



TRUSTLET, OPEN RESEARCH ON TRUST METRICS

PAOLO MASSA, KASPER SOUREN, MARTINO SALVETTI, AND DANILO TOMASONI*

Abstract. A trust metric is a technique for predicting how much a user of a social network might trust another user. This is especially beneficial in situations where most users are unknown to each other such as online communities. We believe the recent tumultuous evolution of social networking demands for a collective research effort. With this in mind we created Trustlet.org, a platform consisting of a wiki for open research on trust metrics. The goal of Trustlet is to collect and distribute trust network datasets and trust metrics code as Free Software, in order to facilitate the comparison of different trust metrics algorithms and a more coherent progress in this field. At present we made available some social network datasets and code for some trust metrics. In this paper we describe Trustlet and report a first empirical evaluation of different trust metrics on the Advogato social network dataset.

Key words: trust metrics, social network analysis, wiki, advogato, free software, data acquisition, science commons

1. Introduction. In our current society it is more and more common to interact with strangers, people who are totally unknown to us. This happens for example when receiving an email asking for collaboration or advice from an unknown person, when we rely on reviews written by unknown people on sites such as Amazon.com, and also when browsing random profiles on social networking sites such as Facebook.com or LinkedIn.com. Even more surprising is the fact a huge amount of commercial exchanges happen now between strangers, facilitated by platforms such as Ebay.com. In all systems in which is possible to interact with unknown people, it is important to have tools able to suggest which other users can be trustworthy enough for engaging with.

Trust metrics and reputation systems [10] have precisely this goal and become even more important, for instance, in systems where people are connected in the physical world such as carpooling systems or hospitality exchange networks (i. e. couchsurfing.com), in which users accept to have strangers into their car or their house. In fact, in all the previous examples, the system can give users the possibility of expressing a trust statement, an explicit statement stating “I trust this person in this context” (for example as a pleasant guest in a house or as a reliable seller of items) [10] and then use this information in order to predict trustworthiness of users. Trust becomes in this way one of the building block of the society [5].

While research about trust issues spanned disciplines as diverse as economics, psychology, sociology, anthropology and political science for centuries, it is only recently that the widespread availability of modern communication technologies facilitated empirical research on large social networks, since it is now possible to collect real world datasets and analyze them [10]. As a consequence, recently computer scientists and physicists started contributing to this new research field as well [13, 3].

Moreover we all start relying more and more on these social networking sites [4], for friendship, commerce, work, ... As this field become more and more crucial, in the past few years many trust metrics have been proposed but there is a lack of comparisons and analysis of different trust metrics under the same conditions. As Sierra and Sabater put it in their complete “Review on Computational Trust and Reputation Models” [15]: “Finally, analyzing the models presented in this article we found that there is a complete absence of test-beds and frameworks to evaluate and compare the models under a set of representative and common conditions. This situation is quite confusing, especially for the possible users of these trust and reputation models. It is thus urgent to define a set of test-beds that allow the research community to establish comparisons in a similar way to what happens in other areas (e.g. machine learning)”. Our goal is to fulfill this void and for this reason we set up Trustlet [1], a collaborative wiki in which we aim to aggregate researchers interested in trust and reputation and build together a lively test-bed and community for trust metrics evaluation. A related project is the Agent Reputation and Trust (ART) Testbed [6]. However ART is more focused on evaluating different strategies for interactions in societies in which there is competition and the goal is to perform more successfully than other players, in a specific context. Our take with Trustlet is about evaluating performances of trust metrics in their ability to predict how much a user could trust another user, in every context. For this reason, we want also to support off-line evaluation of different trust metrics on social network datasets. The two testbeds are hence complementary.

In this paper we describe Trustlet, the reason behind its creation and its goals, we report the datasets we have collected and released and the trust metrics we have implemented and we present a first empirical evaluation of different trust metrics on the Advogato dataset.

*FBK/rst, Via Sommarive, 14, Povo (TN)—Italy, {massa, souren, salvetti, tomasoni}@fbk.eu

$\text{trust}(A,B)=\text{reputation}(B)$ for every user A. This global value is sometimes called reputation [10]. Currently most trust metrics used in web communities are global, mainly because they are simpler to understand for the users and faster to run on central servers since they have to be executed just once for the entire community. However we think that soon users will start asking for systems that take into account their own peculiar points of view and hence local trust metrics, possibly to be run in a decentralized fashion on their own devices.

While research on trust metrics is quite recent, there have been some proposals for trust metrics. We briefly review some of them for later mention in the evaluation presented in Section 4, although our goal is not to provide a complete review of previously proposed trust metrics here.

Ebay web site shows the average of the feedbacks received by a certain user in her profile page. This can be considered as a simple global trust metric, which predicts, as trust of A in B, the average of all the trust statements received by B [11].

In more advanced trust metrics, trust can be extended beyond direct connections. The original Advogato trust metric [8] is global, and uses network flow to let trust flow from a “seed” of 4 users, who are declared trustworthy a priori, towards the rest of the network. The network flow is first calculated on the network of trust statements whose value is Master (highest value) to find who classifies as Master. Then the Journeyer edges are added to this network and the network flow is calculated again to find users who classify as Journeyer. Finally the users with Apprentice status are found by calculating the flow on all but the Observer edges. The untrusted Observer status is given if no trust flow reached a node. By replacing the 4 seed users for an individual user A, Advogato can also be used as a local trust metrics predicting trust from the point of view of A.

The problem of ranking of web pages in the results of a search engine query can be regarded under a trust perspective. A link from page A to page B can be seen as a trust statement from A to B (in this case, the nodes of the trust network are not people but Web pages). This is the intuition behind the algorithm PageRank [2] powering the search engine Google. Trust is propagated with a mechanism resembling a random walk over the trust network.

Moletrust [11] is a local trust metric. Users are ordered based on their distance from the source user, and only trust edges that go from distance n to distance $n + 1$ are regarded. The trust value of users at distance n only depends on the already calculated trust values at distance $n - 1$. The scores that are lower than a specific threshold value are discarded, and the trust score is the average of the incoming trust statements weighted over the trust scores of the nodes at distance $n - 1$. It is possible to control the locality by setting the trust propagation horizon, i.e. the maximum distance to which trust can be propagated.

Golbeck proposed a metric, TidalTrust [7], that is similar to Moletrust. It also works in a breadth first search fashion, but the maximum depth depends on the length of the first path found from the source to the destination. Another local trust metric is Ziegler’s AppleSeed [16], based on spreading activation models, a concept from cognitive psychology.

3. Datasets and Trust Metrics Evaluation. Research on trust metrics started a long time ago, but is somehow still in its infancy. The first trust metric could probably be ascribed to the philosopher John Locke who in 1680 wrote: “Probability then being to supply the defect of our knowledge, the grounds of it are these two following: First, the conformity of anything with our own knowledge, observation and experience. Secondly, the testimony of others, vouching their observation and experience. In the testimony of others is to be considered: (1) The number. (2) The integrity. (3) The skill of the witnesses. (4) The design of the author, where it is a testimony out of a book cited. (5) The consistency of the parts and circumstances of the relation. (6) Contrary testimonies” [9]. This quotation can give an idea of how many different models for representing and exploiting trust have been suggested over the centuries. However of course John Locke in 1680 didn’t have the technological means for empirically evaluating his “trust metric”. Even collecting the required data about social relationships and opinions was very hard in old times. The first contributions in analyzing real social networks can be tracked down to the foundational work of Jacob Moreno [12] (see Figure 2.1) and since then many sociologists, economists and anthropologists have researched on social networks and trust. But the advent of the information age has made it possible to collect, represent, analyze and even build networks way beyond what is possible with pen and paper. Computer scientists and physicists have hence become interested in social networks, now that both huge amounts of data have become available and computing power has advanced considerably [13, 3].

At Trustlet (<http://www.trustlet.org>) we have started a wiki to collect information about research on trust and trust metrics. Our goal is to attract a community of people with interest in trust metrics. The wiki is totally open: anonymous edits are allowed and anybody can register and create an account. We have chosen to

use the Creative Commons Attribution license so that work can easily (and legally) be reused elsewhere. Our effort shares the vision of the Science Commons project¹ which tries to remove unnecessary legal and technical barriers to scientific collaboration and innovation and to foster open access to data. We have also started a repository of the software we create for our analysis, written in Python and available as Free Software under the GNU General Public License (GPL)² so that other researchers can replicate our experiments and reuse our code.

We believe the lack of generally available datasets is inhibiting scientific progress. It's harder to test a hypothesis if it has been tested on a dataset that is not easily available. The other alternative is testing the hypothesis on synthesized datasets, which are hardly representative of real-world situations. Prior to the proliferation of digital networks data had to be acquired by running face-to-face surveys, which could take years to collect data of a mere couple of hundreds of nodes. The proliferation and popularity of on-line social networks [4] has facilitated acquiring data, and the implementation of standards like XFN and common APIs like OpenSocial opens up new possibilities for research [10]. A more widespread availability and controlled release of datasets would surely benefit research and this is one of the goals behind the creation of Trustlet.

We think it is important that research on trust metrics follows an empirical approach and it should be based on actual real-world data. Our goal with Trustlet is to collect as many datasets as possible in one single place and release them in standard formats under a reasonable license allowing redistribution and, at least, usage in a research context. At present, as part of our effort with Trustlet, we collected and released datasets derived from `advogato.org`, `people.squeakfoundation.org`, `robots.net`, `kaitiaki.org.nz` and `epinions.com`³.

We describe in detail the Advogato dataset since our experiments (presented in Section 4) are run on it. Advogato.org is an online community site dedicated to Free Software development, launched in November 1999. It was created by Raph Levien, who also used Advogato as a research testbed for testing his own attack-resistant trust metric, the Advogato trust metric [8]. On Advogato users can certify each other at several levels: Observer, Apprentice, Journeyer or Master. The Advogato trust metric uses this information in order to assign a global certification level to every user. The goal is to be attack-resistant, i. e. to reduce the impact of attackers [8]. Precise rules for giving out trust statements are specified on the Advogato site. Masters are supposed to be principal authors of an "important" Free Software project, excellent programmers who work full time on Free Software, Journeyers contribute significantly, but not necessarily full-time, Apprentices contribute in some way, but are still acquiring the skills needed to make more significant contributions. Observers are users without trust certification, and this is also the default. It is also the level a user certifies another user at to remove a previously expressed trust certification. Notwithstanding the suggestions, users are free to express totally subjective certifications on other users.

For the purpose of this paper we consider these certifications as trust statements [11]. $T(A,B)$ denotes the certification expressed by user A about user B and we map the textual labels Observer, Apprentice, Journeyer and Master in the range $[0,1]$, respectively in the values 0.4, 0.6, 0.8 and 1.0. This choice is arbitrary and considers all the certifications as positive judgments, except for Observer which is used for expressing less-than-sufficient levels. For example, we model the fact raph certified federico as Journeyer as $T(\text{raph}, \text{federico})=0.8$.

The Advogato social network has a peculiarly interesting characteristic: it is almost the only example of a real-world, directed, weighted, large social network. However, besides the leading work of Levien reported in his unfinished PhD thesis [8], we are just aware of another paper using the Advogato dataset which is focused on providing a trust mechanism for mobile devices [14].

There are other web communities using the same software powering Advogato.org and they have the same trust levels and certifications system: `robots.net`, `people.squeakfoundation.org`, `kaitiaki.org.nz`. We collected daily snapshots of all these datasets and made them available on Trustlet but we haven't used them for our analysis in this paper, mainly because they are much smaller than the Advogato dataset. Details about the characteristics of the analyzed Advogato trust network dataset are presented in Section 4.

The other datasets we released on Trustlet are derived from Epinions.com, a website where users can leave reviews about products and maintain a list of users they trust and distrust based on the reviews they wrote [11].

On Trustlet, we released these datasets but our aim is to collectively make it a repository of all the possible datasets useful for research on trust issues. For this reason, we also keep on the Trustlet wiki a list of datasets we are considering for collection and a list of datasets released elsewhere.

¹Science Commons <http://sciencecommons.org>

²GNU General Public License <http://www.gnu.org/licenses/gpl.html>

³See <http://www.trustlet.org/wiki/Trustnetworkdatasets>

Moreover, besides aiming at releasing datasets in a coherent format, we also released on Trustlet.org the Python code we wrote for the trust metrics analyzed in Section 4 under a Free Software license so that code can be reused and inspected.

4. Initial Research Outcomes. In the previous sections we highlighted the reasons for creating Trustlet and the way we aim it can develop into a collaborative environment for the research of trust metrics. As a first example of what we envision Trustlet will be able to bring to research on trust metrics, we report our first investigation and empirical findings.

We chose to start studying the Advogato social network because of its almost unique characteristic: trust statements (certifications) are weighted and this makes it a very peculiar dataset for researching trust metrics, in fact, most other networks just exhibit a binary relationship (either trust is present or not) and the evaluation on trust metrics performances is much less insightful.

In this paper we report experiments performed on the Advogato dataset we downloaded from the web site on May 12th 2008. This dataset is available at Trustlet.org, along with datasets downloaded in other days as well. The Advogato dataset under analysis is a directed, weighted graph with 7294 nodes and 52981 trust relations. There are 17489 Master judgments, 21977 for Journeyer, 8817 for Apprentice and 4698 for Observers. The dataset is comprised of 1 large connected component, comprising 70.5% of the nodes; the second largest component contains 7 nodes. The mean in- and out-degree (number of incoming and outgoing edges per user) is 7.26. The mean shortest path length is 3.75. The average cluster coefficient [13] is 0.116. The percentage of trust statements which are reciprocated (when there is a trust statement from A to B, there is also a trust statement from B to A) is 33%.

While a large part of research on social networks focuses on exploring the intrinsic characteristics of the network [13, 6, 3], on Trustlet we are interested in covering an area that received much less attention, analysis of trust metrics. We have compared several trust metrics through leave-one-out, a common technique in machine learning. The process is as follows: one trust edge (e.g. from node A to node B) is taken out of the graph and then the trust metric is used to predict the trust value A should place in B, i. e. the value on the missing edge. We repeat this step for all edges to obtain a prediction graph, in which some edges can contain an undefined trust value (where the trust metric could not predict the value). The real and the predicted values are then compared in order to derive several evaluation measures: the coverage, which is a measure of the edges that were predictable, the fraction of correctly predicted edges, the mean absolute error (MAE) and the root mean squared error (RMSE). Surely there are other ways of evaluating trust metrics: for instance, it can be argued that an important task for trust metrics is to suggest to a user which other still unknown users are more trustworthy, such as suggesting a user worth following on a social bookmarking site such as del.icio.us or on a music community such as Last.fm. In this case the evaluation could just concentrate on the top 10 trustworthy users. But in this first work we considered only leave-one-out as evaluation technique.

4.1. Evaluation of trust metrics on all trust edges. Table 4.1 reports our evaluation results of different trust metrics on the Advogato dataset. It is a computation of different evaluation measures on every edge of the social network. The reported measures are: Mean Absolute Error (MAE), Root Mean Squared Error (RMSE), fraction of wrong predictions, and coverage. We now describe the compared trust metrics. As already mentioned we released the code and we plan to implement more trust metrics and release them and run more evaluations.

The compared trust metrics are some trivial ones used as baselines such as Random, which predicts simply a random trust score between the 4 possible ones (1.0, 0.8, 0.6, 0.4), or the metrics starting with “Always” which always return the corresponding value as predicted trust score, for example AlwaysApprentice returns 0.6 for every prediction. Other simple trust metrics are OutA which, in predicting the trust user A could have in user B, simply does the average of the trust statements outgoing from user A, and OutB which averages over the trust statements outgoing from user B. These simple trust metrics are considered in order to understand how much and in which cases complex algorithms are useful.

The other trust metrics were already explained in Section2, here we just report the parameters we used in running them. Ebay refers to the trust metric that, in predicting the trust user A could have in user B, simply does the average of the trust statements incoming in user B, i. e. the average of what all the users think about user B. MoletrustX refers to Moletrust applied with a trust propagation horizon of value X. The values returned by PageRank as predicted trust follow a powerlaw distribution, there are few large PageRank scores and many tiny ones. So we decided to rescaled the results simply by sorting them and linearly mapping them in

TABLE 4.1
Evaluation of trust metrics on all trust edges

	Fraction wrong predictions	MAE	RMSE	Coverage
Random	0.737	0.223	0.284	1.00
AlwaysMaster	0.670	0.203	0.274	1.00
AlwaysJourneyer	0.585	0.135	0.185	1.00
AlwaysApprentice	0.834	0.233	0.270	1.00
AlwaysObserver	0.911	0.397	0.438	1.00
Ebay	0.350	0.086	0.156	0.98
OutA	0.486	0.106	0.158	0.98
OutB	0.543	0.139	0.205	0.92
Moletrust2	0.366	0.090	0.160	0.80
Moletrust3	0.376	0.091	0.161	0.93
Moletrust4	0.377	0.092	0.161	0.95
PageRank	0.501	0.124	0.191	1.00
AdvogatoLocal	0.550	0.186	0.273	1.00
AdvogatoGlobal	0.595	0.199	0.280	1.00

the range $[0.4, 1]$, after this we rounded the predicted trust scores. Our implementation of Advogato is based on Pymmetry, whose code is released on Trustlet as well. AdvogatoGlobal refers to the Advogato trust metric run considering as seeds the original founders of Advogato community, namely the users “raph”, “federico”, “miguel” and “alan”. This is the version that is running on the Advogato web site for inferring global certifications for all the users. This version is global because it predicts a trust level for user B which is the same for every user. AdvogatoLocal refers to the local version of Advogato trust metric. For example, when predicting the trust user A should place in user B, the trust flow starts from the single seed “user A”. This version is local because it produces personalized trust predictions which depends on the current source user and can be different for different users. AdvogatoLocal was run on a subset (8%) of all the edges since the current implementation is very slow. In fact, the leave-one-out technique requires the network be different for every evaluation and it has to be rebuilt from scratch for every single trust edge prediction making the entire process very slow.

Since some trust metrics such as Moletrust and PageRank produce trust score predictions in a continuous interval while others just the 4 discrete trust scores, we decided to apply a rounding to the closest possible certification value before the predicted trust scores are compared with the real values so that for example a predicted trust score of 0.746 becomes 0.8 (Journeyer).

The results of the evaluation are reported in Table 4.1. We start by commenting the column “fraction of wrong predictions”. Our baseline is the trust metric named “Random” which produces an incorrect predicted trust score 74% of the times. The best one is Ebay with an error as small as 35% followed by Moletrust2 (36.57%), Moletrust3 (37.60%) and Moletrust4 (37.71%). Increasing the trust propagation horizon in Moletrust allows to increase the coverage but also increases the error. The reason is that users who are nearby in the trust network (distance 2) are better predictors than users further away in the social network (for example, users at distance 4). This is consistent with experiments on other social networks [11].

Note that Moletrust is a local trust metric that only uses information available “near” the source node so it can be run on small devices such as mobiles which only need to fetch information from the (few) trust users and possibly the users trusted by them. This behaviour is tunable through setting the trust propagation horizon to specific values. On the other hand, Ebay, being a global trust metric, must aggregate the entire trust network, which can be costly both in term of bandwidth, memory and computation power. So a local trust metric tends to require less information for producing recommendations which might be a desirable features in some situations.

The AlwaysX metrics depend on the distributions of certifications and are mainly informative of the data distribution.

The fraction of wrong predictions of Advogato (both local and global) is high compared to Ebay and Moletrust. The reason is that Advogato was not designed for predicting an accurate trust value for a specific

pair of users (the trust A should place in B) but to increase attack-resistance [8], i. e. being able to exclude malicious users, while accepting as many valid accounts as possible. A side effect is that it limits the amount of granted global certifications and assigns a large number of Observer certificates. In the case of AdvogatoGlobal, 45% of the predicted global certifications are marked as Observer which obviously has an impact on the leave-one-out evaluation. Different trust metrics might have different goals, that require different evaluation techniques. We could have tuned different parameters of Advogato for making it perform differently, however our intention was to evaluate the original trust metric in the task of predicting trust scores so we decide to run Advogato with the original parameters. Note also that the local version of Advogato is more accurate than the global version. The last metric shown in Table 2.1 is PageRank [2]: the fraction of correct predictions is not too high but again the real intention of PageRank is to rank web pages and not to predict the correct value of assigned trust.

An alternative evaluation measure is the Mean Absolute Error (MAE). The MAE is computed by averaging the difference in absolute value between the real and the predicted trust statement on an edge. There is no need to round values to the closest certification value because MAE computes a meaningful value for continuous values. However, in order to fairly compare trust metrics that return real values and trust metrics that return discrete values, we choose to perform anyway the rounding to the closest possible certification value before computing MAE.

The second column of Table 4.1 reports the MAE for the evaluated trust metrics. The baseline is given by the Random trust metric which incurs in a MAE of 0.2230. These results are the worst besides the trivial trust metrics that always predict the most infrequent certification values. Predicting always Journeyer (0.8) incurs in a small MAE because this value is frequent and central in the distribution of assigned trust scores. Ebay is the trust metric with the best performance, with a MAE of 0.0855. And it is again followed by Moletrust that in a similar way is more accurate with smaller trust propagation horizons than with larger ones.

A variant of MAE is Root Mean Squared Error (RMSE). RMSE is the root mean of the average of the squared differences. This evaluation measure tends to emphasize large errors, which favor trust metrics that remain within a small band of error and don't have many outlying predictions that might undermine the confidence of the user in the system. For example, it penalizes a prediction as Observer when the trust score the source user would have assigned was Master, or vice versa. The baseline trust metric Random has an RMSE of 0.2839. Again Ebay is the best metric with an RMSE of 0.1563 and all the other performances exhibit a pattern similar to the one exposed for the other evaluation measures. However there is one unexpected result: the trivial trust metric OutA is the second best, close to Ebay. Remind that, when asked a prediction for the trust user A should place in user B, OutA simply returns the average of the trust statements going out of A, i. e. the average of how user A judged other users. This trust metric is just a trivial one that was used for comparison purposes. The good performance of OutA in this case is related to the distribution of the data in this particular social setting. The Observer certification has special semantics: it is the default value attributed to a user unless the Advogato trust metric gives a user a higher global certification. So there is little point in certifying other users as Observer. In fact, the FAQ specifies that Observer is "the level to which you would certify someone to remove an existing trust certification". Observer certifications are mainly used when a user changes its mind about another user and wants to downgrade her previously expressed certification as much as possible. This is also our reason for mapping it to 0.4, a less than sufficient level. As a consequence of the special semantics of observer certifications, they are infrequently used. In fact only 638 users used the Observer certification at least once while, for instance, 2938 users used the Master certification at least once. Trust metrics like Ebay and Moletrust work doing averages of the trust edges of the network (from a global point of view for Ebay and only considering the ones expressed by trusted users for Moletrust) and, since the number of Observer edges is very small compared with the number of Master, Journeyer and Apprentice edges, these predicted average tend to be close to higher values of trust. This means that when predicting an Observer edge (0.4) they tend to incur in a large error. This large error is emphasized by the squaring of the RMSE formula. On the other hand, using the average of the outgoing trust edges (like OutA does) happens to be a successful technique for not incurring in large errors when predicting Observer edges. The reason is that a user who used Observer edges tended to use it many times so the average of its outgoing edge certifications is a value that is closer to 0.4 and hence it incurs in lower errors on these critical edges and, as a consequence, in smaller RMSE. This effect can also be clearly seen when different trust metrics are restricted to predict only Observer edges and evaluated only on them. In this case (not shown in Tables), OutA gets the correct value for trust (Observer) 42% of times, while for instance, Ebay only 2.7% of times and Moletrust2 4%. The fact the trivial trust metric

TABLE 4.2
Evaluation of trust metrics on trust edges going into controversial users

	Fraction wrong predictions	MAE	RMSE	Coverage
Random	0.799	0.266	0.325	1.00
AlwaysMaster	0.462	0.186	0.302	1.00
AlwaysJourneyer	0.801	0.202	0.238	1.00
AlwaysApprentice	0.943	0.296	0.320	1.00
AlwaysObserver	0.794	0.414	0.477	1.00
Ebay	0.778	0.197	0.240	0.98
OutA	0.614	0.147	0.199	0.98
OutB	0.724	0.215	0.280	0.92
Moletrust2	0.743	0.195	0.243	0.80
Moletrust3	0.746	0.194	0.241	0.93
Moletrust4	0.746	0.195	0.242	0.95
PageRank	0.564	0.186	0.275	1.00
AdvogatoLocal	0.518	0.215	0.324	1.00
AdvogatoGlobal	0.508	0.216	0.326	1.00

OutA exhibits a so small RMSE supports the intuition that evaluating which conditions a certain trust metric is more suited for than another one is not a trivial task. Generally knowledge about the domain and the patterns of social interaction is useful, if not required, for a proper selection of a trust metric for a specific application and context.

The last column of Table 4.1 reports the coverage of the different trust metrics on the Advogato dataset. For some trust edges, a trust metric might not be able to generate a prediction and the coverage refers to the number of edges that are predictable. The experiment shows that the coverage is always very high. Since local trust metrics use less information (only trust statements of trusted users) their coverage is smaller than the coverage of global trust metrics. Anyway, differently from other social networks [11], it is very high. The Advogato trust network is very dense, so there are many different paths from a user to another user. Even very local trust metrics such as Moletrust2, that only use information from users at distance 2 from the source user, are able to cover and predict almost all the edges.

4.2. Evaluation of trust metrics on controversial users. As a second step in the analysis we devoted our attention to controversial users [11]. Controversial users are users which are judged in very diverse ways by the members of a community. In the context of Advogato, they can be defined as users who received many certifications as Master and many as Apprentice or Observer: the community does not have a single way of perceiving and judging them. The intuition here is that a global average can be very effective when all the users of the community agree that “raph” is a Master, but there can be situations in which something more tailored and user specific is needed, especially when there isn’t a subjective judgment that is shared by all the members of the community.

With this in mind we define controversiality level of an Advogato user as the standard deviation in certifications received by that user, similarly to previous studies [11]. The higher the standard deviation, the more controversial the user is. A user with controversiality level as 0 is not controversial at all since all the other users certify her with the same value. The certification level is not very meaningful when the number of received certifications for an user is small (for example 3); for this reason in the following we are going to report measures on users who received at least 10 or 20 incoming certificates, and for which the standard deviation in received certifications really represents the fact the community does not have a single way of perceiving these popular users.

In Table 4.2 we report the evaluation of the performances of the same trust metrics of Table 4.1 but evaluated only on trust edges going to Advogato users with at least 10 incoming edges and controversiality level of 0.2. In this way we reduce the number of edges considered in the evaluation from 52, 981 to 2, 030, which is still a significant number of edges to evaluate trust metrics on. Figure 4.1 graphically reports the number of edges going into users (who received at least 20 certifications) with at least a certain controversiality level for

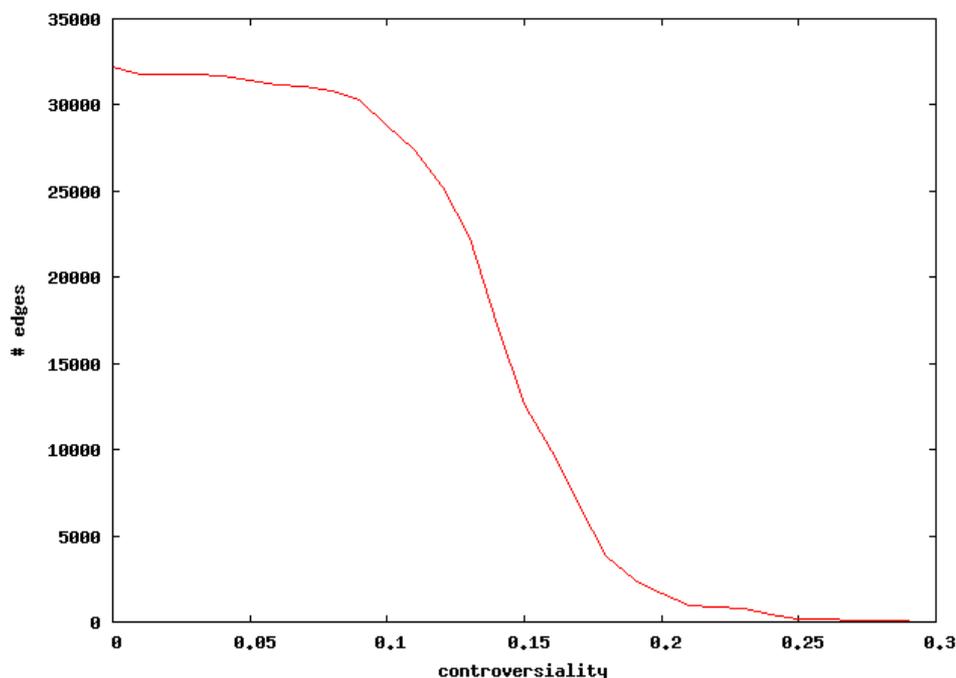


FIG. 4.1. Number of edges per controversy level

all controversy levels from 0 to 0.3. As intuitive, increasing the controversy level of users decreases the number of edges going into users with at least that controversy level.

Figure 4.2 on the other hand shows how at higher controversy levels the percentage of polarized trust scores increases: certifications as Master and Observer becomes more frequent. This means that predicting trust edges going into controversial users is in theory more difficult, since it is important to predict the correct trust score which is not close to an average score. Both figures confirm intuition and are informative of the distribution of trust scores.

Going back to the evaluation measures presented in Table 4.1, we start by commenting the evaluation measures on AlwaysMaster (second row of Table 4.2) because it presents some peculiarities. AlwaysMaster predicts the correct trust value 53.84% (100% 46.16%) of times and, according to the evaluation measure “fraction of correctly predicted trust statements”, seems a good trust metric, actually the best one. However the same trust metric, AlwaysMaster, is one of the less precise when RMSE is considered. A similar pattern can be observed for AdvogatoGlobal. In fact, since in general there is at least one flow of trust with Master certificates going to these controversial users, AdvogatoGlobal tends to predict almost always Master as trust value and since almost half of the edges going into controversial users are of type Master, AdvogatoGlobal often predicts the correct one.

The results presented in Table 4.1 suggest that the same trust metric might seem accurate or inaccurate depending on the choice of the evaluation measure. This fact once more highlights how evaluating trust metrics on real world datasets is a complicated task and a comparison of same trust metrics on many different datasets according to different evaluation methods would be highly beneficial for understanding the situations in which one trust metric is more appropriate and useful than another. We already previously explained why OutA is able to have a so small RMSE, the smallest one on users with controversial level of 0.2: based on how Observer certifications are used in the system, OutA is the only metric that is able to avoid large errors when predicting the Observer edges, which are a relevant percentage since the evaluated users are controversial.

Arriving at a comparison between a global trust metric such as Ebay and a local trust metric such as Moletrust, we were expecting the latter to be significantly more accurate than the first one on controversial users. While on the Epinions dataset, this is what was observed [11], the same is not true here since the two trust metrics incur in very similar performances.

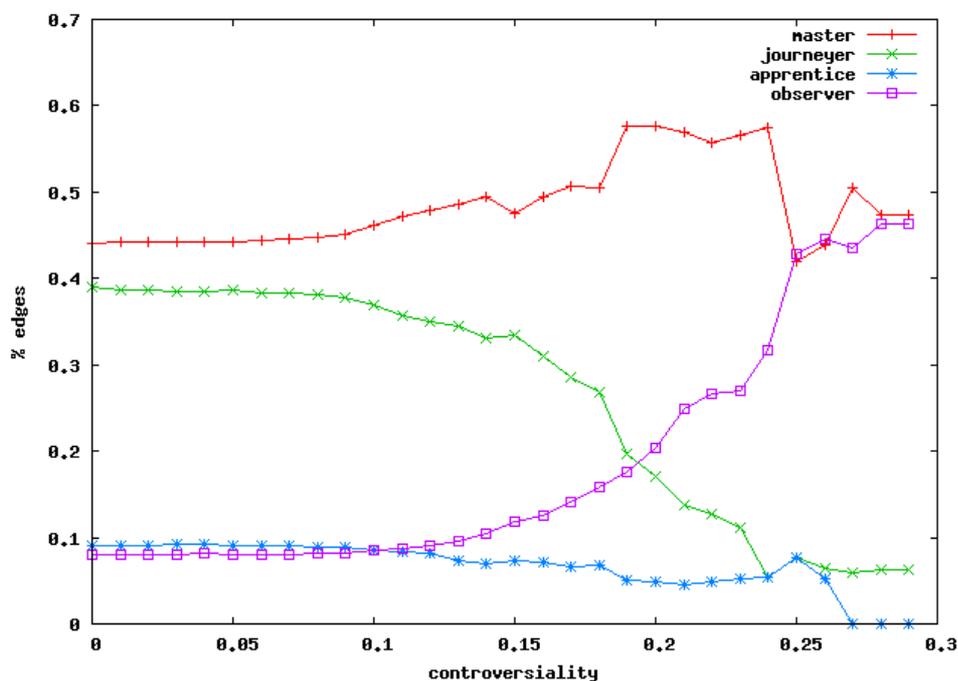


FIG. 4.2. Percentage of edges for each type per controversiality level

Figure 4.3 graphically presents the performances (measured by RMSE) of some selected trust metrics on users with increasing controversiality levels and at least 20 incoming edges. It can be observed that the local trust metrics MoletrustX starts to perform better than Ebay and other metrics when the controversiality levels is larger than 0.25. However the difference is not that large as expected.

The reason for this similarity of performances between Ebay and Moletrust2 is partly that in Epinions, the trust values were binary (either trust or distrust) and it was easier to discriminate. Another reason seems to be that on Advogato the user base is not divided in cliques of users such that users of one clique trust each other and distrust users of other cliques. In fact Advogato users are somehow similar and feel part of one single large community. It is future work to analyze if on a social network with a much higher polarization of opinions (such as for example on essembly.com, a political site, in which users tend to express strong feeling for or against other people based on their political views) the performances of local trust metrics are significantly better than global ones. The study on the Advogato trust network dataset presented in this paper does not allow arguing that local trust metrics and in general complex trust metrics are needed in order to outperform simpler trust metrics. Another future work is exploring different evaluation procedures which might be more informative of the real performances of different trust metrics.

5. Conclusions. In this paper we have presented Trustlet [1], an open environment for research on trust metrics. We have claimed that the rapid development of social networking sites [4] asks for a shared effort in collecting datasets and distributing code of algorithms so that comparisons of different research proposals is easier, replicable and more coherent.

As an initial investigation we have reported our comparison of different trust metrics on the Advogato dataset. The results are partly contradictory and this suggests there is need to run systematically evaluations of different algorithms against a large number of different datasets. As future works we are looking into extending our analysis to more datasets also from different social scenarios, for example the networks of relationships (coediting, talk) among Wikipedia users.

Our goal is to make Trustlet an environment which facilitates this collaborative effort. We believe research on these topics is very needed in a time in which our relationships are starting to move more and more into the “virtual” world and our society and life is affected significantly from the predictions and suggestions produced by many different algorithms.

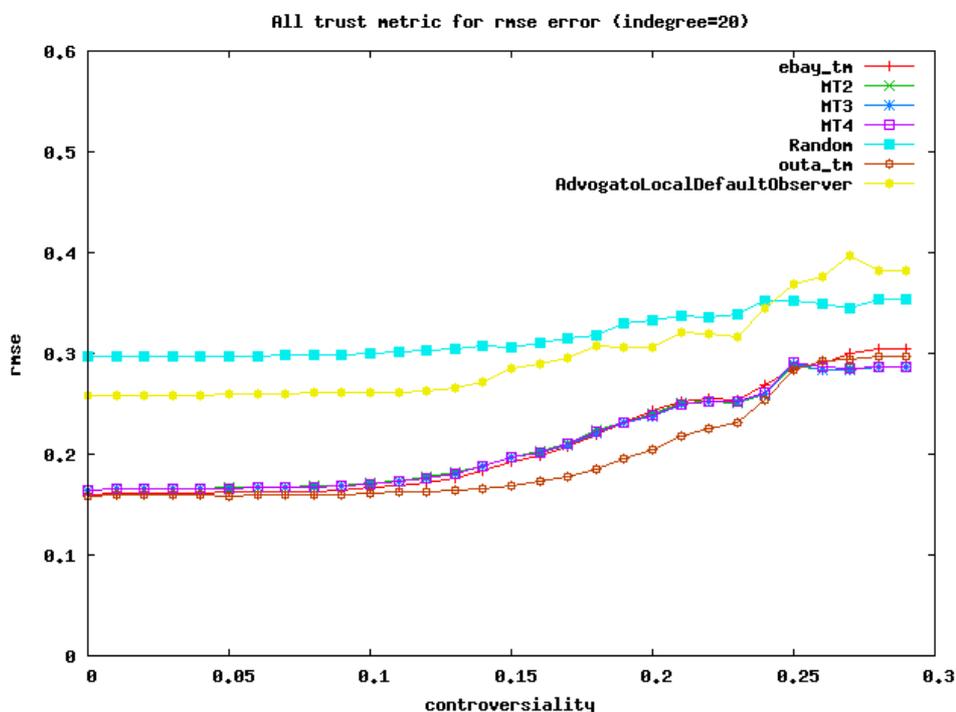


FIG. 4.3. RMSE for some trust metrics for different controversiality levels

REFERENCES

- [1] *Trustlet, collaborative wiki for trust research*. <http://www.trustlet.org>
- [2] D. AUSTIN, *How google finds your needle in the web's haystack*. retrieved on 2008-02-02. <http://www.ams.org/featurecolumn/archive/pagerank.html> 2006.
- [3] A.-L. BARABASI, *Linked: The New Science of Networks*, Perseus, Cambridge, MA, 2002.
- [4] D. M. BOYD AND N. B. ELLISON, *Social network sites: Definition, history, and scholarship*, *Journal of Computer-Mediated Communication*, 13 (2007), p. art. 11.
- [5] F. FUKUYAMA, *Trust: the Social Virtues and the Creation of Prosperity*, Free Press Paperbacks, 1995.
- [6] K. K. FULLAM, T. KLOS, G. MULLER, J. SABATER-MIR, K. S. BARBER, AND L. VERCOUTER, *A specification of the agent reputation and trust (art) testbed*, in *Proceedings of 4th AAMAS Conference*, Utrecht, 2005.
- [7] J. GOLBECK, *Computing and Applying Trust in Web-based Social Networks*, PhD thesis, University of Maryland, 2005.
- [8] R. LEVIEN, *Attack resistant trust metrics*. Ongoing PhD thesis. <http://www.levien.com/thesis/compact.pdf>
- [9] J. LOCKE, *An Essay concerning Human Understanding*, Harvester Press, Sussex, 1680.
- [10] P. MASSA, *Trust in E-Services: Technologies, Practices and Challenges*, Idea Group, Inc, 2006, ch. A survey of trust use and modeling in current real systems.
- [11] P. MASSA AND P. AVESANI, *Trust metrics on controversial user: Balancing between tyranny of the majority and echo chambers*, *International Journal on Semantic Web and Information Systems*, 3 (2007), p. 39-64.
- [12] J. MORENO, *Who Shall Survive? Foundations of Sociometry, Group Psychotherapy and Sociodrama*, Beacon House, New York, 1953.
- [13] M. E. J. NEWMAN, *The structure and function of complex networks*, *SIAM Review*, (2003), pp. 167-256.
- [14] D. QUERCIA, S. HAILES, AND L. CAPRA, *Lightweight distributed trust propagation*, in *Proceedings of the 7th IEEE International Conference on Data Mining*, 2007.
- [15] J. SABATER AND C. SIERRA, *Review on computational trust and reputation models*, *Artificial Intelligence Review*, 24 (2005), pp. 33-60.
- [16] C. ZIEGLER, *Towards Decentralized Recommender Systems*, PhD thesis, Albert-Ludwigs-Universitaet Freiburg, 2005.

Edited by: Dominik Flejter, Tomasz Kaczmarek, Marek Kowalkiewicz

Received: February 9th, 2008

Accepted: March 19th, 2008

Extended version received: August 8th, 2008