

Trust and Transitivity: how trust-transfer works

Rino Falcone and Cristiano Castelfranchi

National Research Council– Institute of Cognitive Sciences and Technologies
Via San Martino della Battaglia, 44 00185 - Roma, Italy
{rino.falcone, cristiano.castelfranchi}@istc.cnr.it

Abstract. Transitivity in trust is very often considered as a quite simple property, trivially inferable from the classical transitivity defined in mathematics, logic, or grammar. In fact the complexity of the trust notion suggests evaluating the relationships with the transitivity in a more adequate way. In this paper, starting from a socio-cognitive model of trust, we analyze the different aspects and conceptual frameworks involved in this relation and show how different interpretations of these concepts produce different solutions and definitions of trust transitivity.

1 Introduction

Many different approaches and models of trust were developed in the last 15 years [1, 2, 3, 4, 5, 6, 7]: they contributed to clarify many aspects and problems about trust and trustworthiness, although many issues remain to be addressed and some elementary but not so trivial trust properties are left in a contradictory form.

One of them is the problem of trust transitivity. If X trusts Y , and Y trusts Z : What about the trust relationship between X and Z ? Different and sometimes diverging answers were given to this problem. The question is not only theoretically relevant; it is very relevant from the practical point of view.

In fact one of the main problems of trusting agents in an open and massive multi-agent world is the necessity of exploiting the cumulated trust by other trustees (who we trust) for trusting agents who they know but for us unknown. But, how does this trust-transfer work? Which are the sophisticated basic cognitive mechanisms for doing this?

In this paper we will present an analysis of the trust transitivity when a socio-cognitive model of trust is taken into consideration. Through this kind of model we are able to evaluate and partially cope with the complexity that the concept of transitivity introduces when applied to the trust relationship.

We cannot in fact just rely on our direct and personal experience; this would restrict too much our interaction possibilities and "market", and also our "trust capital" (how many people and how much trust us, and would be interested in exchanging with us).

This precious information not only is preserved in our memory, from previous interactions; but it is generalized (to groups and categories [14, 15]) for making predictions about new people; it is acquired from the others by observation and use of 'signals' (I see that X exchange with Y), or by communication (advices, recommendations, reputation, etc); it is acquired from Y 's exhibition of his skills, products, virtues, qualities (signaling and "marketing").

So, trust (and distrust) information is circulating a lot as a precious good, and trust has a very active social dynamics in the sense of *transfer* from one trustor to another, from one trustee to another, from context to context, from task to task. This phenomenon is so important that there is the need of clarifying and explicitly modeling the specific mechanism of this trust "circulation".

2 A socio-cognitive model of trust

In our socio-cognitive model of trust [8, 9] we consider trust as a relational construct between the trustor (X), the trustee (Y), about a defined (more or less specialized) task (τ). Introducing a formal operator *Trust*, representing the trust notion, we can write: where are also explicitly present both the X 's goal (g_X , respect to which trust is

$$Trust(X Y C \tau g_X)$$

tested/activated) and the role of the context (C) in which the relationship is going to happens. In fact, the successful performance of the task τ will satisfy the goal g_X . *Task* is the performance X needs and expects from Y ; what X relies on Y for.

Using Meyer, van Linder, van der Hoek et al.'s logics [13] we can introduce the X 's mental ingredients of trust. They mainly are the goal g_X , and a set of *main* trustor's beliefs:

$$Bel(X Can_Y(\tau)), Bel(X Will_Y(\tau)), Bel(X ExtFact_Y(\tau))$$

where:

Bel is the classical modal operator for representing agents' beliefs;

Can_Y(τ) means that Y is potentially able to do τ (in the sense that, under the given conditions, is competent, has the internal powers, skills, know-how, etc); (and in the above formula this is what is believed by X);

Will_Y(τ) means that, under the given conditions, Y has potentially the attributions for being willing, persistent, available, etc., on the task τ (and in the above formula this is what is believed by X);

ExtFact_Y(τ) means that there are a set of external conditions either favoring or hindering Y realizing the task τ (and in the above formula is clear that also this is believed by X).

In our model we also consider that trust can be *graded*: X can have a *strong trust* that Y will realize the task (maybe 0.95 in the range (0,1)); or even just a *sufficient trust* that Y will achieve it (maybe 0.6 with a threshold of 0.55; and so on with other possible values). For this we have introduced a trust quantification, calling it the *Degree of Trust*, **DoT_{XY τ}** , and, in general, a *threshold* (σ) to be overcome from this

DoT_{XY τ} .

Given the previous analysis of the main components of the trust attitude (g_X , **Bel** (**X Can_Y(τ)**), **Bel** (**X Will_Y(τ)**), **Bel** (**X ExtFact_Y(τ)**)), our model is also able to evaluate how a specific degree of trust is, on its turn, resultant from the several *quantifications* of these components:

$$DoT_{XY\tau} = f(DoC_X (Opp_Y(\tau)), DoC_X (Competence_Y(\tau)), DoC_X (Willingness_Y(\tau)))$$

where: f is in general a function that preserves monotonicity;

$DoC_X (Opp_Y(\tau))$ is the X 's degree of credibility about the external opportunities (positively or negatively) interfering with Y 's action in realizing the task τ ;

$DoC_X (Competence_Y(\tau))$ is the X 's degree of credibility about the Y 's ability/competence to perform τ ;

$DoC_X (Willingness_Y(\tau))$ is the X 's degree of credibility about the Y 's willingness to perform τ .

We are ignoring the subjective certainty of the pertinent beliefs (how much sure is X of its evaluative beliefs about that specific quality about either Y 's or the environment, that is a meta-belief; in first approximation we can say that this factor is integrated with the other).

At the same time we are ignoring for now the value of the goal (g_X). In fact, this value (its level of relevance) should have a specific and complex influence on the degree of the threshold (σ): increasing the relevance of the goal both should increase the caution of the trustor of not entrusting a too unreliable trustee (pushing to increase the threshold), and should increase the need of not missing the opportunity to achieve it (pushing to decrease the threshold). Of course, both factors of trustor's personality and the set of potential viable alternatives than that trustee have also to be considered for defining the complex relationships between threshold and goal.

So trivially X will trust Y about the task τ if $DoT_{XY\tau} > \sigma$

that means that a set of analogous conditions must be realized about the other quantitative elements ($DoC_X (Opp_Y(\tau))$, $DoC_X (Competence_Y(\tau))$, $DoC_X (Willingness_Y(\tau))$).

We do not consider in this paper the detailed analysis of how the degree of trust is resulting by the more elementary components. We also omit of considering the potential correlations among the different components.

In addition, we should also say that all the components above showed (competence, willingness, external factors) are in fact not the more elementary ones; they can be described as resultant by more elementary reasons/components. For example, the competence dimension could be considered as constituted by the *know-how*, the *self-confidence*, and the *ability* sub-components; and so on.

Also, introducing the concept of *trustworthiness degree* of an agent we call $Trustworthiness_{YC}(\tau)$

the Y 's trustworthiness about the task τ , in the context C . In general the trustworthiness of an agent is a property of that agent¹. We have in fact two different meanings of this concept:

- an *objective trustworthiness*, what the agent is actually able and willing to do in standard conditions; his/her actual reliability on a more or less specific task.
- a *subjective trustworthiness*, the perceived reliability of the trustee by another agent; it could be different for the different trustors.

The situation is also more complex, because the objective trustworthiness may not be constant with respect to the trustor but the same trustee Y could be differently trustworthy on the same task with different trustors (for example X or Z) (suppose he

¹ In fact the trustworthiness (as trust) derives from specific properties of the trustee that could change during the time both for intrinsic reasons and for external conditions.

has different motivations for helping/serving X and Z). For this we introduce also the variable X in the operator: $Trustworthiness_{YXC}(\tau)$.

Supposing the situation in which the context is constant, if we have: $DoT_{XY\tau} > \sigma$

it derives that: $Bel(X Trustworthiness_{YX}(\tau) > \Sigma)$

Where Σ is the minimum value of trustworthiness for Y (as believed by X) for delegating to him the task τ . In general $\sigma = \Sigma$.

3 Transitivity in Trust

The transitivity property of trust can be presented as: if X trusts Y , and Y trusts Z : What about the trust relationship between X and Z ?

We are interested to translate this problem in our terms of trust.

First of all, we do not consider the unspecified case “ X trusts Y ” because in our model an agent has to trust another agent with respect a task (either very well defined or less defined and abstract); this task directly derives from the goal the trustor has to reach with the trust attribution. So we have to transform “ X trusts Y ” in “ X trusts Y about τ ”. And given the graded qualification of trust we have that: $DoT_{XY\tau} > \sigma$

this means in particular that X believes that Y is potentially *able* and *willing* to do τ and that the *external conditions* in which Y will perform its task are at least not so opposite to the task realization (may be also they are neutral or favorable).

So this Y 's trustworthiness with respect to X (perceived/believed by X) about τ is based on these specific beliefs.

At the same way “ Y trusts Z ” becomes “ Y trusts Z about τ_i ” (about the difference between τ and τ_i , see later) with the same kind of particular Y 's beliefs about Z and the external conditions.

Also in this case we can say that there is a threshold to be overcome and the condition: $DoT_{YZ\tau_i} > \sigma_1$

successfully satisfied in case of trust attribution.

In order to transfer/adopt trust as mere evaluation, esteem, potential expectation (disposition) it is necessary to at least examine the task, that is, "about what" X trusts Y , and Y trusts Z ; and possibly on which bases (the ascribed qualities). But, in order to transfer trust as 'decision' to rely on, as 'act' of trusting (from the decision of Y to trust/rely on Z to the decision of X to entrust Z), something more is necessary: the degree of trust and its thresholds.

If we have to consider the trust relationship between “ X and Z ” as a consequence of the previous trust relationships between “ X and Y ” and between “ Y and Z ” we have to specify the task on which this relationship should be based (question of *assimilation* between τ and τ_i) and the degree of trust that must be overcome even from X :

$DoT_{XZ\tau} > \sigma_2$ with the consideration of the threshold σ_2 .

The role of the trust threshold is quite complex and can have an overlapping with the ingredients of trust. We strongly simplify in this case considering σ as dependent just from the specific intrinsic characteristics of the trustor (those that define an agent intrinsically: prudent, reckless, and so on) independently from the external

circumstances (on the contrary, these factors affect the degree of trust, by affecting the more elementary beliefs above showed).

So, we can say that in this approximation: $\sigma = \sigma_2$

In the case in which all the agents are defined as having the same intrinsic characteristics (possible for artificial entities), we can also say that: $\sigma = \sigma_1 = \sigma_2$

Moreover, as we just saw, not less important in our approach is that trust is an expectation and a bet *grounded on and justified by* certain beliefs about Y . X trusts Y on the basis of the evaluation of Y 's "virtues/qualities", not just on the basis of a statistical sampling, some probability.

The *evaluations* about the needed "qualities" of Y for τ are the *mediator* for the decision to trust Y . This mediation role is fundamental also in trust transitivity.

Let us now consider the case of the differences between the tasks in the different relationships.

For the trust transitivity the two tasks should be the same ($\tau = \tau_1$). Is this enough?

Suppose for example that there are 3 agents: John, Mary and Peter; and suppose that John trusts Mary about "organizing scientific meetings" (task τ), at the same time Mary trusts Peter about "organizing scientific meetings" (again task τ). Can we deduce that, given the transitivity of trust: John trusts Peter about "organizing scientific meetings"? Is in fact transferable that task evaluation? Given the trust model defined in §2 the situation is more complex and there are possible pitfalls lurking: Mary is the central node for that trust transfer and she plays different roles (and functions) in the first case (when her trustworthiness is about *to realize* the task τ , and in the second case (when her trustworthiness is about *evaluating* the Peter's trustworthiness on the task τ).

The situation is even clearer if we split in the example the two kinds of competences: X trusts very much Y as medical doctor; X knows that Y trusts Z as mechanic; will X trust Z as mechanic? Not necessarily at all: if X believes that Y is a good evaluator of mechanics he will trust Z ; but, if X believes that Y is a very naive in this domain, and is frequently swindled by people, he will not trust Z .

So for considering transitivity of trust as a valid property we have to *assimilate* the task with the *evaluation of that task* itself:

$Bel(X \text{ Trustworthiness}_{YX}(\tau) > \sigma) \text{ implies } Bel(X \text{ Trustworthiness}_{YX}(eval(\tau)) > \sigma)$

In words, X has to believe that if Y is sufficiently trustworthy on the task τ , it is also sufficiently trustworthy on the meta-task of evaluating τ .

The reasons for X trusting Y on the task τ have to be based on beliefs that in some way support also the reasons for trusting Y on the meta-task of evaluating τ (or on a different task τ'). More analytically, the qualities (in both the directions: competence and willingness) that X is attributing to Y about the task τ , support, directly or for analogy, also the qualities necessary for trusting Y about the meta-task of evaluating τ (or on a different task τ').

So resuming we have the *more basic case* (case A) of the relationship between *trust* and *transitivity* so defined (we assume $\sigma = \sigma_X = \sigma_Y$):

<p><i>if</i> iA) $DoT_{XY\tau} > \sigma$ (X trusts Y about τ) and iiA) $DoT_{YZ\tau} > \sigma$ (Y trusts Z about τ) and</p>

iiiA) $Bel(X DoT_{YZ\tau} > \sigma)$ (X believes that Y trusts Z about τ) and
ivA) $Bel(X Trustworthiness_{YX}(\tau) > \sigma)$ **implies** $Bel(X Trustworthiness_{YX}(eval(\tau)) > \sigma)$
(Y is sufficiently trustworthy in the realization of the task and in evaluating others performing that task)
then
vA) $DoT_{XZ\tau} > \sigma$ (X trusts Z about τ)

We have to underline that (iiA) should not necessarily be true, the important thing is that is true (iiiA).

In the case (B) in which the *tasks are different* ($\tau \neq \tau'$), we have:

if
iB) $DoT_{XY\tau} > \sigma$ (X trusts Y about τ) and
iiB) $DoT_{YZ\tau} > \sigma$ (Y trusts Z about τ) and
iiiB) $Bel(X DoT_{YZ\tau} > \sigma)$ (X believes that Y trusts Z about τ) and
ivB) $Bel(X Trustworthiness_{YX}(\tau) > \sigma)$ **implies** $Bel(X Trustworthiness_{YX}(eval(\tau')) > \sigma)$
then
vB) $DoT_{XZ\tau'} > \sigma$ (X trusts Z about τ')

In the case B, the transitivity essentially depends from the implication reported in the formula (iiiB); are there elements in the reasons (believed by X) for trusting Y on the task τ that (in X 's view) are sufficient also for trusting Y on the evaluation of a different task τ' ?

3.1 Competence and Willingness in Transitivity

The need for a careful qualitative consideration of the nature of the link between the trustor and the trustee, is even more serious.

Not only it is fundamental (as we have argued) to make explicit and do not forget the specific "task" (activity, and thus "qualities") X is trusting Y or Z about, but it is even necessary to consider the different dimensions/components of the trust disposition (evaluation), decision, and relation.

In our model, trust has two basic nucleuses:

- (i) Y 's *competence*, ability, for correctly performing the delegated task;
- (ii) Y 's *willingness* to do it, to act as expected.

The two dimensions (and 'virtues' of Y) are quite independent on each other: Y might be very well disposed and willing to do, but not very competent or unable; Y might be very expert and skilled, but not very reliable: unstable, unpredictable, not well disposed, insincere, dishonest, etc.

Now, this (at least) double dimension affect transitivity. In fact, even assuming that the *competence* is rather stable (and that Y is a good evaluator of Z 's *competence*) not necessarily Z 's *willingness* is equally stable and transferable from Y to X . This is a more relation-based dimension.

Y was evaluating Z 's *willingness* to do as expected on the basis of their specific *relation*. Is Z a friend of Y ? Is there a specific benevolence, or values sharing, or gratitude and reciprocity, or obligation and hierarchical relation, etc.? Not necessarily the reasons that Z would have for satisfying Y 's expectation and delegation would be present (or equally important) towards X . X 's relation with Z might be very different. Are the reasons/motives motivating Z towards Y , and making him reliable,

transferable or equally present towards X ? Only in this case it would be reasonable for X to adopt Y 's trustful attitude and decision towards Z .

Only certain kinds of relations will be generalized from Y to X ; for example, if Y trusts Z only because it is an economic exchange, only for Z 's interest in money, reasonably X can become a new client of Z ; or if Y relies on Z because Z is a charitable person, generously helping (without any prejudice and discrimination) poor suffering people, and X is in the same condition of Y , then also X can trust in Z .

3.2 Trust Dynamics affects Transitivity

Moreover, we have shown ([8], [9]) that Z 's willingness, and even ability, can be affected, increased, by Y 's trust and reliance (this can affect Z 's commitment, pride, effort, attention, study, and so on). Z 's *trustworthiness* is improved by Y 's trust and delegation. And Y might predict and calculate this in her decision to rely on Z .

However, not necessarily the effect of Y on Z 's trustworthiness will be produced also by another trustor. Thus, also this will affect "transitivity": suppose that Y 's trust and delegation to Z makes him more trustworthy, improves Z 's willingness or ability (and Y trusts and relies on Z on the basis of such expectation); not necessarily X 's reliance on Z would have the same effect. Thus even if X knows that Y reasonably trusts Z (for something) and that he is a good evaluator and decision-maker, not necessarily X can have the same trust in Z , since perhaps Z 's trustworthiness would not be equally improved by X 's reliance.

4 Transitivity and Trust: Related Work

The concept of trust transitivity has been considered in other approaches.

A relevant example is given from the Josang's approach; he introduces the subjective logics (an attempt of overcoming the limits of the classical logics) for taking in consideration the uncertainty, the ignorance and the subjective characteristics of the beliefs. Using this approach Josang addressed the problem of trust transitivity [10], where it is recognized the intrinsic cognitive nature of this phenomenon. However, the main limits of this approach are that trust is in fact the *trust in the information sources*; and the transitivity regards two different tasks (referred to our formalism: $\tau \neq \tau_j$: X has to trust the evaluation of Y (task τ) with respect Z as realizing another task (task τ_j , for example as mechanic). As we showed before, this difference is really relevant for the transitivity phenomenon. In addition, also the first task (Y as evaluator) is just analyzed with respect to the *property of sincerity* (and this is a confirmation of the constrained view of trust phenomenon; they write: " A 's disbelief in the recommending agent B means that A thinks that B consistently *recommends the opposite* of his real opinion about the truth value of x "; where A , B , and x are, in our terms, respectively X , Y and τ_j). But in trusting someone as evaluator of another agent (with respect to a specific task), I have also to consider his *competence* as evaluator, not just his *sincerity*. Trust is based on ascribed qualities. Y could be completely sincere but at the same time completely inappropriate to evaluate that task.

Other authors [11, 12], developed algorithms for inferring trust among agents not directly connected. These algorithms differ from each other in the way they compute trust values and propagate those values in the networks.

References

- [1] Marsh, S.P., (1994), Formalising Trust as a computational concept. PhD thesis, University of Stirling. Available at: <http://www.nr.no/abie/papers/TR133.pdf>.
- [2] Jonker, C., and Treur, J., (1999), Formal Analysis of Models for the Dynamics of Trust based on Experiences, Autonomous Agents '99 Workshop on "Deception, Fraud and Trust in Agent Societies", Seattle, USA, May 1, pp.81-94.
- [3] Barber, S., and Kim, J., (2000), Belief Revision Process based on trust: agents evaluating reputation of information sources, *Autonomous Agents 2000 Workshop on "Deception, Fraud and Trust in Agent Societies"*, Barcelona, Spain, June 4, pp.15-26.
- [4] Jones, A.J.I., Firozabadi, B.S., (2001), On the characterization of a trusting agent: Aspects of a formal approach. In Castelfranchi, C., Tan, Y.H., (Eds.), *Trust and Deception in Virtual Societies*. Pp. 55-90. Kluwer, Dordrecht.
- [5] Resnick P., Zeckhauser, R., (2002), Trust among strangers in internet transactions: Empirical analysis of eBay's reputation systems. In Baye, R. (Editor), *The economic of the internet and e-commerce*. Vol. 11 of *Advances in Applied Microeconomics*. Elsevier Science.
- [6] Yu, B., Singh, M.P., (2003), Searching social networks. In *Proceedings of the second international joint conference on autonomous agents and multi-agent systems (AAMAS)*. Pp. 65-72. ACM Press.
- [7] Sabater, J. (2003), *Trust and Reputation for Agent Societies*, PhD thesis, Universitat Autònoma de Barcelona.
- [8] Falcone R., Castelfranchi C. (2001), The socio-cognitive dynamics of trust: does trust create trust? In *Trust in Cyber-societies: Integrating the Human and Artificial Perspectives* R. Falcone, M. Singh, and Y. Tan (Eds.), LNAI 2246 Springer. pp. 55-72.
- [9] Castelfranchi C., Falcone R., (2010), *Trust Theory: A Socio-Cognitive and Computational Model*, John Wiley and Sons, Chichester, West Sussex, UK, ISBN 978-0-470-02875-9.
- [10] Bhuiyan T., Josang A., Xu Y., An analysis of trust transitivity taking base rate into account, In: *proceeding of the Sixth International Conference on Ubiquitous Intelligence and Computing*, 7-9 July 2009, University of Queensland, Brisbane, 2009.
- [11] Li, X., Han, Z., Shen, C., Transitive trust to executables generated during runtime, *second International Conference on Innovative Computing, Information and Control*, 2007.
- [12] Golbeck J., Hendler J., Inferring binary trust relationships in web-based social networks, *ACM Transactions on Internet Technology*, 6(4), 497-529, 2006.
- [13] Meyer, J.J. Ch., van der Hoek W., (1992), A modal logic for nonmonotonic reasoning. In W. van der Hoek, J.J. Ch. Meyer, Y. H. Tan and C. Witteveen, editors, *Non-Monotonic Reasoning and Partial Semantics*, pages 37-77. Ellis Horwood, Chichester, 1992.
- [14] Falcone R., Piunti, M., Venanzi, M., Castelfranchi C., From Manifesta to Krypta: The Relevance of Categories for Trusting Others, *ACM Transaction on Intelligent Systems and Technology*, (in press 2012).
- [15] Burnett C., Norman T., and Sycara K., 2010, Bootstrapping trust evaluations through stereotypes. In *9th International Conference on Autonomous Agents and Multi-Agent Systems (AAMAS 2010)*, van der Hoek, Kaminka, Lesperance, Luck and Sen, Eds. Toronto, Canada, 241-248.