

A Logic of Trust and Reputation

Emiliano Lorini

(joint work with A. Herzig, J. F. Hübner, L. Vercouter)

IRIT-CNRS, Toulouse, France

Workshop on Trust, Advice, and Reputation
Toulouse, 22-24 October 2009

Motivation: formal analysis of trust and reputation

- many computational models [Sabater & Sierra 02; Huynh et al. 04]
 - quantitative
 - no qualitative analysis
- few logical models
 - mostly focused on trust in information sources [Liau 03; Jones & Firozabadi 01; Dastani et al. 04]
- cognitive model of trust [Castelfranchi & Falcone 98, 01]
 - informal definitions

⇒ formal definitions?

⇒ relations between trust and reputation?

Contribution: definitions in a BDI logic

[Lorini & Demolombe 2008, 2009; Lorini, Falcone, Castelfranchi 2008; Herzig, Lorini, Hübner, Vercoouter 2009]

- 'top-down' qualitative analysis: reduction of trust to more primitive concepts
 - trust \Rightarrow belief, goal, intention, action
 - analysis in terms of logics of action and time, BDI logics
- two versions of trust
 - 'occurrent trust' vs. 'dispositional trust'
 - refinement of Castelfranchi & Falcone's definition
- from dispositional trust to reputation

Trust: the ingredients [Castelfranchi & Falcone]

“ i trusts j to do α in order to achieve φ ”

- truster i ,
- trustee j ,
- action α of j ,
- goal φ of i .
 - important: trust \neq thinking and foreseeing

Informal definition [Castelfranchi & Falcone]

\Rightarrow trust \approx evaluation of the truster about certain properties of trustee (that are relevant for a goal/task φ)

“ i trusts j to do α in order to achieve φ ” if and only if:

- ① i has the *goal* φ ;
- ② i believes that
 - ① j is *capable* to do α ;
 - ② j has the *power* to achieve φ by doing α ;
 - ③ j *intends* to do α .

\Rightarrow trust defined from four more primitive concepts

\Rightarrow logic of goal, ability, power and intention?

- 1 Two kinds of trust: occurrent vs. dispositional
- 2 Logical formalization
- 3 From trust to reputation
- 4 Conclusion

Two kinds of trust: occurrent vs. dispositional

Occurrent trust vs. dispositional trust

two perspectives on trustee's action α :

- truster believes trustee is going to do α *here and now*
⇒ **occurrent trust**
- truster believes trustee is going to do α *whenever some conditions are satisfied*.
⇒ **dispositional trust**

(cf. occurrent belief vs. dispositional belief [Searle 92])

Occurrent trust

i wants φ to be true and believes *j* is going to perform α here and now so that φ will be ensured

$$\text{OccTrust}(i, j, \alpha, \varphi) \stackrel{\text{def}}{=} \text{OccGoal}(i, \varphi) \wedge \\ \text{Belief}(i, \\ \text{OccCap}(j, \alpha) \wedge \\ \text{OccPower}(j, \alpha, \varphi) \wedge \\ \text{OccIntends}(j, \alpha))$$

\Rightarrow predicates to be defined: *Belief*, *OccGoal*, *OccCap*, *OccPower*, *OccIntends*

Example

1's *occurrent trust* in 2 to send a certain product in view of satisfying 1's goal of possessing the product:

- 1 wants to possess the product,
- 1 believes that
 - 2 is capable to send the product,
 - 2's act of sending the product will result in 1 possessing it,
 - 2 has the intention to send the product.

Dispositional trust

i expects that there will be some situations κ in which he will have the goal of achieving φ (*potential goal*), and believes that in all these situations *j* will ensure φ by doing action α

$$\begin{aligned}
 \text{DispTrust}(i, j, \alpha, \varphi, \kappa) &\stackrel{\text{def}}{=} \text{PotGoal}(i, \varphi, \kappa) \wedge \\
 &\text{Belief}(i, \text{Henceforth} \\
 &((\kappa \wedge \text{OccGoal}(i, \varphi)) \rightarrow \\
 &(\text{OccCap}(j, \alpha) \wedge \\
 &\text{OccPower}(j, \alpha, \varphi) \wedge \\
 &\text{OcclIntends}(j, \alpha)))
 \end{aligned}$$

\Rightarrow predicates to be defined: *PotGoal*, *Belief*, *Henceforth*, *OccGoal*, *OccCap*, *OccPower*, *OcclIntends*

Example

1's *dispositional trust* in 2 (a mechanic) to repair his car so that the car will be in order, in the circumstances in which 1 will ask the mechanic to repair his car

Logical formalization

Occurrent trust and dispositional trust: components

- definiendum:
 - belief,
 - occurrent goal,
 - occurrent capability,
 - occurrent power,
 - occurrent intention.
- definiens: formulas of a modal logic of time, belief, preference, and action
 - 'spell out' predicates $Belief(i, \varphi)$, $OccGoal(i, \varphi)$, $PotGoal(i, \varphi, \kappa)$, $OccCap(j, \alpha)$, $OccPower(j, \alpha, \varphi)$, $OccIntends(j, \alpha)$ using modal operators of time, belief, preference, and action

convention:

- Bel_i , etc. = logical operators

Temporal operators

Henceforth φ = “ φ henceforth holds”

$$\text{Eventually } \varphi \stackrel{\text{def}}{=} \neg \text{Henceforth } \neg \varphi$$

Eventually φ = “ φ eventually holds”

- will be used to define $\text{OccGoal}(i, \varphi)$, etc.

Belief operators

$\text{Bel}_i \varphi$ = “agent i believes that φ ”

$$\text{Poss}_i \varphi \stackrel{\text{def}}{=} \neg \text{Bel}_i \neg \varphi$$

$\text{Poss}_i \varphi$ = “agent i thinks that φ is possible”

- epistemic and doxastic logic [Hintikka 62]
- express truster’s beliefs about trustee’s properties

$$\text{Belief}(i, \varphi) \stackrel{\text{def}}{=} \text{Bel}_i \varphi$$

Goals as preferences about the future

$\text{Pref}_i \varphi$ = “agent i wants φ to be true”

- binary preferences [Cohen & Levesque 90]
- positive introspection: $\vdash \text{Pref}_i \varphi \rightarrow \text{Bel}_i \text{Pref}_i \varphi$
- negative introspection: $\vdash \neg \text{Pref}_i \varphi \rightarrow \text{Bel}_i \neg \text{Pref}_i \varphi$
- realism: $\vdash \text{Bel}_i \varphi \rightarrow \text{Pref}_i \varphi$

$$\text{OccGoal}(i, \varphi) \stackrel{\text{def}}{=} \text{Pref}_i \text{Eventually } \varphi$$

Capability in dynamic logic

$\text{After}_{i:\alpha} \varphi$ = “ φ will be true after *every possible* execution of action α by agent i ”

- propositional dynamic logic (PDL)
- formulas = φ, ψ, \dots = state descriptions
- actions = α, β, \dots = state transition descriptions
 - action \neq formula
 - $i : \alpha$ = action α with author i
- $\text{After}_{i:\alpha} \perp$ = “ i cannot do α ”

$$\text{OccCap}(j, \alpha) \stackrel{\text{def}}{=} \neg \text{After}_{j:\alpha} \perp$$

- “ j is capable to perform α ” = “ j can do α ”

Power in dynamic logic

$$\text{OccPower}(j, \alpha, \varphi) \stackrel{\text{def}}{=} \text{After}_{j:\alpha} \varphi$$

- relates j 's action α with i 's goal φ
- missing: epistemic aspect of power (\Rightarrow 'knowing how to play')

Intention-to-do as preferred action

Does $_{j:\alpha} \varphi$ = “agent i is going to do α and φ will be true afterwards”

- Does $_{j:\alpha} \top = i$ does α

$$\text{OcclIntends}(j, \alpha) \stackrel{\text{def}}{=} \text{Pref}_j \text{Does}_{j:\alpha} \top$$

- Some axioms:

- 1 Does $_{j:\alpha} \varphi \rightarrow \neg \text{After}_{j:\alpha} \neg \varphi$ ($\vdash \text{Does}_{j:\alpha} \top \rightarrow \neg \text{After}_{j:\alpha} \perp$)
- 2 ($\text{Pref}_j \text{Does}_{j:\alpha} \top \wedge \neg \text{After}_{j:\alpha} \perp$) $\rightarrow \text{Does}_{j:\alpha} \top$
- 3 Does $_{j:\alpha} \top \rightarrow \text{Pref}_j \text{Does}_{j:\alpha} \top$

Theorem

$$\vdash (\text{Pref}_j \text{Does}_{j:\alpha} \top \wedge \neg \text{After}_{j:\alpha} \perp) \leftrightarrow \text{Does}_{j:\alpha} \top$$

Formal definition of occurrent trust

$$\text{OccTrust}(i, j, \alpha, \varphi) \stackrel{\text{def}}{=} \text{Pref}_i \text{Eventually } \varphi \wedge \\ \text{Bel}_i (\neg \text{After}_{j:\alpha} \perp \wedge \text{After}_{j:\alpha} \varphi \wedge \\ \text{Pref}_j \text{Does}_{j:\alpha} \top)$$

Example

1's *occurrent trust* in 2 to send a certain product in view of satisfying 1's goal of possessing the product

$$\begin{aligned}
 & \text{OccTrust}(1, 2, \text{send}P1, \text{has}P1) \stackrel{\text{def}}{=} \\
 & \text{Pref}_1 \text{Eventually } \text{has}P1 \wedge \\
 & \text{Bel}_1(\neg \text{After}_{2:\text{send}P1} \perp \wedge \\
 & \text{After}_{2:\text{send}P1} \text{has}P1 \wedge \\
 & \text{Pref}_2 \text{Does}_{2:\text{send}P1} \top)
 \end{aligned}$$

Some properties of occurrent trust

Theorem

- $\vdash \text{OccTrust}(i, j, \alpha, \varphi) \leftrightarrow$
 $(\text{Pref}_i \text{Eventually } \varphi \wedge \text{Bel}_i (\text{Does}_{i:\alpha} \top \wedge \text{After}_{j:\alpha} \varphi))$
- $\vdash \text{OccTrust}(i, j, \alpha, \varphi) \leftrightarrow \text{Bel}_i \text{OccTrust}(i, j, \alpha, \varphi)$
- $\vdash \text{OccTrust}(i, j, \alpha, \varphi) \rightarrow \text{Bel}_i \text{Eventually } \varphi$

Potential goal

$PotGoal(i, \varphi, \kappa) \stackrel{\text{def}}{=}$

$Poss_i \text{ Eventually } (\kappa \wedge Pref_i \text{ Eventually } \varphi)$

Formal definition of dispositional trust

$DispTrust(i, j, \alpha, \varphi, \kappa) \stackrel{\text{def}}{=}$

$Poss_i \text{ Eventually } (\kappa \wedge Pref_i \text{ Eventually } \varphi) \wedge$
 $Bel_i \text{ Henceforth } ((\kappa \wedge Pref_i \text{ Eventually } \varphi) \rightarrow$
 $(\neg \text{After}_{j:\alpha} \perp \wedge \text{After}_{j:\alpha} \varphi \wedge Pref_j \text{ Does}_{j:\alpha} \top))$

Example

1's *dispositional trust* in 2 (a mechanic) to repair his car so that the car will be in order, in the circumstances in which 1 will ask the mechanic to repair his car

$DispTrust(1, 2, rep1c, OK1c, ask12) \stackrel{def}{=}$

$$Poss_1 \text{ Eventually } (ask12 \wedge Pref_1 \text{ Eventually } OK1c) \wedge \\ Bel_1 \text{ Henceforth } ((ask12 \wedge Pref_1 \text{ Eventually } OK1c) \rightarrow \\ (\neg After_{2:rep1c} \perp \wedge After_{2:rep1c} OK1c \wedge Pref_2 \text{ Does}_{2:rep1c} \top))$$

Some properties of dispositional trust

Theorem

- $\vdash \text{DispTrust}(i, j, \alpha, \varphi, \kappa) \leftrightarrow$
(Poss_i Eventually_i ($\kappa \wedge \text{Pref}_i \text{ Eventually}_i \varphi$) \wedge
Bel_i Henceforth_i (($\kappa \wedge \text{Pref}_i \text{ Eventually}_i \varphi$) \rightarrow
(Does_{i:\alpha} $\top \wedge \text{After}_{j:\alpha} \varphi$)))
- $\vdash \text{DispTrust}(i, j, \alpha, \varphi, \kappa) \leftrightarrow \text{Bel}_i \text{DispTrust}(i, j, \alpha, \varphi, \kappa)$

From dispositional to occurrent trust

Theorem

$$\vdash (\mathit{DispTrust}(i, j, \alpha, \varphi, \kappa) \wedge \\ \mathit{Pref}_i \mathit{Eventually} \varphi \wedge \mathit{Bel}_i \kappa) \rightarrow \\ \mathit{OccTrust}(i, j, \alpha, \varphi)$$

Discussion: other definition of trust

- Deutsch's definition (1958): trust involves risk perception (trusting is doing a real bet on the trustee)
 - Uncertain beliefs are needed to capture it
- Genuine trust (Trust game, Baier 1986, Holton 1994): i 's belief that j has goodwill towards him (j is willing to sacrifice his wellbeing for i)
 - It is just a special case of our definition
- Jones's definition (2002): rule belief+conformity belief
 - Trust \neq thinking or foreseeing

From trust to reputation

Reputation: the building blocks

$Rep(I, j, \alpha, \varphi, \kappa)$ = “ j has reputation in group I to ensure φ by doing α in the circumstances κ ”

- Reputation as the collective counterpart of trust
- TRUST = agent i 's individual evaluation (belief) about some properties of agent j that are relevant for a goal of i
- REPUTATION = group I 's collective evaluation about some properties of agent j that are relevant for a (group) goal of I
- to be defined: collective evaluation, group goals

Public facts

For every I such that $|I| > 2$,

$\text{Public}_I \varphi =$ “ φ is public in the group of agents I ” (“ φ is said in I ”)

- “ φ is public in I ” \neq common belief in I that φ
- “ φ is public in I ” does not imply “every agent in I believes φ ”

Logic of public facts: some principles

- 1 $\neg(\text{Public}_I \varphi \wedge \text{Public}_I \neg \varphi)$
- 2 $\text{Public}_I \varphi \rightarrow \text{Public}_J \text{Public}_I \varphi$
if $J \subseteq I$
- 3 $\neg \text{Public}_I \varphi \rightarrow \text{Public}_J \neg \text{Public}_I \varphi$
if $J \subseteq I$
- 4 $\text{Public}_I \varphi \rightarrow \text{Bel}_i \text{Public}_I \varphi$
if $i \in I$
- 5 $\neg \text{Public}_I \varphi \rightarrow \text{Bel}_i \neg \text{Public}_I \varphi$
if $i \in I$
- 6 $\text{Public}_I (\bigwedge_{i \in I} \text{Bel}_i \varphi \rightarrow \varphi)$

Group goals

$\text{GroupPref}_I \varphi$ = “the agents in I want that φ ”

- group goals: broad sense (individual preference aggregation)
- weaker than joint goals and joint intentions [Grosz & Kraus 96]

Some examples of $\text{GroupPref}_I \varphi$ definition:

- \exists -group goal: $\bigvee_{i \in I} \text{Pref}_i \varphi$
- \forall -group goal: $\bigwedge_{i \in I} \text{Pref}_i \varphi$
- Group goal based on majority: $\bigvee_{J \subseteq I, |J| > |\wedge J|} \bigwedge_{i \in J} \text{Pref}_i \varphi$

Formal definition of reputation

$$\mathit{Rep}(I, j, \alpha, \varphi, \kappa) \stackrel{\text{def}}{=}$$

$$\mathit{PotGoal}(I, \varphi) \wedge$$

$$\text{Poss}_I \text{Henceforth}((\kappa \wedge \text{GroupPref}_I \text{Eventually } \varphi) \rightarrow (\neg \text{After}_{j:\alpha} \perp \wedge \text{After}_{j:\alpha} \varphi \wedge \text{Pref}_j \text{Does}_{j:\alpha} \top))$$

- $\text{Poss}_I \varphi \stackrel{\text{def}}{=}} \neg \text{Public}_I \neg \varphi$

- $\text{Poss}_I \varphi$ = “according to the group I , φ is possible”

- $\mathit{PotGoal}(I, \varphi) \stackrel{\text{def}}{=}}$

$$\text{Poss}_I \text{Eventually}(\kappa \wedge \text{GroupPref}_I \text{Eventually } \varphi)$$

Conclusion

Conclusion

- Contribution.
 - formal definitions of occurrent trust and dispositional trust
 - formal definition of reputation and its relation with trust
- Our related works
 - semantics for the different modalities
 - mathematical properties of the logic of belief, preference, action and time (soundness, completeness)
 - trust in information sources and in communication systems
 - from binary trust to graded trust
 - trust in agents vs. trust in roles
 - Implementation of a BDI agent reasoning about trust (see Vercouter's presentation)

THANK YOU!