projecting on  $S_i$  space yields:

$$E_{S_{i,t}}P_i(t+n) = \Pi_i M^n(P(t)) + \Pi_i \sum_{l=1}^n M^{l-1} E_{S_{i,t}}(V X_i(t+n-l) + W P_{ext}(t+n+1-l))$$

where  $\Pi_i$  is the projection operator on agent  $i$  space.

## 8 Appendix B. Optimization problem

The optimization problem for the strategic agent  $S_i$  is :

$$\begin{aligned} 0 &= E_{S_i,t} \frac{\partial}{\partial X_i(t)} \left( \sum_{n \geq 0} \beta_i^n V_i(t+n) \right) \\ &= E_{S_i,t} \left( \sum_{n \geq 1} \beta_i^n \frac{\partial P_i(t+n)}{\partial X_i(t)} (B_i X_i(t+n) - C_i P_i(t+n)) \right) - A_i X_i(t) + B_i P_i(t) \end{aligned}$$

Using the expression for  $\frac{\partial P_i(t+n)}{\partial X_i(t)}$  in the text leads to the following equation:

$$0 = E_{S_i,t} \beta_i \left( \sum_{n \geq 1} \Pi_i (\beta_i M)^{n-1} V (B_i X_i(t+n) - C_i P_i(t+n)) \right) - A_i X_i(t) + B_i P_i(t) \quad (9)$$

All  $S_i$  forecasts at time  $t$  can be built from  $(P(t))$ . Moreover  $S_i$ 's optimization problem at time  $t$  does not depend on past actions  $\therefore$  at each period, the past is unaccounted for. The situation is reset to 0. Ultimately, all equations in our problem are linear. Consequently,  $S_i$  will linearly choose  $X_i(t)$  as a linear function of  $(P(t))$  :  $X_i(t) = G_i(P(t))$ .

To find  $G_i$ , we first replace  $X_i(t)$  in  $E_{S_i,t}(P(t+n))$ :

$$\begin{aligned} E_{S_i,t}(P(t+n)) &= E_{S_i,t} M(P(t+n-1)) \\ &\quad + V E_{S_i,t} X_i(t+n-1) + W E_{S_i,t} P_{ext}(t+n) \\ &= E_{S_i,t} (M + V G_i)(P(t+n-1)) + W E_{S_i,t} P_{ext}(t+n) \end{aligned}$$

This is solved recursively to yield:

$$\begin{aligned} E_{S_i,t}(P(t+n)) &= (M + V G_i)^n (P(t)) \quad (10) \\ &\quad + E_{S_i,t} \sum_{l=1}^n (M + V G_i)^{l-1} W P_{ext}(t+n+1-l) \quad (11) \end{aligned}$$

Turning back to the optimization equation and introducing 10 in 9 gives:

$$\begin{aligned} 0 &= \beta_i \left( \sum_{n \geq 1} \Pi_i (\beta_i M)^{n-1} V (B_i G_i - C_i \Pi_i) \left[ [(M + V G_i)^n + M_n^{ext}] (P(t)) \right] \right) \\ &\quad - A_i X_i(t) + B_i P_i(t) \end{aligned}$$

where  $M_n^{ext} = \sum_{l=1}^n (M + V G_i)^{l-1} W F_{t,t+n}^{ext} \Pi^{ext}$  and  $\Pi^{ext}$  is the projector sending  $(P(t))$  on  $P_{ext}(t)$ .

Isolating  $X_i(t)$  is straightforward

$$X_i(t) = D_i P_i(t) + \beta_i \Pi_i \left( \sum_{n=0}^{\infty} (\beta_i M)^n V (B_i G_i - C_i \Pi_i) (M + V G_i)^{n+1} \right) (P(t))$$

and yields, after identification with  $G_i(P(t))$ :

$$G_i = D_i \Pi_i + \beta_i \Pi_i \left( \sum_{n=0}^{\infty} (\beta_i M)^n V (B_i G_i - C_i \Pi_i) \left( (M + V G_i)^{n+1} + M_{n+1}^{ext} \right) \right)$$

as claimed in the text.

## 9 Appendix C. The strategic agent in the 2 structure model

The vector  $(P(t))$  is three dimensional, since:

$$(P(t)) = \begin{pmatrix} P_1(t) \\ P_2(t) \end{pmatrix} = \begin{pmatrix} a_0 \\ a_2(t) \\ a_1(t) \end{pmatrix}.$$

Thus, the optimal action  $a_2(t) = G_2(P(t))$  encompasses three parameters. We therefore let  $G_2 = \begin{pmatrix} a & b & c \end{pmatrix}$ .

We first find  $(a, b, c)$  : using the parameters values as defined in the text yields the matrices needed for the identification of  $G_2$  .

$$B_2 = D_2 = \alpha, C_2 = \gamma, D_1 = \begin{pmatrix} 1 & -1 \end{pmatrix}$$

$$V = \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix}, \Pi_2 = \begin{pmatrix} 0 & 0 & 1 \end{pmatrix}, M = \begin{pmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 1 & -1 & 0 \end{pmatrix}$$

$$M + VG_2 = \begin{pmatrix} 0 & 0 & 0 \\ a & b & c \\ 1 & -1 & 0 \end{pmatrix}, V(B_2G_2 - C_2\Pi_2) = \begin{pmatrix} 0 & 0 & 0 \\ a\alpha & b\alpha & c\alpha - \gamma \\ 0 & 0 & 0 \end{pmatrix}$$

$$F_{t,t+n}^{ext} \Pi^{ext} = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{pmatrix}, M_n^{ext} = \frac{1 - (M + VG_2)^n}{1 - (M + VG_2)} F_{t,t+n}^{ext} \Pi^{ext}$$

Note that  $(\beta_i M)^n = 0$  for  $n \geq 2$ .

Replacing these matrices in the equation defining  $G_2$

$$G_2 = D_2\Pi_2 + \beta_2\Pi_2 \left( \sum_{n=0}^{\infty} (\beta_2 M)^n V(B_2G_2 - C_2\Pi_2) \left( (M + VG_2)^{n+1} + M_{n+1}^{ext} \right) \right)$$

leads to the system of equations for  $a, b, c$ :

$$\begin{aligned} a &= \beta_2^2 ((1-a)(\gamma - c\alpha) - b\alpha(c + ab + a) - a\alpha) \\ b &= \beta_2^2 (b\alpha(c - b^2) - b(\gamma - c\alpha)) \\ c &= \alpha - \beta_2^2 (c\alpha b^2 + c(\gamma - c\alpha)) \end{aligned}$$

Given that  $\alpha^2 - \gamma < 0$ , one can check that there are two solutions for the vector  $(a, b, c)$ :

$$c = \frac{(\gamma\beta_2^2 + 1) \pm \sqrt{(\gamma\beta_2^2 + 1)^2 - 4\alpha(\alpha\beta_2^2)}}{2(\alpha\beta_2^2)}, b = 0, a = \frac{\beta_2^2(\gamma - c\alpha)}{1 + \beta_2^2(\alpha + \gamma - c\alpha)}$$



Moreover, taking into account  $c \simeq \alpha$  for  $\beta_2 \rightarrow 0$ , the solution is :

$$\begin{aligned} a &= \frac{\sqrt{(\gamma\beta_2^2 + 1)^2 - 4\alpha^2\beta_2^2} + \gamma\beta_2^2 - 1}{\sqrt{(\gamma\beta_2^2 + 1)^2 - 4\alpha^2\beta_2^2} + 2\alpha\beta_2^2 + (\gamma\beta_2^2 + 1)} \\ b &= 0 \\ c &= \frac{(\gamma\beta_2^2 + 1) - \sqrt{(\gamma\beta_2^2 + 1)^2 - 4\alpha(\alpha\beta_2^2)}}{2(\alpha\beta_2^2)} \end{aligned}$$