

# Trustworthy Advice

Nardine Osman and Patricia Gutierrez and Carles Sierra <sup>1</sup>

**Abstract.** We propose a novel trust model for assessing the trustworthiness of advice. We say an advice is composed of a plan to execute and a goal to be fulfilled. The expectation of an advice's outcome is calculated by assessing the probabilities of committing to and executing the plan, and the probability of the executed plan fulfilling the intended goal. The probabilities are learned from similar past experiences.

## 1 The CONSUASOR Model

'If you go to Ferran Adria's restaurant you will have the time of your life!' 'If you study daily you will get good marks next semester.' These are examples of advice. Advice is what one relies on to help making decisions. But how can one choose which advice to follow and which to discard? CONSUASOR is a trust model that assesses the trustworthiness of advice.

We say an advice has two main components: a plan to execute and the goal to be achieved. In dynamic logic, an advice may be formalised as:  $[P_\eta]G$ . That is, if  $\eta$  performs plan  $P$ , then goal  $G$  will be achieved. The question then is: how much should  $\alpha$  trust an advice  $[P_\eta]G$  recommended by  $\rho$ ?

CONSUASOR is based on the concept that good advice takes into consideration three issues: (1) the *compliance* of the person being advised with following an advice; (2) the *honour* of the person being advised with respect to executing an advice that he has accepted, and (3) the *causality* describing whether a given plan can cause the intended goal.

We say a trust measure reflects the expectation about an advice's outcome, and we model this expectation as the conditional probability of observing goal  $G$  being achieved given  $\rho$ 's advice  $[P_\eta]G$ , which we define accordingly:

$$p(\text{Observe}(\alpha, G) \mid \text{Commit}(\rho, [P_\eta]G)) = \sum_{P', P'' \in \mathcal{P}} \frac{p(\text{Observe}(\alpha, G) \mid \text{Observe}(\beta, P'')) \cdot p(\text{Observe}(\beta, P'') \mid \text{Commit}(\eta, P'))}{p(\text{Commit}(\eta, P') \mid \text{Commit}(\rho, [P_\eta]G))} \quad (1)$$

Equation 1 defines this probability as a product of the probabilities of compliance (the probability of  $\eta$  committing to  $P'$ , given that  $P$  was recommended by  $\rho$ ), honour (the probability of observing  $\eta$  executing  $P''$ , given that it committed to  $P'$ ), and causality (the probability of observing the goal  $G$  being achieved, given that the plan  $P''$  was executed). An aggregation is then used for considering all possible plans  $\mathcal{P}$  that  $\eta$  may commit to and execute for goal  $G$  to be fulfilled.

Of course, different goals may also be achieved. The probability distribution describing all possible outcomes becomes:

$$\mathbb{P}(\text{Observe}(\alpha, X) \mid \text{Commit}(\rho, [P_\eta]G)) = \left\{ \begin{array}{l} p(\text{Observe}(\alpha, G') \mid \text{Commit}(\rho, [P_\eta]G)), \\ p(\text{Observe}(\alpha, G'') \mid \text{Commit}(\rho, [P_\eta]G)), \\ \dots \end{array} \right\} \quad (2)$$

These probabilities are built by learning from past experiences. We define a single experience  $\mu$  as follows:

$$\langle \text{Commit}(\rho, [P_\eta]G), \text{Commit}(\eta, P'), \text{Observe}(\beta, P''), \text{Observe}(\alpha, G') \rangle$$

**Simplification.** We will use the notation  $\mathbb{P}_C(X \mid [P_\eta]G)$  to describe the probability distribution of compliance ( $\mathbb{P}(\text{Commit}(\eta, X) \mid \text{Commit}(\rho, [P_\eta]G))$ ),  $\mathbb{P}_H(X \mid P_\eta)$  to describe the probability distribution of honour ( $\mathbb{P}(\text{Observe}(\beta, X_\eta) \mid \text{Commit}(\eta, P))$ ), and  $\mathbb{P}_R(X \mid P_\eta)$  to describe the probability distribution of causality ( $\mathbb{P}(\text{Observe}(\alpha, X) \mid \text{Observe}(\beta, P_\eta))$ ). We use the notation  $\mathbb{P}_{C|H|R}(X \mid \cdot)$  to refer to any of those three distributions.

**Initialisation.** At the initial time  $t_I$  when no experiences have been considered yet, we choose the uniform distribution  $\mathbb{F}$  to describe ignorance:  $\mathbb{P}_{C|H|R}^{t_I}(X \mid \cdot) = \mathbb{F}$ . As new experiences are encountered, the probabilities are modified as follows.

**Update.** To update a probability distribution  $\mathbb{P}_{C|H|R}^{t_I}(X \mid \cdot)$  with respect to an experience  $\mu$ , the following is performed:

1. The relevance of the experience  $\mu$  is calculated accordingly:

- The relevance of an experience  $\mu = \langle \text{Commit}(\rho', [P'_\eta]G'), \text{Commit}(\eta, P''), \cdot, \cdot \rangle$  with respect to updating  $\mathbb{P}_C^{t_n}(X \mid [P'_\eta]G)$  is calculated as  $R_\mu(\mathbb{P}_C^{t_n}(X \mid [P'_\eta]G)) = (\zeta_g \cdot \text{Sim}(G', G) + \zeta_p \cdot \text{Emp}(P', P)) / (\zeta_g + \zeta_p)$ , which describes that an experience is considered relevant if the goals are similar (specified by  $\text{Sim}(G', G)$ ) and the past plan empowers the newly recommended plan (specified by  $\text{Emp}(P', P)$ ). The similarity between two goals  $G'$  and  $G$  is measured as their semantic distance in the domain ontology. A plan  $P'$  empowers plan  $P$  if the capabilities required to execute plan  $P$  are implied by the capabilities required for plan  $P'$ . The parameters  $\zeta_g$  and  $\zeta_p$  help specify the weight of each these measure.
- The relevance of an experience  $\mu = \langle \cdot, \cdot, \text{Commit}(\eta, P''), \text{Observe}(\beta, P''), \cdot \rangle$  with respect to updating  $\mathbb{P}_H^{t_n}(X \mid P_\eta)$  is calculated as  $R_\mu(\mathbb{P}_H^{t_n}(X \mid P_\eta)) = \text{Emp}(P', P)$ . Note that goals are no longer relevant in this context.
- The relevance of an experience  $\mu = \langle \cdot, \cdot, \text{Observe}(\beta, P''), \text{Observe}(\alpha, G') \rangle$  with respect to updating  $\mathbb{P}_R^{t_n}(X \mid P_\eta)$  is calculated as  $R_\mu(\mathbb{P}_R^{t_n}(X \mid P_\eta)) = \text{Sim}(P', P)$ . Note that plan similarity is used as opposed to plan empowerment in this context, as peers' capabilities are not relevant for assessing causal relations between plans and goals.

2. The experience  $\mu$  is used to modify the probability of a

<sup>1</sup> IIIA-CSIC, Spain, email: {nardine,patricia,sierra}@iia.csic.es

single expectation accordingly:

$$p_{C|H|R}^{t_n}(X=x|_-) = p_{C|H|R}^{t_{n'}}(X=x|_-) + (1 - p_{C|H|R}^{t_{n'}}(X=x|_-)) \cdot \epsilon \cdot R_\mu(\mathbb{P}_{C|H|R}(X|_-)) \quad (3)$$

where the past probability  $p_{C|H|R}^{t_{n'}}(X=x|_-)$  (calculated at time  $t_{n'} < t_n$ ) is increased by a fraction ( $\epsilon \cdot R_\mu(\mathbb{P}_{C|H|R}(X|_-))$ ) of the total amount that the probability is allowed to increase ( $1 - p_{C|H|R}^{t_{n'}}(X=x|_-)$ ). This fraction is defined by a fixed percentage that is specified via  $\epsilon$ , and it is then tuned further by the relevance of  $\mu$  with respect to  $x$  ( $R_\mu(\mathbb{P}_{C|H|R}(X|_-))$ ). We note that when assessing compliance or honour,  $x$  is chosen based on its semantic distance to what was recommended ( $P$ ) such that this distance is equal to the semantic distance between  $P'$  and  $P''$  of the experience  $\mu$  (that is,  $Sim(P, x) = Sim(P', P'')$ ). In the case of assessing causality,  $x$  is the observed goal  $G'$  of  $\mu$ .

3. The probability distribution is updated w.r.t. the experience  $\mu$  following the minimum relative entropy approach:

$$\mathbb{P}_{C|H|R}^{t_n}(X|_-) = \arg \min_{\mathbb{P}(X|_-)} KL(\mathbb{P}_{C|H|R}^{t_{n'}}(X|_-), \mathbb{P}(X|_-)) \quad (4)$$

such that  $p(X=x|_-) = p_{C|H|R}^{t_n}(X=x|_-)$

where  $KL$  calculates the Kullback-Leibler distance, or the relative entropy, between  $\mathbb{P}_{C|H|R}^{t_{n'}}(X|_-)$  and the argument  $\mathbb{P}(X|_-)$ , and  $p(X=x|_-) = p_{C|H|R}^{t_n}(X=x|_-)$  specifies the constraint that the argument  $\mathbb{P}(X|_-)$  should satisfy. In other words, we look for distributions that satisfy the newly calculated point  $p_{C|H|R}^{t_n}(X=x|_-)$  and are at a minimal distance from the original distribution  $\mathbb{P}_{C|H|R}^{t_{n'}}(X|_-)$ .

**Trust Measure.** After calculating Equation 2, which we refer to as  $\mathbb{P}^{t_n}(X|[P_\eta^\rho]G)$  for simplification, the question now is: How do we calculate a *trust measure* based on this expectation that is expressed as a probability distribution? One proposed approach is to calculate the distance between the distribution  $\mathbb{P}^{t_n}(X|[P_\eta^\rho]G)$  and the distribution representing  $\rho$ 's promised outcome:  $\mathbb{P}_P(X|[P_\eta]G) = \{1, \text{ if } X = G; 0, \text{ otherwise}\}$ . The final trust measure is calculated as:

$$trust^{t_n}(\alpha, \rho, [P_\eta]G) = 1 - emd(\mathbb{P}_P(X|[P_\eta]G), \mathbb{P}^{t_n}(X|[P_\eta^\rho]G)) \quad (5)$$

where  $emd$  measures the earth mover's distance (with the range  $[0, 1]$ ) between two probability distributions.

## 2 Evaluation

We consider an action meronomy  $\mathcal{M}$  and a goal ontology  $\mathcal{O}$  of 10 terms each. A set of plans  $\mathcal{P} \in 2^{\mathcal{M}}$  and goals  $\mathcal{G} \in 2^{\mathcal{O}}$ . A causality function  $f: \mathcal{P} \rightarrow \mathcal{G}$  that describes what goal does each plan achieve. A single user  $\eta$  by the tuple  $\langle G_\eta, c_\eta, h_\eta, d_1, d_2 \rangle$ , where  $G_\eta$  specifies  $\eta$ 's goal,  $c_\eta: \mathcal{P} \rightarrow \mathcal{P}$  and  $h_\eta: \mathcal{P} \rightarrow \mathcal{P}$  describe  $\eta$ 's compliance and honour (when  $P$  is recommended to  $\eta$ ,  $\eta$  commits to  $c_\eta(P)$ ; and when  $\eta$  commits to  $P$ ,  $\eta$  executes  $h_\eta(P)$ ), tuned by distances  $d_1, d_2 \in [0, 1]$  such that  $\forall P \in \mathcal{P}: Sim(c_\eta(P), P) \geq 1 - d_1$  and  $Sim(h_\eta(P), P) \geq 1 - d_2$ . And a set of recommenders, where each recommender  $\rho$  is defined by the tuple  $\langle \{c_{\rho, \eta}^{-1}\}_{\forall \eta}, \{h_{\rho, \eta}^{-1}\}_{\forall \eta}, f_\rho^{-1}, d_3, d_4, d_5 \rangle$  such that:

$$\begin{aligned} \forall P \in \mathcal{P}: & Sim(c_{\rho, \eta}^{-1}(P), P') \geq 1 - d_3 \text{ and } c_\eta(P') = P \\ & Sim(h_{\rho, \eta}^{-1}(P), P') \geq 1 - d_4 \text{ and } h_\eta(P') = P \\ & Sim(f_\rho^{-1}(P), P') \geq 1 - d_5 \text{ and } f(P') = P \end{aligned}$$

where  $d_3, d_4, d_5 \in [0, 1]$  and  $c_{\rho, \eta}^{-1}$ ,  $h_{\rho, \eta}^{-1}$ , and  $f_\rho^{-1}$  describe what  $\rho$  believes the inverse functions of  $c_\eta$ ,  $h_\eta$ , and  $f$  are.

We compare 3 trust strategies: selecting advise randomly, selecting the advice whose adviser is ranked top by eigentrust [1],<sup>2</sup> and selecting the advice ranked top by CONSUASOR. Each experiment runs for 100 timesteps, and in each timestep: (1)  $\eta$  asks for advice for  $G_\eta$ ; (2) *each* recommender  $\rho$  suggests a plan  $P = c_{\rho, \eta}^{-1}(h_{\rho, \eta}^{-1}(f_\rho^{-1}(G)))$ ; (3)  $\eta$  selects an advice following the experiment's trust strategy, and (4) the experience  $\mu = \langle P, c_\eta(P), h_\eta(c_\eta(P)), f(h_\eta(c_\eta(P))) \rangle$  is generated, the success of  $[P_\eta]G$  is calculated as  $Sim(G, f(h_\eta(c_\eta(P))))$ , and the current success of  $\eta$  ( $Succ_\eta^t$ ) is an aggregation of the success of all its previously adopted advices.

Figure 1 shows results with good users ( $d_1 = d_2 = 0$ ) and bad users ( $d_1 = d_2 = 1$ ). For bad users,  $Succ_\eta^t$  remains low because it is not always possible to find plans that achieve  $G_\eta$  and that  $\eta$  is willing to commit to or execute. Cases (a) and (c) have 30 recommenders with medium knowledge ( $\bar{d} = 0.8$ , where  $\bar{d} = (d_1 + d_2 + d_3)/3$ ). Cases (b) and (d) have 5 recommenders ranging from fully knowledgeable ( $\bar{d} = 0.6$ ) to ignorant ( $\bar{d} = 1$ ). Eigentrust obtains high levels of success when there is at least one good recommender, but when recommenders are not very knowledgeable or users are not fully compliant/honorable, its success diminish. CONSUASOR is able to learn which advices are trustworthy and always reaches high levels of success.

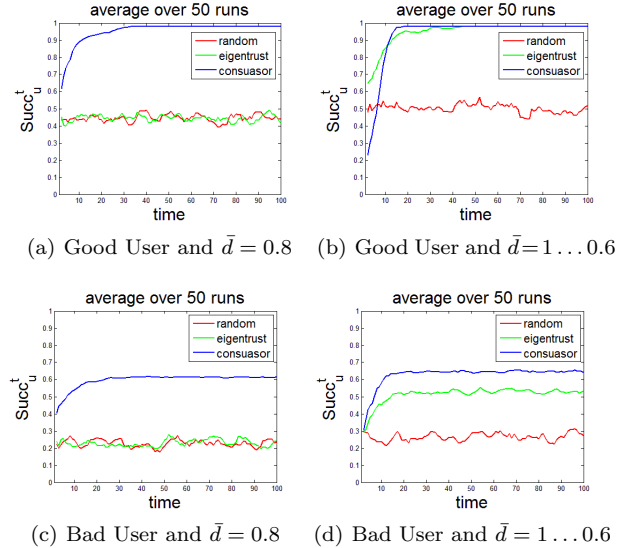


Figure 1. Success over time

## 3 Acknowledgments

This work is supported by the Agreement Technologies project (CONSOLIDER CSD2007-0022, INGENIO 2010) and the PRAISE project (EU FP7 grant number 388770).

## REFERENCES

- [1] Sepandar D. Kamvar, Mario T. Schlosser, and Hector Garcia-Molina, 'The eigentrust algorithm for reputation management in p2p networks', in *Proceedings of WWW-03*, (2003).

<sup>2</sup> In the eigentrust case, we compute the normalized local trust. The notion of transitive trust is not needed since neither users nor recommenders provide advice among themselves.