

Rationality, Autonomy and Coordination: the Sunk Costs Perspective

Matteo Bonifacio¹, Paolo Bouquet¹, Roberta Ferrario¹, and Diego Ponte¹

University of Trento

Via Belenzani, 12 – I-38100 Trento, Italy

`bonifacio@itc.it`, `bouquet@dit.unitn.it`, `ferrix@cs.unitn.it`, `ponte@itc.it`

Abstract. Our thesis is that an agent¹ is autonomous only if he is capable, within a non predictable environment, to balance two forms of rationality: one that, given goals and preferences, enables him to select the best course of action (means-ends), the other, given current achievements and capabilities, enables him to adapt preferences and future goals. We will propose the basic elements of an economic model that should explain how and why this balance is achieved: in particular we underline that an agent's capabilities can often be considered as partially sunk investments. This leads an agent, while choosing, to consider not just the value generated by the achievement of a goal, but also the lost value generated by the non use of existing capabilities. We will propose that, under particular conditions, an agent, in order to be rational, could be led to perform a rationalization process of justification that changes preferences and goals according to his current state and available capabilities. Moreover, we propose that such a behaviour could offer a new perspective on the notion of autonomy and on the social process of coordination.

1 Rationality in Traditional Theories of Choice

Traditional theories of choice are based upon the paradigm that choosing implies deciding the best course of action in order to achieve a goal [23]. Goals are generally considered as given or, at least, they are selected through an exogenous preference function which assigns an absolute value to each possible state of the world [21]. Potential goals, once ordered according to preferences, are selected by comparing each absolute value with the cost of its achievement. In particular, the agent will commit to the goal that maximizes the difference between the absolute benefit of the goal and the cost of using the capabilities that are needed. This means-ends paradigm subtends a type of rationality that March defines as anticipatory, causative, consequential, since an agent anticipates the consequences of his actions through a knowledge of cause-effect relationships [7] [18]. Here, as underlined by Castelfranchi, autonomy is viewed in the restrictive sense of *executive autonomy*: the only discretionality the agent possesses is about the

¹ In this paper we do not intentionally draw any distinction between artificial and human agents, but we rather discuss the concept of agent in general.

way in which a goal is to be achieved and not about which kind of goal should be preferable; in this sense, even if an agent selects a goal, he is unable to direct the criteria of the selection. The interest of the agent is always reconducible to the one of the designer and, as Steels concludes referring to Artificial Agents, “AI systems built using the classical approach are not autonomous, although they are automatic . . . these systems can never step outside the boundaries of what was foreseen by the designer because they cannot change their own behaviour in a fundamental way.” [28]. Sometimes, as we will propose, autonomy and rationality lie in the possibility to change our mind on what is good and what is bad on the basis of current experience; basically, this is equivalent to the possibility to decide not just how to achieve a goal, but rather which goal is to be achieved and, moreover, which is preferable.

2 Another Perspective on Rationality: Ex-post Rationalization

Another way to look at rationality, that March defines as ex-post rationality or rationalization offers an opposite perspective on decision making [19](see also [31]). At the extreme, it envisions an agent as somebody who first acts and then justifies his actions defining appropriate goals and preferences in order to be consistent to his current achievements. More realistically, it presents an agent not as somebody who is only able to be rational in the sense of setting appropriate courses of action, but also in the sense of changing his mind about what is preferable when planned achievements become unrealistic [20]. Such an agent is able to learn not just in terms of finding better ways to achieve a goal but also in terms of finding goals that are more appropriate to his capabilities. As we will see afterwards, in an environment characterized by a non predictable evolution, an agent who has a partial and perspective view of the world [4] [15] will often come to situations in which ex-post rationalization is more rational than setting a plan for the achievement of given goals [21]. We will propose that this process hides an economic principle of reuse and conservation that could lead an agent to try to fit the world, rather than pretending the world to be appropriate to him.

Moreover, if non predictability is the main reason to be rational and autonomous in the sense just stated, ex post rationalization is also an opportunity for the agent to be like this. In particular, whenever an environment is ambiguous and undefined, equally ambiguous and undefined is the definition of what is good and what is bad. More simply, we often describe a situation as good or bad not because it is so in itself, but rather because of our interpretation and our convenience; as commonly said, it is a question of perspective [9]. Here “rationalization” appears as an opportunity, since it can hide a powerful tool to learn from experience, which produces as outcome the possibility of seeing the world from different perspectives. As underlined by [21], this view represents decisions as constructive interpretations, since they “are often reached by focusing on reasons that justify the selection of one option over another. Different frames,

contexts, and elicitation procedures highlight different aspects of the options and bring forth different reasons and considerations that influence decisions”. More simply, thanks to rationalization, an agent can understand that, under a different point of view, a mistake or an unlucky event could become an opportunity to learn. In this sense the “value” of a goal appears to be a choice rather than an evidence [18].

3 The Rationale of Rationalization: Sunk Costs, Economies of Reuse and Irreversibility

In this work we propose that ex-post rationalization can find a rational justification in the sunk cost effect which derives from the co-occurrence of two conditions that could characterize an agent’s capability: economies of reuse and partial irreversibility. To start, an agent has a set of means that we can view both as capabilities when used to perform actions and as resources when used to develop or acquire new capabilities. Under this perspective, at a given moment, an agent’s set of capabilities can be interpreted as the result of an investment of resources. Traditional theories of choice assume that when calculating the net benefit of a decision, we take into account just the costs of those resources and capabilities that will be used in order to achieve the goal [17]. Said differently, the value of non used means are not to be considered when selecting ends.

The observation of the process of decision shows something different. In particular:

- when calculating the cost of a decision, we consider not just the cost of those capabilities that we use, but also the cost generated by the non use of some other one. This is because the generation of a capability implies the sustenance of some fixed costs that can be amortized through its repeated use. In fact, in presence of fixed costs, reuse implies a decrease in the unitary cost of each re-utilization. In economic theory, this effect is called the economies of reuse effect. Consequently, not using a resource implies a loss of value generated by the lost opportunity of a cost saving. Moreover, each time we use a capability we exploit its economies of reuse effect and at the same time we loose the correspondent effect of those we do not use;
- each capability, when considered as a resource that can be used to acquire another capability, displays a rate of irreversibility. This is because when trying to transform it into another, we can sustain a loss in value if the resource is difficult to manipulate or if it is difficult to find a buyer on the market. In general, if a resource is totally reversible (for example currency), it can be sold on the market and the economies of reuse effect has only marginal impact on the decision of the agent. On the other hand, if totally irreversible, the resource will completely display its economies of reuse effect; if this is not used, its owner will suffer the loss of a potential cost saving [10].

In all these cases, a resource which is characterized both by economies of reuse and a rate of irreversibility is considered a sunk cost and it generates the

effect that “paying for the right to use a good will increase the rate at which the good will be utilized *ceteris paribus*” [30]. This hypothesis will be referred to as the sunk cost effect [2]² and, as underlined by [24], “This tendency (the sunk cost consideration) contradicts a basic principle in economics that past costs and benefits should be irrelevant to current decisions [13]”. In [17] Johnstone writes: “For decision-making purposes, sunk costs are strictly irrelevant. This is a law of economic logic justified by the argument that because no action (current or future) can avert or reduce a sunk cost, no sunk cost can be attributed to or have any relevance to current or future action. It is evident, however, that for many of us, the edict that sunk costs must be ignored is hard to accept, if not as a matter of logic then at least in application.”³. It is common sense that each resource displays some of these effects, generating on the one hand an incentive to reuse and, on the other, an incentive to create markets to enhance reversibility. Moreover, it is important to notice how such effects can affect decisions. In fact, when deciding, an agent will consider not just the costs currently sustained, but also those losses in value generated by the non use of sunk investments. Now the point is that in a non predictable environment, while pursuing a goal, an agent can be led to develop and acquire capabilities that, to some extent, have no use in order to achieve his current goal. This is probable in case of very turbulent environments and it is enhanced when the agent is in an advanced stage of his life or is particularly experienced in the domain of the decision he is facing; the former circumstance leads the agent to unforeseen situations and to the generation of redundant capabilities; the latter to the growing accumulation of sunk investments and costs, since those that are more reversible were probably used during the earlier stages of the agent’s life or experience.

4 Generating Preferences and Goals

In this section we will give evidence on how a decision making process that considers sunk costs, can lead to an ex-post rationalization whereby an agent manipulates preferences in order to justify his current state. As we will propose, through this endogenous process of preferences formation, an agent becomes able to select and pursue goals which are not predictable a-priori.

4.1 Generating Preferences

As [14] suggests, the “commitment to a current course of action is a function of the comparison between the perceived utility of continuing with the action and the perceived utility of withdrawal and/or changing the action”. If sunk costs are taken into consideration in determining which option is preferable in a decision, an agent can face a situation in which the cost of changing his mind about what

² Some authors also call this “escalation effect” [25] [26] and it has been applied in studying political decisions constrained by the presence of sunk costs (such as military escalation).

³ See also [22] [29].

is preferable is lower than the cost of going on in the pursuit of his intentions. In particular, in the decision function the weight of sunk costs overcomes the one of current opportunities. Some authors have remarked that this tendency can lead to irrational behaviours such as the “irrational escalation” [26], whereby social agents could irrationally justify the current failure in order to explain past choices and “save their face” [5]. In this sense, decisions generating investments influence future choices that are constrained by the need to preserve past investments. On the other hand, some other has remarked that such a behaviour is also to be considered as a manifestation of coherence and rationality [11]. In fact, seen the other way round, (i.e. as a rational behaviour), such a situation offers the agent the opportunity to do something qualitatively different: instead of reasoning on how to pursue an unrealistic goal, he could realistically consider his current state as appropriate to his capabilities. In this case, a current unexpected situation could be viewed by the agent as a proper ex-post goal, and remaining where he is could be more rational than moving. As argued before, instead of reasoning about means necessary to achieve ends which were shown to be irrational, he rationalizes his current state as an end which is appropriate to his means. Under this perspective, the sunk cost effect is an attempt to demonstrate the rationality of behaviours that are otherwise not explained and thus labelled as “irrational” by traditional theories of rationality ⁴. In this sense, for example, [3] argued that the escalation effect, which is typically adducted as an example of irrational course of action, stems from a “don’t waste” decision rule (see also [1]). We suggest that this attitude leads to a process of retrospective self-justification [27] that implies a change in preferences. In fact, in order to be consistent with his history, an agent who rationalizes his current state needs to change his preferences accordingly. As a matter of fact, in order to justify the (even ex-post) adoption of a goal, a rational agent needs to express such a goal as desirable. If not, the agent would display the inconsistent behaviour of choosing a goal which is not desirable. This necessity leads the agent to invert his reasoning process on preferences which are turned from fixed tools used to select goals to variable matters that are adapted to (now fixed) current achievements. As clearly stated by [11], “When people realize they are in situations that they have never considered before, they do not judge themselves to be irrational. Instead, they simply try to decide what beliefs and preferences to adopt (if any)”. In other words, it is rational for the agent to perform a counterfactual process [12] [16] that could be expressed by a sentence such as: “What should I have preferred in

⁴ Traditional studies on the sunk cost effect predicts that the more a resource is irreversible, the more a sunk cost effect will be displayed. According to this view, the agent will decide irrationally since he will consider the value of something that cannot be reversed. Differently, the authors are currently involved in an experiment to demonstrate that both complete and null reversibility leads a rational agent not to consider sunk costs in decisions as predicted by the classic rational model. On the other hand, when the rate of reversibility is partial and ambiguous, the agent has the opportunity to reuse and thus to exploit the value of a sunk investment. In general we will propose that considering sunk costs in calculating decisions is a rational strategy when an agent is facing ambiguity.

order to be satisfied with the state of the world I am currently in?” or “What should I have preferred in order to desire to reach a goal that is consistent with the current state of the world?”.

4.2 Generating Goals

As anticipated, through this counterfactual process, the agent asks himself which goal he should have been committed to in order to be, given his resources/capabilities, satisfied with what he currently is. In a particular sense, such a process represents the first attempt for an agent to endogenously generate a goal; the goal is the already achieved current state that, only ex-post, can be viewed as a goal. That is to say, we propose that the first manifestation of *goal autonomy* is the rationalization of an unexpected and undesired current state, turned into a desired one. Since rationalization is exactly driven by sunk costs, this first goal, by definition, will display the peculiarity of exploiting the value of current sunk investments. We underline that this original process of goal and preferences creation is not an abstract process of imagining new possible worlds and preferences but a concrete exercise that uses the presence of a goal (the current state) as a tool to derive a proper set of preferences.

At a first sight, the behaviour of this agent could appear to be intrinsically conservative. Once the weight of past investments overcomes the weight of immediate value, the agent stops where he is, due to the retrospective justification of his state as a desired one. Even if we think that, in time, conservative behaviours are an underlying tendency of the agent, we propose that, within this tendency, new deliberative behaviours can emerge. Here we give just an example of how new goals can emerge, derived by the idea of Castelfranchi [8] of social adoption. In fact, given the new set of preferences, the agent is now able to assign new “values” to every state of the world that is accessible to his knowledge. In this way an agent is able to reorder the states of the world on his new preference scale and set new goals. On the other hand, he has now changed his beliefs on means-ends relations and on how a particular set of capabilities (the ones used to reach the current state) can be used in order to achieve a goal. In particular, he changes his beliefs on what is preferable and on which means are needed in order to get to a particular goal (the current state) [11]. For example, we might argue that the agent, observing other agent’s situations, discovers that another state of the world displays a net benefit (considering the new preferences and existing sunk costs) which is higher than the one of preserving the current state. Now he will adopt the new possible state as the new goal. Again, as above, this process can lead to the acquisition of new resources and to the possibility that, on the path to the goal, the agent happens again to be in unexpected states of the world that might influence, through the evaluation of sunk investments, his preferences.

5 Conclusions: Autonomy and Coordination

This conclusion leads us to some considerations on the notion of autonomy. We agree with the one proposed by Castelfranchi [6] whereby an agent is autonomous if he is able to choose goals on the basis of a personal interest. Here we underline the need that such interest is an endogenous production of the agent rather than something exogenously given by a designer (in the case of artificial agents) or by another human or metaphysical entity (in the case of human agents). Moreover, Castelfranchi remarks that the definition of autonomy currently used in artificial agents literature is referable to the weaker notion of *executive autonomy* (as opposed to *goal autonomy*): an agent is autonomous if he is able to choose among alternative courses of action. As he underlines, this kind of autonomy could resolve both in a type of slavery (from some external utility function) and in a form of irrationality (pursuing some other's interest when this is conflicting with our own is irrational). We strongly believe as Castelfranchi in the idea that an agent, if not goal autonomous, is not autonomous at all and, moreover, potentially irrational. Now the question becomes how such a type of autonomy can emerge in order to design, if possible, agents that can display *goal autonomy* through the generation of endogenous preferences and the consequent adoption of non a-priori predictable goals. In this work, we sketch the lines of a model that could give an answer. In particular our thesis is that an agent, in order to be rational, endogenously develops preferences and goals that are consistent to his "emerging" interest. This interest is the consequence of an unforeseen evolution of his life that led to the generation of sunk costs that need to be considered in decision making. Such an evolution, assuming a non predictable environment, leads to the autonomous formation of preferences that are not predictable a-priori, and that are the rational consequence of an economic principle of reuse. Through preference formation, new goals become desirable while old ones are abandoned. In this sense we say that the agent, at a certain stage of his life, in order to be rational, needs to become autonomous (and form new preferences).

One last point addresses the way in which this approach could be used to interpret some fundamental aspects of an agent's sociality, in particular those aspects that involve coordination with other agents. Specifically, if we consider coordination efforts as investments that display a sunk cost effect, we can explain the persistency of social relations. We refer to the observation that social relations among social agents are less prone to opportunism than what is predicted by traditional utilitaristic theories. In fact, whenever the current value of a relation is lower than the cost of keeping it, an agent should break such relation. As a matter of fact, social relations seem to be more persistent than this. A way to interpret such persistency without recurring to exogenous factors (such as social norms) [8], is provided by the perspective of sunk costs. Here a social relationship is viewed as a resource and capability that displays an economies of reuse effect (as a consequence of the initial investment in creating the relation) and a rate of irreversibility (since a social relationship cannot always be transformed into another). As a consequence, an agent, in order to achieve his goal, will tend to reuse and justify current established relations before creating new ones.

References

1. H. R. Arkes. The psychology of waste. *Journal of Behavioral Decision Making*, 9, 1996.
2. H. R. Arkes and P. Ayton. The Sunk Cost and Concorde Effect: are Humans Less Rational Than Lower Animals? *Psychological Bulletin*, 125, 1999.
3. H. R. Arkes and C. Blumer. The psychology of sunk cost. *Organizational Behavior and Human Performance*, 35:129–140, 1985.
4. M. Benerecetti, P. Bouquet, and C. Ghidini. Contextual reasoning distilled. *Journal of Theoretical and Experimental Artificial Intelligence*, 12(3):279–305, 2000.
5. J. Brockner, J. Z. Rubin, and E. Lang. Face-saving and entrapment. *Journal of Experimental Social Psychology*, 17:68–79, 1981.
6. C. Castelfranchi. Guarantees for autonomy in cognitive agent architecture. In M. Wooldridge and N. R. Jennings, editors, *Intelligent Agents: Theories, Architectures, and Languages (LNAI Volume 890)*, pages 56–70. Springer-Verlag: Heidelberg, Germany, 1995.
7. M. D. Cohen, J. G. March, and J. P. Olsen. A garbage can model of organizational choice. *Administrative Science Quarterly*, 17:1–25, 1972.
8. R. Conte and C. Castelfranchi. *Cognitive and social Action*. UCL Press, 1995.
9. R. L. Daft and K. E. Weick. Toward a model of organizations as interpretation systems. *Academy of Management Review*, 9(2):284–295, 1984.
10. B. di Bernardo and E. Rullani. *Il management e le macchine*. il Mulino, 1990.
11. J. Doyle. Rationality and its roles in reasoning. *Computational Intelligence*, 8(2):376–409, 1992.
12. R. Ferrario. Counterfactual reasoning. In *Proceedings of the 3th International and Interdisciplinary Conference on Modelling and using Context*, volume 2116. Springer-Verlag, 2001.
13. R. H. Frank. *Microeconomics and Behavior*. McGraw Hill: New York, 1994.
14. H. Garland and S. Newport. Effects of Absolute and Relative Sunk Cost on the Decision to Persist with a Course of Action. *Organizational Behavior and Human Decision Processes*, 48:55–69, 1991.
15. C. Ghidini and F. Giunchiglia. Local models semantics, or contextual reasoning = locality + compatibility. *AI*, 127(2):221–259, April 2001.
16. M. L. Ginsberg. Counterfactuals. *AI*, 30(1):35–79, 1986.
17. D. Johnstone. The reverse sunk cost effect and explanation: rational and irrational. <http://www.departments.bucknell.edu/management/apfa/papers/17Johnstone.pdf>, 2000.
18. J. G. March. How decisions happen in organizations. *Human Computer Interaction*, 6:95–117, 1991.
19. J. G. March. *A primer on decision making : how decisions happen*. The Free Press, 1994.
20. J. W. Payne, R. Bettman, and R. J. Johnson. Behavioral decision research: a Constructive Processing Perspective. *Annual Review of Psychology*, 4:87–131, 1992.
21. E. Shafir, I. Simonson, and A. Tversky. Reason-based choice. *Cognition*, 49, 1993.
22. M. Shefrin. *Beyond greed and fear: understanding behavioral finance and the psychology of investing*. Harvard Business School Press, 2000.
23. H. A. Simon. *Reason in human affairs*. Stanford University Press, 1983.
24. D. Soman. The Mental Accounting of Sunk Time Costs: Why Time is not Like Money. *Journal of Behavioral Decision Making*, 14:169–185, 2001.

25. B. Staw. Attribution of the causes of performance: an alternative interpretation of cross-sectional research on organizations. *Organizational Behaviour and Human Performance*, 13:414–432, 1975.
26. B. Staw. Knee-deep in the big muddy: a study of escalating commitment to a chosen course of action. *Organizational Behaviour and Human Performance*, 16:27–44, 1976.
27. B. Staw and J. Ross. Understanding behavior in escalation situations. *Science*, 246:216–220, 1989.
28. L. Steels. When are robots intelligent autonomous agents? *Journal of Robotics and Autonomous Systems*, 15:3–9, 1995.
29. R. Thaler. Toward a theory of consumer choice. *Journal of Economic Behaviour and organization*, 1:39–60, 1980.
30. R. Thaler. *Quasi rational economics*. Russel Sage foundation, 1994.
31. E.K. Weick. *The social psychology of organizing*. McGraw-Hill, Inc., 1979.