# xEM: Explainable Entity Matching in Customer 360

Sukriti Jaitly*
Carnegie Mellon University
Pittsburgh, USA
sjaitly@andrew.cmu.edu

Deepa Mariam George*
IBM Data and AI
Bengaluru, India
deepa.george@ibm.com

Balaji Ganesan
IBM Research
Bengaluru, India
bganesa1@in.ibm.com

Muhammad Ameen
IBM Data and AI
Bengaluru, India
muhammed.abdul.majeed.ameen@ibm.com

Srinivas Pusapati
IBM Data and AI
Bengaluru, India
srinivas.pusapati@in.ibm.com

## ABSTRACT

Entity matching in Customer 360 is the task of determining if multiple records represent the same real world entity. Entities are typically people, organizations, locations, and events represented as attributed nodes in a graph, though they can also be represented as records in relational data. While probabilistic matching engines and artificial neural network models exist for this task, explaining entity matching has received less attention. In this demo, we present our Explainable Entity Matching (xEM) system and discuss the different AI/ML considerations that went into its implementation.

## CCS CONCEPTS

• **Information systems** → *Data mining*; • **Applied computing** → *Enterprise data management*.

## KEYWORDS

explainability, entity matching, graph neural networks

## 1 INTRODUCTION

Entity matching is the task of predicting if two entities belong to the same real world entity. This task is critical for managing *master data* in enterprises, governments and many commercial applications. Master data refers to the critical customer data that organizations maintain. Master Data Management (MDM) refers to a group of products that help organizations manage this master

---
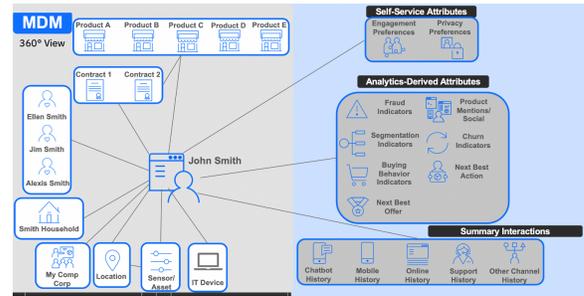
*Work done while interning at IBM Data and AI, India

**Figure 1: Customer 360 provides a 360 degree view of a customer in an enterprise data fabric**
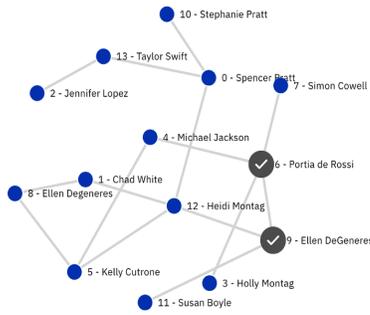
data. [9] introduces master data management and Customer 360 in great detail.

As shown in Figure 1, customer 360 provides a 360 degree view of the customer. Entity Matching is a core component of Customer 360 with multiple use-cases as we describe in this demo. People and organization entities are of particular interest though entity matching techniques are applicable to locations, events, products, and even abstract ideas like compliance clauses and law points in legal documents.

In Customer 360, entity matching is transitive. If a record A matches with record B; and record B matches with record C; all three records A, B and C will be linked together. This transitive linking is very useful and it helps in matching the records with partial matches. But this comes with the problem of some false positives and creating large entities of records in the system. Understanding and explaining these large entities have been a challenge so far.

In real customer scenarios, occurrence of entities of size 1000 and above is not uncommon, though most entities have fewer records. Visualization and explanation of such large entities using path based approaches [4] is hard because of scale. Even when the entities are not large, we need to explain why different records have been assigned to an entity.

We propose a solution to this problem by treating the relational data in a Customer 360 instance as graphs. Each record in the data becomes a node in our graph. The edges are only between records within an entity. Typically, there is an edge between a record and the representative record for the entity. Once the graph is in place, we use node embeddings from Graph Neural Networks [4], the

**(a) Graph structure around candidates to be matched**



**(b) Comparison of node attributes**

**Figure 2: Entity Matching could be explained both by node attributes and the adjoining graph structure**

scores from the probabilistic matching engine [9], to explain entity matching.

Some of the benefits of explaining entity matching include, identifying weak links or alternatively gluing members in the entity formation, identifying false positives, identification of matching on anonymous values, and applying manual unlink rules.

But before we proceed to our solution, we'll describe our current system in Section 3 and two other baselines that we compared our solution against. Our demo described in Section 4, is available from our research group page [1].

## 2    RELATED WORK

[9] describes the Probabilistic Matching Engine that is at the core of our entity matching solution. A number of heuristics have been developed over years, for name, addresses, phone numbers, identification numbers that are typical node attributes in Enterprise graphs. From finding edit distance, to complex statistical models, each attribute is handled differently. [7] presented a deep neural model for entity matching. [10] proposed SystemER, an active learning based approach for entity resolution.

A closely related problem to entity matching in graphs is node similarity which we have described in [8] and [3]. Our explainability techniques using GNN models and explainers can also be applied to the node similarity task.

GNN explainability techniques in the literature include [17], [19], [16], [12] and [6]. These techniques typically produce a subgraph as an explanation for the predicted node class or a link between two nodes.

We have in prior works attempted to substantiate GNN model predictions using information retrieval [4], path ranking [5] and reasoner based explanations [2]. In [14], we had used a random forest model for post-hoc explanations, while in [13], we used ideas from Anchor Explanations [11] and tried on graphs models. In [15], we sought to automate the evaluation of these explainability techniques, since explanations are subjective and human evaluation is cost prohibitive.

## 3    BASELINES

In this section, we'll first describe our current system which we seek to improve. Match 360 is an IBM product for the Customer 360 use-case and more generally for master data management.

We then implemented the models in DeepMatcher, LEMON baselines, and our own solution based on GCN and GNN Explainer. We use a combination of all three techniques to explain entity matches in different scenarios.

We evaluated the state of the art [1] solution on both the Amazon-Google dataset and a synthetic organization dataset that we have created. Our synthetic dataset consists of records business name and address details with more than 10 thousand tuples.

### IBM Match 360

As discussed in the previous section, transitive linking in entity matching for Master Data Management(MDM) solutions can lead to problems of false positives and large entities. Match360 is a modern MDM based solution by IBM that works with enterprise data to perform indexing, matching and linking of data from different sources (e.g. CRM, Experian, Salesforce, Web Portal), creating a 360 degree view of customer data.

Matching record pair data in Match360 requires comparing different record attributes (e.g. Name, Address, DOB, Identifier) from each pair of records to determine if they match and should subsequently be linked, based on a series of mathematically derived statistical probabilities and complex weight tables.

Such solutions that rely on Probabilistic Matching Engines(PME) for entity matching often provide very little insight into the entities making it difficult for customers to understand why such an entity was formed, or to explain them. An example of the attributes in a pair of matched records is as shown in Figure 2b.

### DeepMatcher

DeepMatcher [7] performs matching on labelled tuple pairs by training a neural network to predict matching. DeepMatcher adapts the RNN architecture to aggregate the attribute values and then compares/aligns the aggregated attribute representations. Deep Matcher trains the word embeddings using FastText.

---

[1]https://researcher.watson.ibm.com/researcher/view_group.php?id=11043

Out[36]:

Prediction: Not match (6%)

| | |
|---|---|
| name | sony pink cyber-shot 7.2 megapixel digital camera dscw120p |
| description | sony pink cyber-shot 7.2 megapixel digital camera dscw120p 7.2 megapixel 4x optical zoom 2.5 ' tft lcd 15 mb internal memory face detection super steadyshot image stabilization smile shutter mode smart zoom pink finish |
| price | |

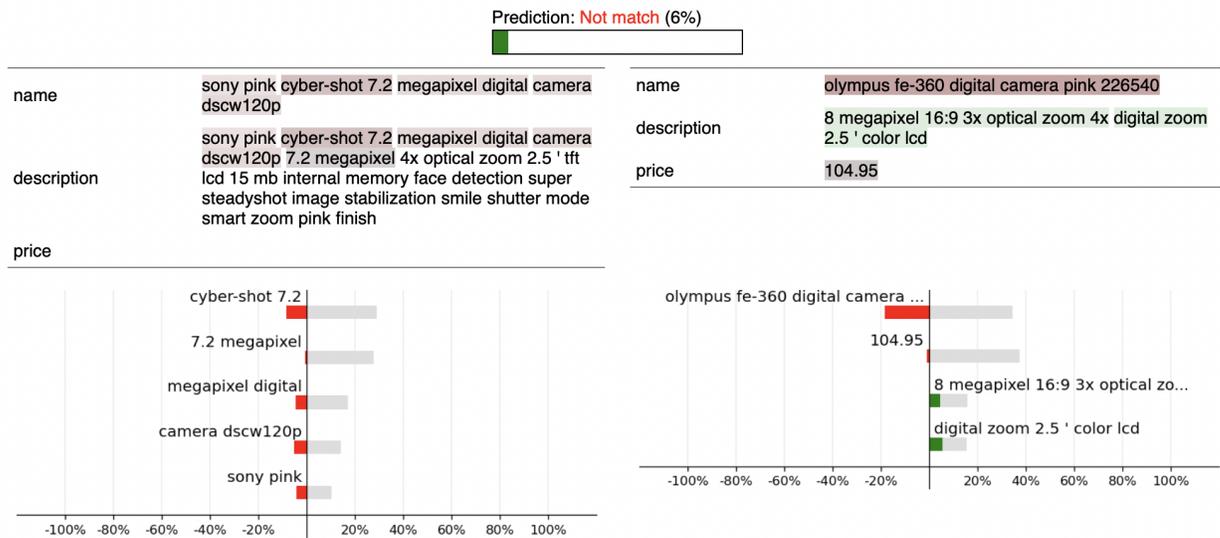| | |
|---|---|
| name | olympus fe-360 digital camera pink 226540 |
| description | 8 megapixel 16:9 3x optical zoom 4x digital zoom 2.5 ' color lcd |
| price | 104.95 |

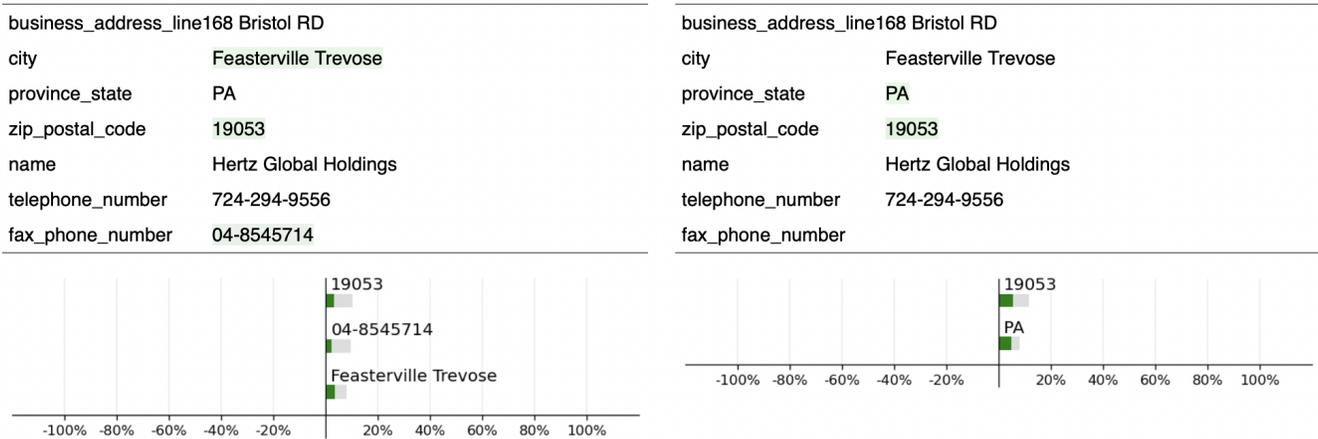Figure 3: Explaining a non-match in Amazon-Google dataset using LEMON

Figure 4: Explaining a match in a synthetic organization dataset using LEMON

## LEMON

LEMON [1] is a model-independent and schema-flexible approach for evaluating explainable entity matching. This approach is effective at communicating to the user the location of the decision border, particularly in the case of non-matches.

| Dataset | Precision | Recall | F1 |
|---|---|---|---|
| Amazon-Google | 0.79 | 0.38 | 0.52 |
| Synthetic Org Dataset | 0.37 | 0.37 | 0.37 |

Table 1: LEMON results on entity matching datasets

As shown on Table 1, we are able to roughly reproduce the performance of LEMON on the Amazon-Google dataset. But our performance on the synthetic dataset remains poor because of various challenges including the absence of any long text columns in the dataset. A typical Customer 360 dataset is similar to the examples shown in Figure 4 and hence this is a limitation of LEMON like models from being used for this task.

Entity Matching plays a critical role in data fabric in general and data marketplace in particular. By matching records from one source to another source, we can determine relationships between datasets that may otherwise not be feasible only using metadata. We observe that LEMON explanations are particularly suited to this usecase. As shown in Figures 3 and 4, we can use the explainable entity matching techniques to verify the relationship between datasets in a data fabric. We show examples of both a match and non-match. Each bar shows the positive or negative contribution of the feature to the eventual match or not-match prediction.
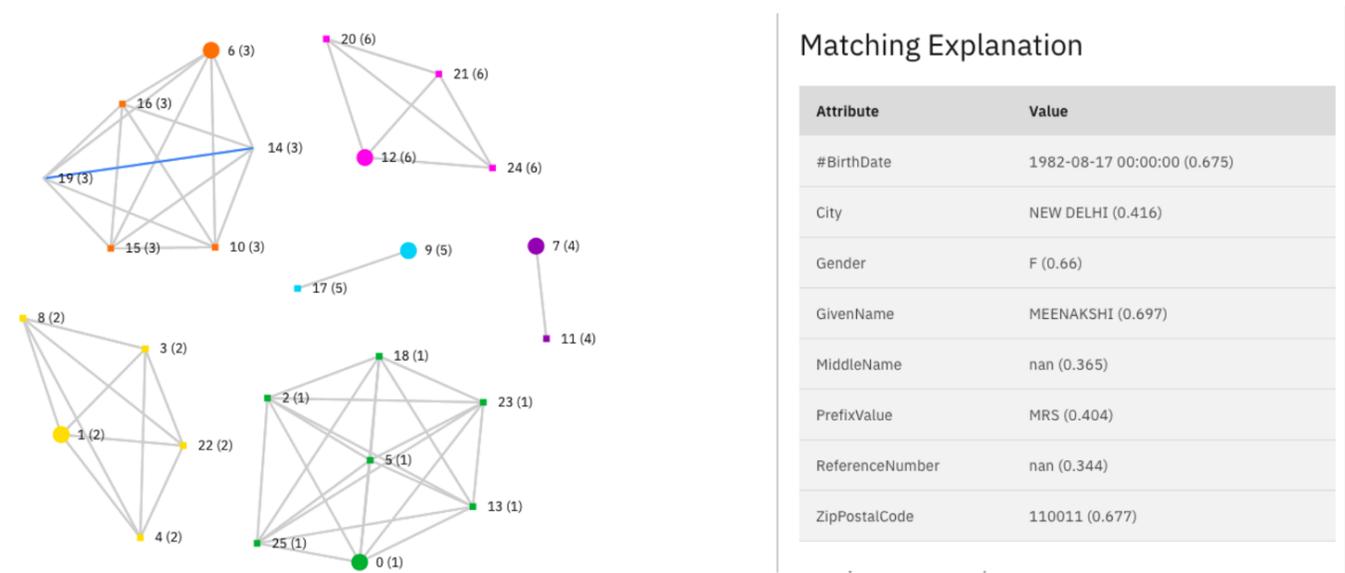
**Figure 5: Explaining entities in Customer 360**

## 4 DEMO

Our solution for explaining entities in Customer 360 is to augment LEMON and other pair wise entity matching explanations with Graph Neural Network model predictions and explanations. We treat records as nodes, and edges of this graph indicate that two records match pair-wise. A cluster of all such records makes up an entity.

Our solution involves training a GNN network in batch mode on the output of the Match360, where the training would be using the pairwise comparison scores of pairs of records that PME considered during matching the dataset. We use a Graph Convolution network as implemented in [18] to train our model. Our entity matching system is treated as a blackbox and the trained GCN model acts a proxy for the underlying entity matching system.

After training the GNN network, during inference time, we can use the GNN model to make predictions on a limited number of pairs of nodes from the entity we intent to look into. We then proceed to explain each of these predictions using a GNN explainability solution like GNNExplainer[17] and identify important features. These important features and their values are then passed as results to customers, helping them to examine different parts of entities and also identify any errors in the entity matching.

We are able to explain why a record is part of this entity, by explaining why it is related to another record in the entity. By way of explanation, we highlight the important node features. Unlike a typical knowledge graph or a social network graph, there are no *friend*, *parent* and other kinds of relations. Hence we do not show the edge masking, but only highlight the important features using feature masking. This is in contrast to typical GNN Explanations which are sub-graphs highlighting both node and edge masks. We believe converting an explanation subgraph into a tabular form like in Figure 5 makes it easier for end users to understand the explanation.

Simultaneously the different records that make up the nodes are displayed below the graph for any relational querying (not shown in the figure). There are multiple use-cases where this solution is being used. The default use-case is explaining why an entity has been formed by resolving multiple records (nodes) into an entity. Semantic matching where records from one source (dataset) are matched to another dataset is another use-case.

Our solution is deployed on the IBM Cloud using a Code Engine. Front end is a reactJS application while the backend is a python flask application serving REST API endpoints. Both the front end and the backend are deployed as IBM Code Engine applications. They can be bundled together onto a single container.

While these containers can be deployed both as separate microservices run on-prem, they can also be used as a cloud service on IBM Cloud. Coupled with Watson Knowledge Studio, Watson ML and Openscale, our explainability solution can be used by any customers who have Customer 360 workloads. A third monolith software approach for some of the legacy use-cases is also being considered.

In [2], we had discussed ways to evaluate a typical explainability solution using neuro-symbolic reasoning. This is because, unlike entity matching or matching between records in different datasets, entity resolution in Customer 360 is typically a clustering problem. We leave both the neuro-symbolic evaluation of explanations and graph clustering explanations for future work.

## CONCLUSION

In this demo, we discussed the state of the art explainability techniques for entity matching and showed how the explanations from existing literature are inadequate for the Customer 360 use-case. We then introduced our GNN based Explainable Entity Matching (xEM) system and discussed the different AI/ML considerations that went into its implementation.

## REFERENCES

[1] Nils Barlaug. 2022. LEMON: explainable entity matching. *IEEE Transactions on Knowledge and Data Engineering* (2022).

[2] Vanya BK, Balaji Ganesan, Aniket Saxena, Devbrat Sharma, and Arvind Agarwal. 2021. Towards Automated Evaluation of Explanations in Graph Neural Networks. arXiv:2106.11864 [cs.AI]

[3] Jaspreet Singh Dhani, Ruchika Bhatt, Balaji Ganesan, Parikshet Sirohi, and Vasudha Bhatnagar. 2021. Similar cases recommendation using legal knowledge graphs. *arXiv preprint arXiv:2107.04771* (2021).

[4] Balaji Ganesan, Gayatri Mishra, Srinivas Parkala, Neeraj R Singh, Hima Patel, and Somashekar Naganna. 2020. Link Prediction using Graph Neural Networks for Master Data Management. *arXiv preprint arXiv:2003.04732* (2020).

[5] Balaji Ganesan, Hima Patel, and Sameep Mehta. 2020. Explainable Link Prediction for Privacy-Preserving Contact Tracing. SpicyFL 2020 Workshop at NeurIPS 2020.

[6] Dongsheng Luo, Wei Cheng, Dongkuan Xu, Wenchao Yu, Bo Zong, Haifeng Chen, and Xiang Zhang. 2020. Parameterized explainer for graph neural network. *Advances in neural information processing systems* 33 (2020), 19620–19631.

[7] Sidharth Mudgal, Han Li, Theodoros Rekatsinas, AnHai Doan, Youngchoon Park, Ganesh Krishnan, Rohit Deep, Esteban Arcaute, and Vijay Raghavendra. 2018. Deep learning for entity matching: A design space exploration. In *Proceedings of the 2018 International Conference on Management of Data*. 19–34.

[8] Phillipp Müller, Xiao Qin, Balaji Ganesan, Nasrullah Sheikh, and Berthold Reinwald. 2020. An Integrated Graph Neural Network for Supervised Non-obvious Relationship Detection in Knowledge Graphs.. In *EDBT*. 379–382.

[9] Martin Oberhofer, Eberhard Hechler, Ivan Milman, Scott Schumacher, and Dan Wolfson. 2014. *Beyond big data: Using social MDM to drive deep customer insight*. IBM Press.

[10] Kun Qian, Lucian Popa, and Prithviraj Sen. 2019. Systemer: A human-in-the-loop system for explainable entity resolution. (2019).

[11] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. 2018. Anchors: High-precision model-agnostic explanations. In *Thirty-Second AAAI Conference on Artificial Intelligence*.

[12] Michael Sejr Schlichtkrull, Nicola De Cao, and Ivan Titov. 2020. Interpreting graph neural networks for nlp with differentiable edge masking. *arXiv preprint arXiv:2010.00577* (2020).

[13] Anjali Singh, Balaji Ganesan, et al. 2021. Reimagining GNN Explanations with ideas from Tabular Data. *arXiv preprint arXiv:2106.12665* (2021).

[14] Lingraj S Vannur, Lokesh Nagalapatti, Balaji Ganesan, and Hima Patel. 2020. Data Augmentation for Personal Knowledge Graph Population. *arXiv preprint arXiv:2002.10943* (2020).

[15] BK Vanya, Muhammed Abdul Majeed Ameen, Balaji Ganesan, Devbrat Sharma, and Arvind Agarwal. 2021. Automated Evaluation of GNN Explanations with Neuro Symbolic Reasoning. In *Annual Conference on Neural Information Processing Systems*.

[16] Minh N. Vu and My T. Thai. 2020. PGM-Explainer: Probabilistic Graphical Model Explanations for Graph Neural Networks. arXiv:2010.05788 [cs.LG]

[17] Zhitao Ying, Dylan Bourgeois, Jiaxuan You, Marinka Zitnik, and Jure Leskovec. 2019. Gnnexplainer: Generating explanations for graph neural networks. In *Advances in neural information processing systems*. 9244–9255.

[18] Jiaxuan You, Rex Ying, and Jure Leskovec. 2019. Position-aware graph neural networks. In *International conference on machine learning*. PMLR, 7134–7143.

[19] Hao Yuan, Jiliang Tang, Xia Hu, and Shuiwang Ji. 2020. Xgnn: Towards model-level explanations of graph neural networks. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. 430–438.