

VU Research Portal

Network metrics for assessing the quality of entity resolution between multiple datasets

Idrissou, O.A.K.; van Harmelen, Frank; van den Besselaar, P.A.A.

published in

Knowledge Engineering and Knowledge Management
2018

DOI (link to publisher)

[10.1007/978-3-030-03667-6_10](https://doi.org/10.1007/978-3-030-03667-6_10)

document version

Publisher's PDF, also known as Version of record

document license

Article 25fa Dutch Copyright Act

[Link to publication in VU Research Portal](#)

citation for published version (APA)

Idrissou, O. A. K., van Harmelen, F., & van den Besselaar, P. A. A. (2018). Network metrics for assessing the quality of entity resolution between multiple datasets. In A. Napoli, C. Ghidini, Y. Toussaint, & C. Faron Zucker (Eds.), *Knowledge Engineering and Knowledge Management: 21st International Conference, EKAW 2018, Nancy, France, November 12-16, Proceedings* (pp. 147-162). (Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics); Vol. 11313). Springer Nature Switzerland AG. https://doi.org/10.1007/978-3-030-03667-6_10

General rights

Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the public portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the public portal ?

Take down policy

If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.

E-mail address:

vuresearchportal.ub@vu.nl



Network Metrics for Assessing the Quality of Entity Resolution Between Multiple Datasets

Al Koudous Idrissou^{1,2(✉)}, Frank van Harmelen¹, and Peter van den Besselaar²

¹ Department of Computer Science, Vrije Universiteit Amsterdam, Amsterdam, Netherlands

{o.a.k.idrissou, frank.van.harmelen, p.a.a.vanden.besselaar}@vu.nl

² Department of Organization Sciences, Vrije Universiteit Amsterdam, Amsterdam, Netherlands

Abstract. Matching entities between datasets is a crucial step for combining multiple datasets on the semantic web. A rich literature exists on different approaches to this entity resolution problem. However, much less work has been done on how to *assess* the quality of such entity links once they have been generated. Evaluation methods for link quality are typically limited to either comparison with a *ground truth dataset* (which is often not available), *manual work* (which is cumbersome and prone to error), or *crowd sourcing* (which is not always feasible, especially if expert knowledge is required). Furthermore, the problem of link evaluation is greatly exacerbated for links between more than two datasets, because the number of possible links grows rapidly with the number of datasets. In this paper, we propose a method to estimate the quality of entity links between multiple datasets. We exploit the fact that the links between entities from multiple datasets form a network, and we show how simple metrics on this network can reliably predict their quality. We verify our results in a large experimental study using six datasets from the domain of science, technology and innovation studies, for which we created a gold standard. This gold standard, available online, is an additional contribution of this paper. In addition, we evaluate our metric on a recently published gold standard to confirm our findings.

Keywords: Entity resolution · Data integration · Network metrics

1 Introduction

Matching entities between datasets (known as entity resolution) is a crucial step for the use of multiple datasets on the semantic web. There exists a fair amount of entity resolution tools for *generating* links between pairs of resources: AGDIS-TIS [15], LIMES [12] Linkage Query Writer [7, 8], SILK [16], etc. However, much fewer methods exist for *validating* the links produced by these methods. Currently, only three validation options are available for such validation: (1) *ground truth*, which is often not available; (2) *manual work*, which is a cumbersome

task prone to error; (3) *crowd sourcing*, which is not always feasible especially if specialist knowledge is required. Furthermore, the problem of link evaluation is greatly exacerbated for entity resolution between more than two datasets, because the number of possible links grows rapidly with the number of datasets. Therefore, it is important to investigate ***the accurate automated evaluation of discovered links***. Any answer to this question should generalise beyond the setting of just two datasets, and be applicable to the general setting of links between multiple datasets. In such a multi-dataset scenario, linked resources cluster in small groups that we call *Identity Link Networks (ILNs)*. The goal of this paper is not to propose any new method for entity resolution but instead to provide a method to estimate the quality of an identity link network, and consequently validate a set of discovered links. To do so, ***we hypothesize that the structure of an identity link network correlates with its quality***. We test our hypothesis in two experiments where we show that the proposed metrics indeed reliably estimates the quality of an identity network. We also test our hypothesis on recently published experimental data from ESWC 2018 (see Sect. 8). Here too, the results confirm that our quality metric reliably predicts human assessment of entity links.

In summary, our contributions is a method that estimates the quality of an identity network. It is tested against human judgement in three large experiments and correctly classifies large amount of ILNs available online.¹

This paper begins with a short motivation in Sect. 2. Section 3 discusses the related work and Sect. 4 describes the proposed metric. In Sect. 5 we describe the datasets involved in our experiments. Sections 6, 7 and 8 describe our three experiments, and Sect. 9 concludes.

2 Identity Link Networks

We assume the well known setting of a real-world entity that has one or more digital representations in multiple datasets. The task of entity resolution is to discover which entity (or entities) in each dataset denotes the same real world entity. An Identity Link Network (ILN) is a network of links between entities from a number of datasets that are found by one or more entity resolution algorithms to represent the same real world entity. An ILN can be derived directly from entity resolution results (Sects. 6 and 7), or it may be generated by sophisticated clustering methods as in our experiment in Sect. 8. In this work we do *not* propose any new entity resolution algorithm. Instead, we propose a method to automatically *evaluate* discovered links, particularly when they involve more than two datasets. Unfortunately, gold standards in initiatives such as OAEI do not go beyond two datasets.

Figure 1 shows two examples of such ILNs that have been generated by an entity resolution algorithm between entities from six datasets taken from the field of Science, Technology and Innovation studies (STI) (more details in Sect. 5). Figure 1a shows the ILN for the real world entity University of Trier, Fig. 1b shows

¹ <https://github.com/alkoudouss/Identity-Link-Network-Metric>.

the same for the National Chung Cheng University. *In this paper, we hypothesise that the structure of these ILNs is a reliable indicator for the correctness of the links in the network they form.*

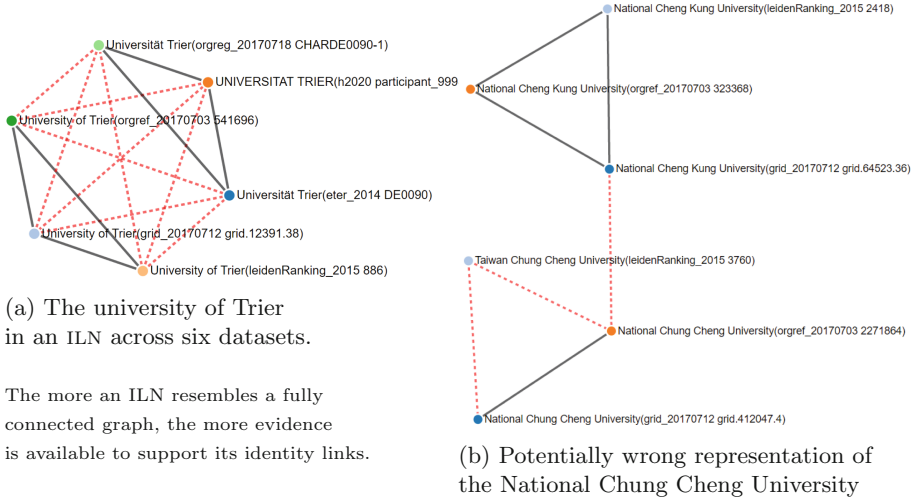


Fig. 1. Two real life examples of Identity Link Networks (ILNs); dotted lines indicate links with a low confidence.

3 Related Work

We briefly discuss a number of related areas from the literature, and indicate how our work differs from these in aim and scope.

Schema Matching. Much work in the literature focuses on ontology matching, especially schema matching [5]. Some rely on concept distance or an extended version of it [3, 10, 17]. Some rely on alignment similarities [4], others relies on formal logical conflicts between ontologies to detect and possibly repair mappings at a schema-level [9]. The current paper does not aim to match ontologies, nor does it critically rely on using ontological or schema information. We only assume the existence of external entity resolution algorithms for suggesting links between entities. Such algorithms may or may not exploit ontological information, but this does not affect our central hypothesis.

Information Gain. The work in [14] also uses network structure to evaluate link quality, but in a very different way. The main intuition there is that an individual link in an ILN is more reliable when it leads to a greater information gain. The paper does not consider the structure of the ILN as a whole, as we do in this paper.

Entity Clustering. Part of the literature also uses clustering of the digital representations of the same real world entity in one or multiple sources. While their data sources are mainly unstructured [1, 2], our interest lies in clusters

derived from the mappings of entities exclusively across knowledge-bases. In addition, they also do not consider the structure of the ILN as a whole. Another part of the literature specifically focuses on clustering algorithms. The FAMER [13] framework for example provides and compares seven different link-based entity clustering approaches. The aim of our work is different from all of these. Whereas these works use clustering algorithms to *construct* entity resolutions, we show how a cluster-based metric can be used to *assess* the quality of a network of entity links, irrespective of how these links were generated.

Network Metrics. The work by Guéret et al. [6] is one of the few papers to our knowledge that uses network metrics to assess the quality of links. The key point that separates this work from ours is that it uses *local* network features, i.e. only the direct neighbours of a single node, while we employ *global* network features. [11] also addresses the same challenge. It evaluates a given cluster G by comparing it to a reference cluster R based on the number of splits and merges required to go from G to R. Our proposed metric does not need such a reference cluster, and is hence more easily applicable.

4 Network Properties and Quality of a Link-Network

Figure 2 illustrates a set of six simple network topologies over the same number of nodes. Our proposed metric is based on the intuition that multiple links provide corroborating evidence for each other, suggesting that in the case of an ILN, the ideal topology is a **fully connected** network. It illustrates a total agreement between all resources (not the case for any other topology), and it does not require any intermediate resource to establish an identity-link between two resources (again, not the case for any other topology). Hence, intuitively, the amount of redundancy between paths in an ILN is an indicator for the quality of the links in the ILN. We will capture these and similar intuitions using three different global graph features over ILNs: *Bridge*, *Diameter* and *Closure*.

We will now first define and explain the rationale behind each metric, then normalise each metric to values² between 0 and 1, and finally average the sum of all metrics to obtain the metric which we will use for estimating the quality of the ILN.

Bridge Metric. A bridge (also known as an isthmus or a cut-edge) in a graph is an edge whose removal increases the number of connected components of the graph, or equivalently, an edge that does not belong to any cycle. The intuition for this measure is that a bridge in an ILN suggests a potentially problematic link which is not corroborated by any other links. As a

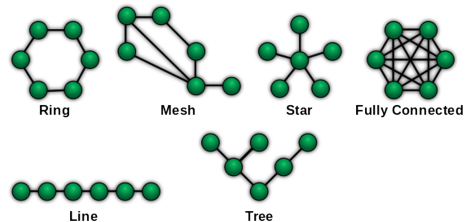


Fig. 2. Example of network topologies. Source: https://en.wikipedia.org/wiki/Network_topology

² The metric value indicates the negative impact of one or more missing links in an ILN.

graph with n nodes contains at most $n - 1$ bridges (e.g. in a **Line** network), the bridge value is normalised as $n_b = \frac{B}{n-1}$, where B is the number of bridges. An ideal link network would have no bridge ($n_b = 0$). As n_b is sensitive to the total number of nodes in the graph (it decreases for large graphs, even when the number of bridges is constant), we “soften” the value of n_b with a sigmoid function: $n'_b = \max(n_b, \text{sigmoid}(B))$, where the function $\text{sigmoid}(x) = \frac{x}{|x|+1.6}$ helps stabilising the impact of the size of the graph by providing a minimal value for n'_b . The value 1.6 is a hyper-parameter that has been determined empirically.

Diameter Metric. The diameter D of a graph with n nodes is the maximum number of edges (distance) in a shortest path between any pair of vertices (i.e. the longest shortest path). In an ideal scenario, if three resources A, B and C are representations of the same real world object, there would be no need for an intermediate resource for confirming the identity of any of the resource in the network. In a fully connected graph of n nodes, the diameter $D = 1$. The longest diameter is observed in a **Line** network structure, with $D = n - 1$ for a line network of n nodes. To scale to the $[0,1]$ interval, the diameter is normalised as $n_d = \frac{D-1}{(n-1)-1}$. Like the bridge, because the diameter is also sensitive to the number of nodes, the normalised diameter is calculated as $n'_d = \max(n_d, \text{sigmoid}(D - 1))$.

Closure Metric. In a connected graph of n nodes, the closure is the ratio of the number of arcs A in the graph over the total number of possible arcs $\frac{1}{2}n(n - 1)$. In a complete graph, this ratio has value 1. Hence, to evaluate how far the observed graph is from the ideal (complete) one, we normalise the closure metric as $n_c = 1 - \frac{A}{\frac{1}{2}n(n-1)}$. The minimum number of connections is $n - 1$, as observed in **Line** and **Star** network structures.

Estimated Quality Metric. All of these metrics capture the same intuition: the more an ILN resembles a fully connected graph, the higher the quality of the links in the ILN. Of course, these three metrics are not independent: $n_c = 0$ or $n'_d = 0$ implies $n'_b = 0$. However, using only n_c or n'_d would be too uninformative since the converse of the implication does not hold. Table 1 shows that each of n_c , n'_d and n'_b capture different (though related) amounts of redundancy in the ILN and that each metric by itself fails to properly discriminate between the seven ILNs depicted in Fig. 2. For example, n_c and n'_c treat a *Tree*, *Star* and *Line* as qualitatively equal but disagree on whether a *Full Mesh* is as good as a *Ring*. Consequently, to compute an overall estimated quality e_Q of an identity link network, we combine the three separate metrics by taking their average, and invert them so that the value 1 indicates the highest quality: (We apply e_Q to ILNs of size ≥ 3 as it is the smallest network where redundancy can be observed.)

$$e_Q = 1 - \frac{n'_b + n'_d + n_c}{3}.$$

Discrete Intervals. The e_Q metric scores all ILNs on a continuous value in the $[0,1]$ interval. To automatically discriminate potentially good networks from bad ones, we divide this interval into three segments: ILNs with values $0.9 \leq e_Q \leq 1$

will be rated as **good**, with values $0.75 < e_Q < 0.9$ as **undecided**, and with values $0 \leq e_Q \leq 0.75$ as **bad**. These boundaries are empirically determined, and can be adjusted depending on the use-case. The specific values of these boundaries does not affect the essence of our hypothesis.

Hypothesis. We can now state our hypothesis more formally: “*The e_Q intervals defined above are predictive of the quality of the links in an entity link network between multiple datasets*”.

Example. By way of illustration, Table 1 gives the value of our e_Q metric for the six networks from Fig. 2, and shows that the metric does indeed capture redundancy in a network.

In the following sections, we will test this hypothesis against human evaluation on hundreds of ILNs containing thousands of links in three experiments using between three to six datasets.

Table 1. Metrics values for each of the topologies from Fig. 2.

Link-Network Quality Estimation				
ILN	Bridge	Diameter	Closure	Est. Quality
Ring	$B = 0 \quad n_b = 0.00$	$D = 3 \quad n_d = 0.56$	$C = 0.40 \quad n_c = 0.60$	$e_Q = 0.61$
Mesh	$B = 1 \quad n_b = 0.38$	$D = 3 \quad n_d = 0.56$	$C = 0.47 \quad n_c = 0.53$	$e_Q = 0.51$
Star	$B = 5 \quad n_b = 1.00$	$D = 2 \quad n_d = 0.38$	$C = 0.33 \quad n_c = 0.67$	$e_Q = 0.32$
Full Mesh	$B = 0 \quad n_b = 0.00$	$D = 3 \quad n_d = 0.00$	$C = 1.00 \quad n_c = 0.00$	$e_Q = 1.00$
Line	$B = 5 \quad n_b = 1.00$	$D = 1 \quad n_d = 1.00$	$C = 0.33 \quad n_c = 0.67$	$e_Q = 0.11$
Tree	$B = 5 \quad n_b = 1.00$	$D = 4 \quad n_d = 0.38$	$C = 0.33 \quad n_c = 0.67$	$e_Q = 0.34$

5 Datasets

We considered using datasets and gold standards from the OAEI³ initiative, but none of these go beyond links between two datasets. We therefore created our own gold standard on realistic datasets taken from the domain of social science, more specifically from the field of Science, Technology and Innovation studies. We consider this to be an important contribution of this paper. All datasets and our gold standard are available online at the locations given in later paragraphs.

Entities of interest to the STI domain of study are (among others) universities and other research-related organisations, such as R&D companies and funding agencies. Our six datasets are widely used in the field, and describe organisations and their properties such as name, location, type, size and other features.⁴

³ <http://oaei.ontologymatching.org/>.

⁴ The information provided here about the datasets was collected in January 2018. The datasets themselves are of earlier dates: Grid: 2017.07.12; Orgref: 2017.07.03; OpenAire: 2018.08.16; OrgReg: 2017.07.18; Eter: 2014; Leiden Ranking 2015: 2017.6.16; and Cordis-H2020: 2016.12.22. All these datasets are available on the RISIS platform at <http://datasets.risis.eu/>.

Grid⁵ describes 80248 organisations across 221 countries using 12308 relationships. All organisations are assigned an address, while 96% of them have an organisation type, and only 78% have geographic coordinates.

OrgRef⁶ collates data about the most important worldwide academic and research organisations (31000) from two main sources: Wikipedia and ISNI.

The Leiden Ranking dataset⁷ offers scientific performance indicators of more than 900 major universities. These universities are only included when they are above the threshold of 1000 fractionally counted Web of Science indexed core publications. This explains its coverage across only 54 worldwide countries.

Eter⁸ is a database on European Higher Education Institutions that not only includes research universities, but also colleges and a large number of specialized schools. The dataset covered 35 countries in 2015.

OrgReg⁹ is based on Eter but adds to the about 2700 HE institutions some 500 public research organizations and university hospitals. Collected between 2000 and 2016, its organisations are distributed across 36 countries.

The European Organisations' Projects H2020 database¹⁰ documents the Horizon 2020 participating organisations.

6 e_Q Put to the Test

We test our hypothesis on a real life case study that revolves around the six datasets described in Sect. 5, with as goal to investigate the coverage of OrgReg (coverage analysis of datasets is a typical question asked by social scientists before including a dataset in their studies). This is done by comparing the entities in OrgReg to those in the other five datasets (Fig. 3).

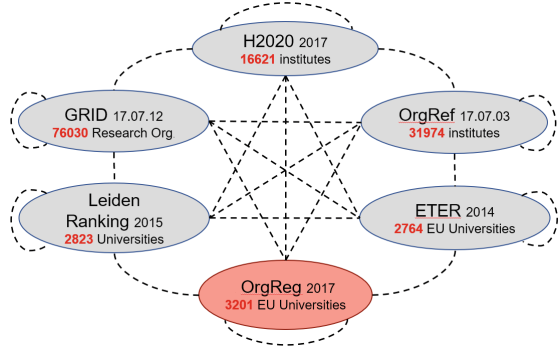


Fig. 3. Disambiguating OrgReg. To evaluate e_Q , all possible links are evaluated. So, the lack of one or more links is considered a potential evidence for suggesting the corresponding entities being different.

⁵ <https://www.grid.ac>.

⁶ <http://www.orgref.org>.

⁷ <http://www.leidenranking.com/>.

⁸ <https://www.eter-project.com/>.

⁹ <http://risis.eu/orgreg/>.

¹⁰ <http://www.gaeu.com/sv/item/horizon-2020>.

6.1 Experiment Design

Organizations are linked across or within datasets using an approximate string matching on their names with minimal similarity threshold 0.8. Based on this, we generate links between each pair of datasets, resulting in 21 sets of links (including linking a dataset to itself in order to detect duplicate entities in the dataset). We then take the union of all 21 sets of links, resulting in a collection of ILN’s of varying size (see Fig. 4).

Now that we have constructed a large collection of multi-dataset ILNs, we will compute the e_Q value for all of them. Then, the machine-predicted good/bad categories (using e_Q) will be checked against the ground truth by a non-domain expert (the first author of this paper) and further verified by a domain expert (the third author). This ground truth is available online.¹¹

Notice that we have deliberately used a very weak entity resolution algorithm in this experiment (approximate string matching). This produces links of both very high and rather low quality, providing a genuine test for our e_Q metric to distinguish between them.

6.2 Results of First Evaluation

Ideally, we would find only ILNs of size 6 if each OrgReg entity were linked with one and only one entity in each of the five other datasets. With less than 100% coverage of OrgReg, we also expect to find ILNs of size < 6 . Figure 4 shows that we also find a substantial number of ILNs of size > 6 . This is due to (a) duplicates occurring in a single dataset, resulting in links in the ILN between two items from the same dataset, and (b) an imperfect matching algorithm (in our case approximate name matching), resulting in incorrect links in the ILN.

Due to the high number of ILNs generated¹², we evaluate only the 846 ILNs of size 5 to 10, with the following frequencies: 391 (size 5), 224 (6), 96 (7), 66 (8), 45 (9) and 24 (10). We predict a ‘good’ or ‘bad’ score based on the e_Q interval values for each of the 846 ILNs, and then compare the scores against those of a human expert, resulting in F_1 scores. In red, Fig. 4 displays the F_1 value for each ILN size. Overall, our e_Q metric resulted in high F_1 values ($0.806 \leq F_1 \leq 0.933$). We also pitched our e_Q metric against a Majority Class Classifier. Table 2 shows that our e_Q metric outperforms the Classifier on F_1 measure, Accuracy (ACC) and Negative Predicted Value (NPC) for ILNs of all sizes.

All of these findings show the very strong predictive power of our e_Q metric for the quality of ILNs when compared to human judgement.

¹¹ <https://github.com/alkoudouss/Identity-Link-Network-Metric>.

¹² On a 6th Gen Intel®Core™i7 notebook with 8 GB RAM, it takes about 1:40 min to automatically evaluate all 4398 clusters of size three and above (see Fig. 4).

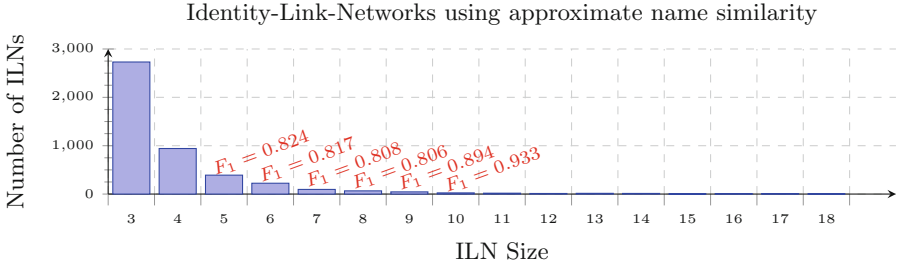


Fig. 4. Overview of the generated Identity Link Networks.

6.3 Results of Second Evaluation

For a further evaluation by a Dutch domain expert from the field of STI (the third author of this paper), we selected 148 ILNs (ranging from size 3 to 10 as depicted in Table 2) in which at least one entity is located in the Netherlands. The expert deviated from the first evaluation in only 12 out of 148 cases. Although the changes slightly affect the ground truth for each ILN size, the F_1 values computed here are even higher ($0.848 \leq F_1 \leq 1$) as compared to the previous experiment. This shows that the non-expert nature of the first human judgement was not detrimental to our results.¹³ This second experiment confirms our finding in the first experiment that e_Q is a reliable predictor of ILN quality.

Table 2. Network-metric (e_Q) results compared to the MCC baseline using non expert Ground Truth (left), and Expert sampled Ground Truth (right).

Majority Class Classifier (Baseline) vs Network Metric (e_Q)								
<i>MajorityClassClassifier</i>								
<i>NetworkMetrics</i>								
$GT_P =$ Ground Truth Positive $GT_N =$ Ground Truth Negative								
Size	$GT_P GT_N$	F_1	ACC	NPV	$GT_P GT_N$	F_1	ACC	NPV
3					56 8	<u>0.933</u> 0.931	<u>0.875</u> 0.875	— 0.5
4					19 5	<u>0.884</u> 0.878	<u>0.792</u> 0.792	— 0.5
5	272 119	<u>0.821</u> 0.824	<u>0.696</u> 0.747	— 0.598	14 1	<u>0.966</u> 0.929	<u>0.933</u> 0.867	— 0
6	139 85	<u>0.766</u> 0.817	<u>0.621</u> 0.768	— 0.709	14 5	<u>0.848</u> 0.848	<u>0.737</u> 0.737	— —
7	50 56	<u>0.685</u> 0.808	<u>0.521</u> 0.792	— 0.810	10 2	<u>0.909</u> 1.0	<u>0.833</u> 1.0	— 1.0
8	35 31	<u>0.693</u> 0.806	<u>0.530</u> 0.803	— 0.765	4 0	<u>1.0</u> 1.0	<u>1.0</u> 1.0	— —
9	21 24	— 0.894	<u>0.533</u> 0.889	<u>0.533</u> 1	8 1	<u>0.941</u> 1.0	<u>0.889</u> 1.0	— 1.0
10	8 16	— 0.933	<u>0.667</u> 0.958	<u>0.667</u> 0.941	1 0	<u>1.0</u> 1.0	<u>1.0</u> 1.0	— —

¹³ However, the very imbalanced character of the ground truth makes it hard to always outperform the baseline as illustrated in Table 2.

6.4 Analysis

Both of the evaluations of e_Q above resulted in very high F_1 average values of 0.847 and 0.961 respectively. Furthermore, e_Q outperformed a majority-class classifier in the first experiment (not in the second because of the highly imbalanced distribution). All this supports our hypothesis that our e_Q measure is strongly predictive of the quality of the links between the entities in an Identity Link Network.

7 e_Q Estimations in Noisy Settings

The previous experiment created links between entities using a rather weak entity resolution heuristic. This was an interesting setting because such weak matching strategies are a fact of daily life on the semantic web (and in data integration in general). In the next experiment, we will use e_Q to evaluate ILN's that have been constructed using a more sophisticated matching heuristic, where we can control the amount of incorrect links in the ILNs. We will see that also in this case, e_Q is strongly predictive of human judged link quality.

The stronger matching heuristic that we use in this second experiment combines organisation names with the geo-location of the organisation. The experiment is run over Eter, Grid and OrgReg as they are the only datasets at our disposal that contain such geo-coordinates for organisations. To test the performance of the e_Q metric at various levels of noise, we implement three sub-experiments where noise (the number of false positive links) is introduced by decreasing the name similarity threshold from 0.8 (experiment 1) to 0.7 and by increasing the geographic proximity distance threshold as described in the next sub-section.

7.1 Experiment Design

This subsection describes in three phases how the experiment is conducted.

Phase-1: Create Links. The first phase links organizations across the three datasets whenever they are located within a radius of 50 m, 500 m and 2 km. This creates nine sets of links (three for each radius).

Phase-2: Refine Links. Each set of links is then refined by applying an approximate name comparison over the linked resources with a threshold of 0.7.

By now, we have **geo-only** (without name comparison) and **geo+names** sets of links, organised in three subgroups (50 m, 500 m and 2 km) each.

Phase-3: Combine Links. To generate the final ILNs, the sets of links within each subgroup are combined using the union operator. The goal of this is to compare, within a specified distance, ILNs that were generated without name matching to those generated with name matching.

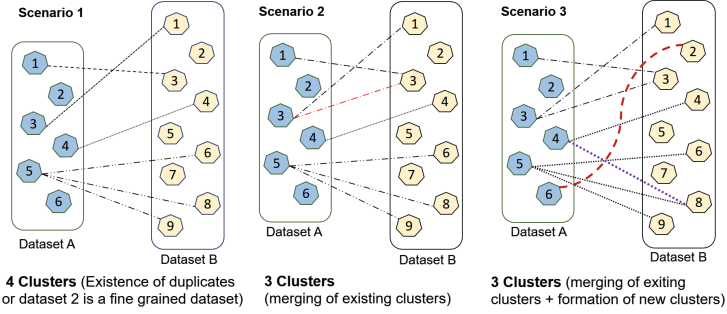


Fig. 5. Decrease/Increase of ILNs

7.2 Strict vs. Liberal Clustering

To understand how link-networks are formed as we increase the geo-similarity distance, Fig. 5 illustrates how ILNs may evolve as we move from strict constraints (scenario 1) to liberal constraints (scenario 3). First, in **scenario 1**, four ILNs are derived from the six links: $c_1 = \{\{a_1\}, \{b_3\}\}$, $c_2 = \{\{a_3\}, \{b_1\}\}$, $c_3 = \{\{a_4\}, \{b_4\}\}$ and $c_4 = \{\{a_5\}, \{b_6, b_8, b_9\}\}$. Then, the new link between a_3 and b_3 in **scenario 2** forces c_1 and c_2 to *merge*. We now have a total of three ILNs: $c_1 = \{\{a_1, a_3\}, \{b_1, b_3\}\}$, $c_3 = \{\{a_4\}, \{b_4\}\}$ and $c_4 = \{\{a_5\}, \{b_6, b_8, b_9\}\}$. Finally, in **scenario 3**, two new links appear. The first link between a_4 and b_8 causes the merging of c_3 and c_4 while the second link connecting a_6 to b_2 causes the creation of a new ILN. Thereby, the total number of ILNs remains 3. These scenarios show that, as the ILN constraints become more liberal, the number of links discovered increases while the number of ILNs may increase, remain equal, or even decrease. In other words, when the matching conditions become liberal or less strict, two types of event may happen: (1) formation of new ILNs and/or (2) merging of ILNs. Table 3, shows that, in experiment 2, phenomenon (1) overtakes (2), which explains the increase in the number of ILNs as the near-by distance increases.

7.3 Result and Analysis

Overall, as illustrated in Table 3, the number of ILNs generated in this experiment increases with the increase of the geo-similarity radius. Within a radius of 50 m, a total of 230 ILNs are generated based on geo-distance only. This number reached 841 ILNs at a 2 km radius. After performing name matching, many links are pruned. Depending on the matching radius, the number of ILNs then varies from 36 to 371.

Due to manpower limitations we restrict our evaluation efforts to networks of size 3. These ILNs cover 86% of the overall ILNs within 50 m radius and 92% within 500 m and 2k radius. Table 4 shows the results of pitching our e_Q metric against the human evaluation of the ILNs under both the geo-only and the geo+names conditions.

Table 3. Link-network overview.

Statistics on ILNs of size > 2						
	50 meters		500 meters		2 kilometres	
Size	geo-only	geo+names	geo-only	geo+names	geo-only	geo+names
≥ 3	230	36	738	168	841	371

As an example, the values $F_1 = 0.803$ and $F_1 = 0.912$ detail the machine quality judgements versus human evaluations of the networks generated within 2km radius under respectively geo-only and geo+names conditions.¹⁴

Table 4. Automated flagging versus human evaluation.

	50 meters		500 meters		2 kilometres	
Size	geo-only	geo+names	geo-only	geo+names	geo-only	geo+names
= 3	92	31	249	155	198	342
Machine statistics on ILN's of size 3						
Machine	$M_{good}: 45$ $M_{maybe}: 0$ $M_{bad}: 47$	$M_{good}: 19$ $M_{maybe}: 12$ $M_{bad}: 0$	$M_{good}: 115$ $M_{maybe}: 0$ $M_{bad}: 134$	$M_{good}: 127$ $M_{maybe}: 0$ $M_{bad}: 28$	$M_{good}: 81$ $M_{maybe}: 0$ $M_{bad}: 117$	$M_{good}: 279$ $M_{maybe}: 0$ $M_{bad}: 63$
Human evaluation on ILN's of size 3						
Human	$H_{good}: 31$ $H_{maybe}: 4$ $H_{bad}: 57$	$H_{good}: 27$ $H_{maybe}: 1$ $H_{bad}: 3$	$H_{good}: 64$ $H_{maybe}: 7$ $H_{bad}: 176$	$H_{good}: 148$ $H_{maybe}: 1$ $H_{bad}: 6$	$H_{good}: 61$ $H_{maybe}: 3$ $H_{bad}: 134$	$H_{good}: 322$ $H_{maybe}: 8$ $H_{bad}: 12$
F_1 measures						
	$F_1 = 0.693$	$F_1 = 0.826$	$F_1 = 0.682$	$F_1 = 0.909$	$F_1 = 0.803$	$F_1 = 0.912$

Analysis. In this experiment, we test the behaviour of the proposed e_Q metric in both noisy (*proximity only*) and noise-less (*proximity plus name*) scenarios. The proposed e_Q metric is in general able to exclude poor networks in noisy environments and to include good networks in noise-less environments. In addition, on the one hand, the relatively low F_1 measures displayed in Table 5 in noisy scenarios, highlight that for the data at hand, proximity alone is not a good enough criterion for identity. On the other hand, the relatively high F_1 measures in noise-less scenarios is an indication of stability and consistency that is in line with results outlined in experiment 1.

The results depicted in Table 5 show an uneven distribution of the candidate-sets. In a relatively balanced candidate-set scenario, our approach works well as can be seen in the first experiment and in the *proximity only* scenario. However,

¹⁴ All confusion matrices supporting the analysis can be found on the RISIS project website at <http://sms.risis.eu/assets/pdf/metrics-link-network.pdf>.

even though in extreme cases (*proximity plus name*) the Majority Class Classifier takes the lead, the network metric does not fall far behind.

Table 5. Network-metric (e_Q) result versus the MCC baseline.

Majority Class Classifier (Baseline) vs Network Metrics (e_Q)						
<i>MajorityClassClassifier</i>						
<i>NetworkMetrics</i>						
GT = Ground Truth	GT_P = Ground Truth Positive	GT_N = Ground Truth Negative				
50m geo-only	GT=92	$GT_P=30$	$GT_N=62$	$F_1 : \frac{-}{0.693}$	ACC: $\frac{0.674}{0.75}$	NPV: $\frac{0.674}{0.915}$
500m geo-only	GT=249	$GT_P=66$	$GT_N=183$	$F_1 : \frac{-}{0.682}$	ACC: $\frac{0.735}{0.779}$	NPV: $\frac{0.735}{0.978}$
2km geo-only	GT=198	$GT_P=61$	$GT_N=137$	$F_1 : \frac{-}{0.803}$	ACC: $\frac{0.692}{0.859}$	NPV: $\frac{0.692}{0.966}$
50m geo+names	GT=31	$GT_P=27$	$GT_N=4$	$F_1 : \frac{0.931}{0.826}$	ACC: $\frac{0.871}{0.742}$	NPV: $\frac{-}{0.333}$
500m geo+names	GT=155	$GT_P=148$	$GT_N=7$	$F_1 : \frac{0.977}{0.909}$	ACC: $\frac{0.955}{0.839}$	NPV: $\frac{-}{0.179}$
2km geo+names	GT=342	$GT_P=322$	$GT_N=20$	$F_1 : \frac{0.97}{0.912}$	ACC: $\frac{0.942}{0.845}$	NPV: $\frac{-}{0.238}$

As in the first experiment, for further evaluation, we extracted a sample based on ILNs in which at least one organisation originates from the Netherlands. Out of the **107** sampled ILNs, the domain expert deviated from the first evaluation in only 1 case.

8 e_Q Put to a Ranking Test

The authors of the recently published paper [13] compared seven algorithms for clustering entities from multiple sources at different string similarity thresholds. They evaluated the quality of the clusters that these algorithms generated on three gold standard datasets¹⁵, one manually built (referred here as GT1), and two syntactically generated. We take the evaluation results from [13] on GT1, and then test if our e_Q score is able to correctly predict the ranking of the algorithms as found in the reported evaluation. In contrast to the earlier experiments (where we use e_Q to assess the quality of clusters), we are now testing if e_Q can be used to correctly rank different clustering algorithms across datasets.

A slightly complicating factor is that the evaluation in [13] relies on F_1 values computed on *true pairs of entities found*. Since e_Q evaluates entire clusters (i.e. *sets* of pairs of entities) of size greater than 2 ($S > 2$), we recompute the F_1 values based on *true clusters found* ($S > 2$) and plot these performance measures for each algorithm in Fig. 6 as *Baseline*. The resulting plot is comparable to the original one in [13]. We then ran the e_Q metric over the outputs of each algorithm at the same thresholds, displayed in Fig. 6 as e_Q *Evaluation*.

¹⁵ https://dbs.uni-leipzig.de/de/research/projects/object_matching/famer.

The results show that the ranking of the algorithms by e_Q (e_Q *Evaluation*) does not significantly deviate from the recomputed ranking (*Baseline*). This illustrates the usefulness of the e_Q metric as it demonstrates its potential to rank algorithms whenever they show *significant performance differences*.

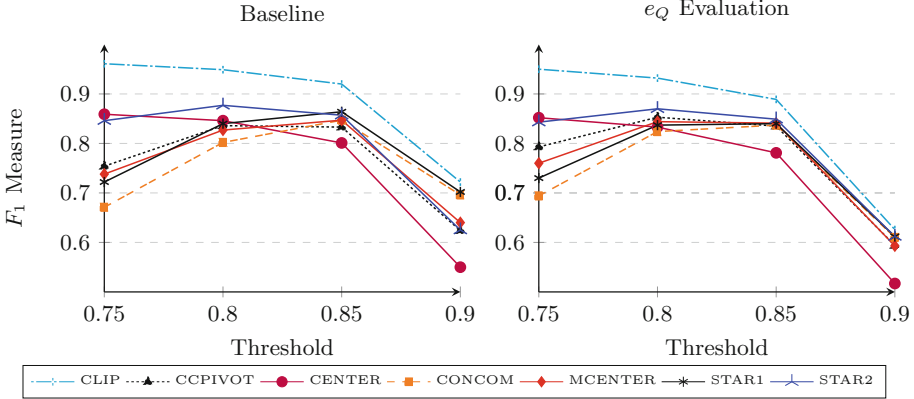


Fig. 6. Evaluation of e_Q on the ranking from [13]

9 Conclusions and Future Work

9.1 Conclusion

Entity resolution is an essential step in the use of multiple datasets on the semantic web. Since entity resolution algorithms are far from being perfect, the links they discover must often be human validated. Since this is both a costly and an error-prone process, it is desirable to have computer support that can accurately estimate the quality of links between entities. In this paper, we have proposed a metric for precisely this purpose: it estimates the quality of links between entities from multiple datasets, using a combination of graph metrics over the network (>2) formed by these links. Our metric captures the intuition that high redundancy in such a linking-network correlates with high quality.

We have tested our metric in three different scenarios. Using a collection of six widely used social science datasets in the first two experimental settings, we compared the predictions of link quality by our metric against human judgements on hundreds of networks involving thousands of links. In both evaluations, our metric correlated strongly with human judgement ($0.806 \leq F_1 \leq 1$), and it consistently beats the Majority Class Classifier baseline (except in cases where this is numerically near impossible because of a highly skewed class distribution). In the experimental condition where we deliberately constructed noisy and non-noisy link-networks, we showed that our metric is in general able to exclude poor networks in noisy environments and to include good networks in noiseless environments. With the last experiment, we also show that our metric is

able to rank entity resolution algorithms on their quality, using an externally produced dataset and corresponding ground truth. All this amounts to testing the e_Q metric on a dozen different algorithms and parameter settings. Across these different experimental conditions, our quality metric consistently agrees with human judgement.

To encourage replication studies and extensions to our work, all the datasets used in these experiments are available online.

9.2 Future Work

Including Link Strength. The metric e_Q is based on the presence and absence of links, but does not consider any strength associated with these links. We are currently working on refinements of e_Q that use link confidence scores produced by entity resolution algorithms.

Dynamic Link Adjustment. The current work simply takes the output of an entity resolution algorithm as given, and tries to estimate the quality of that output. A closer coupling between our metric and an entity resolution algorithm would allow the algorithm to dynamically adjust its output based on the e_Q quality estimates. Similarly, embedded in a user-interface, the score of our metric could help the user to give the final judgement to accept or reject an ILN.

Parameter Tuning. In this work, we empirically determined the 1.6 sigmoid hyper-parameter, the discrete e_Q intervals and the string similarity thresholds. Experimenting on fine-tuning these parameters using the current ground truths and data from other domains would help understanding how and when different choices could lead to an increase or a decrease of the metrics' predictive power.

Acknowledgement. We kindly thank *Paul Groth* for his constructive comments and proofreading, *Alieh Saeedi* for sharing her experiments data and supporting the reproducibility of their experiments, and the *EKAW reviewers* for constructive comments. This work was supported by the European Union's 7th Framework Programme under the project RISIS (GA no. 313082).

References

1. Baron, A., Freedman, M.: Who is who and what is what: experiments in cross-document co-reference. In: Proceedings of the Conference on Empirical Methods in Natural Language Processing, pp. 274–283. Association for Computational Linguistics (2008)
2. Cucerzan, S.: Large-scale named entity disambiguation based on Wikipedia data. In: Proceedings of the 2007 Joint Conference on EMNLP-CoNLL (2007)
3. David, J., Euzenat, J.: Comparison between ontology distances (Preliminary Results). In: Sheth, A., et al. (eds.) ISWC 2008. LNCS, vol. 5318, pp. 245–260. Springer, Heidelberg (2008). https://doi.org/10.1007/978-3-540-88564-1_16

4. David, J., Euzenat, J., Šváb-Zamazal, O.: Ontology similarity in the alignment space. In: Patel-Schneider, P.F., et al. (eds.) ISWC 2010. LNCS, vol. 6496, pp. 129–144. Springer, Heidelberg (2010). https://doi.org/10.1007/978-3-642-17746-0_9
5. Euzenat, J., Shvaiko, P.: *Ontology Matching*, 2nd edn. Springer, Heidelberg (2013)
6. Guéret, C., Groth, P., Stadler, C., Lehmann, J.: Assessing linked data mappings using network measures. In: Simperl, E., Cimiano, P., Polleres, A., Corcho, O., Presutti, V. (eds.) ESWC 2012. LNCS, vol. 7295, pp. 87–102. Springer, Heidelberg (2012). https://doi.org/10.1007/978-3-642-30284-8_13
7. Hassanzadeh, O., Kementsietsidis, A., Lim, L., Miller, R.J., Wang, M.: A framework for semantic link discovery over relational data. In: 18th ACM Conference on Information and Knowledge Management, pp. 1027–1036. ACM (2009)
8. Hassanzadeh, O., Xin, R., Miller, R.J., Kementsietsidis, A., Lim, L., Wang, M.: Linkage query writer. *Proc. VLDB Endow.* **2**(2), 1590–1593 (2009)
9. Li, W., Zhang, S., Qi, G.: A graph-based approach for resolving incoherent ontology mappings. In: *Web Intelligence*, vol. 16, pp. 15–35. IOS Press (2018)
10. Maedche, A., Staab, S.: Measuring similarity between ontologies. In: Gómez-Pérez, A., Benjamins, V.R. (eds.) EKAW 2002. LNCS (LNAI), vol. 2473, pp. 251–263. Springer, Heidelberg (2002). https://doi.org/10.1007/3-540-45810-7_24
11. Menestrina, D., Whang, S.E., Garcia-Molina, H.: Evaluating entity resolution results. *Proc. VLDB Endow.* **3**(1–2), 208–219 (2010)
12. Ngomo, A.-C.N., Auer, S.: Limes—a time-efficient approach for large-scale link discovery on the web of data. In: IJCAI, pp. 2312–2317 (2011)
13. Saeedi, A., Peukert, E., Rahm, E.: Using link features for entity clustering in knowledge graphs. In: Gangemi, A., et al. (eds.) ESWC 2018. LNCS, vol. 10843, pp. 576–592. Springer, Cham (2018). https://doi.org/10.1007/978-3-319-93417-4_37
14. Sarasua, C., Staab, S., Thimm, M.: Methods for intrinsic evaluation of links in the web of data. In: Blomqvist, E., Maynard, D., Gangemi, A., Hoekstra, R., Hitzler, P., Hartig, O. (eds.) ESWC 2017. LNCS, vol. 10249, pp. 68–84. Springer, Cham (2017). https://doi.org/10.1007/978-3-319-58068-5_5
15. Usbeck, R., et al.: AGDISTIS - graph-based disambiguation of named entities using linked data. In: Mika, P., et al. (eds.) ISWC 2014. LNCS, vol. 8796, pp. 457–471. Springer, Cham (2014). https://doi.org/10.1007/978-3-319-11964-9_29
16. Volz, J., Bizer, C., Gaedke, M., Kobilarov, G.: Discovering and maintaining links on the web of data. In: Bernstein, A., et al. (eds.) ISWC 2009. LNCS, vol. 5823, pp. 650–665. Springer, Heidelberg (2009). https://doi.org/10.1007/978-3-642-04930-9_41
17. Vrandečić, D., Sure, Y.: How to design better ontology metrics. In: Franconi, E., Kifer, M., May, W. (eds.) ESWC 2007. LNCS, vol. 4519, pp. 311–325. Springer, Heidelberg (2007). https://doi.org/10.1007/978-3-540-72667-8_23