

Prolegomena for a logic of trust and reputation

Emiliano Lorini

Institut de Recherche en Informatique de Toulouse (IRIT), France

NorMAS 2008, 15th July 2008

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; Demolombe 01; 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

- ‘top-down’ qualitative analysis: reduction of trust to more primitive concepts
 - trust \Rightarrow belief, goal, intention, action
 - analysis in terms of existing 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]

⇒ 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 .

⇒ trust defined from four more primitive concepts

⇒ logic of goal, ability, power and intention?

- 1 Two kinds of trust
- 2 Occurrent trust
- 3 Dispositional trust
- 4 From dispositional trust to reputation
- 5 Related works and conclusion

Two kinds of trust: occurrent trust and dispositional trust

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 believes *j* is going to perform α *here and now*

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

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

Dispositional trust

i believes *henceforth*, *j* will perform α *under some conditions*

$$\begin{aligned}
 \text{DispTrust}(i, j, \alpha, \varphi) &\stackrel{\text{def}}{=} \text{PotGoal}(i, \varphi) \wedge \\
 &\text{Believes}(i, \text{CondCap}(j, \alpha)) \wedge \\
 &\text{Believes}(i, \text{CondPower}(j, \alpha, \varphi)) \wedge \\
 &\text{Believes}(i, \text{CondIntends}(j, \alpha))
 \end{aligned}$$

\Rightarrow predicates to be defined: *PotGoal*, *CondCap*, *CondPower*, *CondIntends*

\Rightarrow relationship between occurrent and dispositional trust:

$$\vdash (\text{DispTrust}(i, j, \alpha, \varphi) \ \& \ \text{conditions}) \rightarrow \text{OccTrust}(i, j, \alpha, \varphi)$$

Occurrent trust

Occurrent 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 $Believes(i, \psi)$, $OccGoal(i, \varphi)$ using modal operators of time, belief, preference, and action

convention:

- Bel_i , etc. = logical operators

Temporal operators

Henceforth φ = “ φ henceforth holds”

Eventually φ = “ φ eventually holds”

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

Belief operators

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

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

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

Goals as preferences about the future

$\text{Pref}_i \varphi$ = “agent i prefers that φ ”

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

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

- add a negative condition: $\text{Pref}_i \text{Eventually } \varphi \wedge \neg \text{Bel}_i \varphi$

Capability in dynamic logic

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

- 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 ” = “ α can be executed by j ”

Power in dynamic logic

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

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

- $\neg \text{After}_{i:\alpha} \perp = i$ can do α
- Does $_{i:\alpha} \top = i$ does α
- logical relationships [Lorini & Demolombe DEON'08]:
 $\text{Does}_{i:\alpha} \varphi \rightarrow \neg \text{After}_{i:\alpha} \neg \varphi$ ($\text{Does}_{i:\alpha} \top \rightarrow \neg \text{After}_{i:\alpha} \perp$)

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

- present-directed intention
- vs. future-directed intention:
 $\text{Pref}_j \text{Eventually} \text{Does}_{j:\alpha} \top$ [Bratman 87, Lorini & Herzig 08]

Formal definition of occurrent trust

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

Dispositional trust

Dispositional trust: components

- definiendum:
 - belief,
 - potential goal,
 - conditional capability,
 - conditional power,
 - conditional intention.
- definiens: formulas of a modal logic of time, belief, preference, and action
- conditional capability, etc. = conditioned occurrent capability etc., prefixed by 'henceforth'

Potential goal

$$PotGoal(i, \varphi) \stackrel{\text{def}}{=} \neg Bel_i \text{ Henceforth } \neg OccGoal(i, \varphi)$$

Conditional capability

$CondCap(j, \alpha) \stackrel{\text{def}}{=} \text{Henceforth}(\kappa_{OccCap(j, \alpha)} \rightarrow OccCap(j, \alpha))$

- $\kappa_{OccCap(j, \alpha)}$ = circumstances under which i expects j has (occurrent) capability to perform α
- **executability precondition of α**

Conditional power

$CondPower(j, \alpha, \varphi) \stackrel{\text{def}}{=} \text{Henceforth}(\kappa_{OccPower(j, \alpha, \varphi)} \rightarrow OccPower(j, \alpha, \varphi))$

- $\kappa_{OccPower(j, \alpha, \varphi)}$ = circumstances under which i expects that j has the power to obtain φ via α
- **effect precondition for α causing φ**

Conditional intention

$CondIntends(j, \alpha) \stackrel{\text{def}}{=}$

$Henceforth (\kappa_{OccIntends(j, \alpha)} \rightarrow OccIntends(j, \alpha))$

- $\kappa_{OccIntends(j, \alpha)}$ = circumstances under which i expects j will intend to perform α
- $CondIntends(j, \alpha)$ = “willingness”
- different forms of willingness [Lorini & Demolombe DEON’08]:
 - when j believes that i wants him to do α , j forms the intention to do α ($\Rightarrow i$ is adoptive toward j)
 - when j believes to be obliged to do α , j forms the intention to do α ($\Rightarrow i$ is obedient).

Formal definition of dispositional trust

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

$$\begin{aligned} & \neg \text{Bel}_i \text{ Henceforth } \neg \text{Pref}_j \text{ Eventually } \varphi \wedge \\ & \text{Bel}_i \text{ Henceforth } (\kappa_{\text{OccCap}}(j, \alpha) \rightarrow \neg \text{After}_{j:\alpha} \perp) \wedge \\ & \text{Bel}_i \text{ Henceforth } (\kappa_{\text{OccPower}}(j, \alpha, \varphi) \rightarrow \text{After}_{j:\alpha} \varphi) \wedge \\ & \text{Bel}_i \text{ Henceforth } (\kappa_{\text{OccIntends}}(j, \alpha) \rightarrow \text{Pref}_j \text{ Does}_{j:\alpha} \top) \end{aligned}$$

- $\kappa_{\text{OccCap}}(j, \alpha)$, $\kappa_{\text{OccPower}}(j, \alpha, \varphi)$, $\kappa_{\text{OccIntends}}(j, \alpha)$ = nonlogical conditions
- should be jointly true when i has occurrent goal that φ
- replace $\text{PotGoal}(i, \varphi)$ by:

$$\neg \text{Bel}_i \text{ Henceforth } \neg (\text{Pref}_j \text{ Eventually } \varphi \wedge \kappa_{\text{OccCap}}(j, \alpha) \wedge \kappa_{\text{OccPower}}(j, \alpha, \varphi) \wedge \kappa_{\text{OccIntends}}(j, \alpha)).$$

From dispositional to occurrent trust

Theorem

Suppose Bel_i and Henceforth are normal modal operators, and suppose Henceforth obeys $\text{Henceforth } \varphi \rightarrow \varphi$. Then

$$\vdash \left(\begin{array}{l} \text{DispTrust}(i, j, \alpha, \varphi) \\ \wedge \text{OccGoal}(i, \varphi) \\ \wedge \text{Believes}(i, \kappa_{\text{OccCap}}(j, \alpha)) \\ \wedge \text{Believes}(i, \kappa_{\text{OccPower}}(j, \alpha, \varphi)) \\ \wedge \text{Believes}(i, \kappa_{\text{OccIntends}}(j, \alpha)) \end{array} \right) \rightarrow \text{OccTrust}(i, j, \alpha, \varphi)$$

Formal definition of reputation

Reputation: the building blocks

$Rep(I, j, \alpha, \varphi)$ = “ j has reputation to do α in order to achieve φ in group I ”

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

Group beliefs and group goals

$\text{GroupBelief}_I \varphi$ = “the group of agents I believes that φ ”
[Tuomela 02; Gaudou et al. 06; Lorini et al. 08]

- “ I has group belief that φ ” = “ φ is publicly established in the group of agents I ”
- group belief \neq common belief
 - group belief does not imply individual belief
- “ I has group belief that φ ” \approx “shared voice in I that φ ”
[Conte & Paolucci 02]

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

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

Formal definition of reputation

$$\begin{aligned}
 \mathit{Rep}(I, j, \alpha, \varphi) &\stackrel{\text{def}}{=} \mathit{PotGoal}(I, \varphi) \wedge \\
 &\quad \text{GroupBelief}_I \mathit{CondCap}(j, \alpha) \wedge \\
 &\quad \text{GroupBelief}_I \mathit{CondPower}(j, \alpha, \varphi) \wedge \\
 &\quad \text{GroupBelief}_I \mathit{CondIntends}(j, \alpha)
 \end{aligned}$$

- $\mathit{PotGoal}(I, \varphi) \stackrel{\text{def}}{=} \neg \text{GroupBelief}_I \text{Henceforth} \neg \text{GroupPref}_I \text{Eventually } \varphi$
- $\mathit{CondCap}(j, \alpha)$, $\mathit{CondPower}(j, \alpha, \varphi)$ and $\mathit{CondIntends}(j, \alpha)$ defined as before

Related works and conclusion

Existing implemented models of reputation: a comparison (1)

Cognitive (Cog): mental states like beliefs are used to define the reputation.

Evaluation group (Grp): a target agent can have different reputations for different groups.

Group goal (GG): the reputation is with respect to a goal of the group.

Action capability (Cap): the definition of reputation considers the agent's capabilities to do a certain action.

Action power (Pow): the definition of reputation considers the power of the agent to achieve the group goal by means of an action.

Intention (Int): the definition of reputation considers the agent's intentions to perform the action or goal for the group.

Existing computational models of reputation: a comparison (2)

<i>Definition</i>	Cog	Grp	GG	Cap	Pow	Int
Our model	yes	yes	yes	yes	yes	action
Repage [Sabater et al. 06]	yes	yes	yes	no	no	goal
FIRE [Huynh et al. 04]	no	no	no	no	no	action
LIAR [Muller & Vercouter 05]	no	no	yes	no	no	no
Regret [Sabater & Sierra 02]	no	yes	no	yes	no	no
e-Bay	no	no	no	no	no	no

Table: Properties of definitions of reputation

Conclusion

- Contribution.
 - formal definitions of occurrent trust, dispositional trust, and reputation
 - evaluation of existing reputation systems w.r.t. the formal definition
- Our related works [DEON'08, KRAMAS'08].
 - mathematical properties of the logic of belief, preference, action and time (soundness, completeness)
 - interactions between the different modalities
 - trust in information sources and communication systems

THANK YOU!