

Applying triplet loss to entity resolution models

Charissa Ding, Lisa Kim, Derek Zhao

1 Abstract / Introduction

Entity resolution refers to the task of identifying records from different data sources that refer to the same entity. It is one of the most crucial stages in data integration, and has been spent the past few decades studying various methods to make this process both accurate and efficient.

The most fundamental technique used in entity resolution is to manually design and calculate similarity metrics between a pair of records and setting a threshold to decide whether the pair is a match [1]. Additional efforts have been built upon this technique, including machine-learning algorithms such as decision tree ensembles [2, 3], as well as combinations of active learning and crowdsourcing [4, 5]. Some of the most recent research efforts on the subject explore techniques such as deep learning algorithms [6] and tree-based adversarial generation. While these models can achieve near-perfect accuracy on datasets that are well structured and have little noise, they struggle to perform well with datasets that are noisy and have unstructured attributes [6].

In recent years, researchers have made significant progress in face recognition utilizing a deep convolutional network trained with triplet loss function, setting a new record accuracy of 99.63% on the widely used Labeled Faces in the Wild (LFW) dataset with more than 13,000 face images [7]. In this project, we propose to apply this method to the entity resolution problem, aiming to build a model that is both accurate and robust as the level of complexity and amount of noise grow in a given dataset.

This problem is relevant because improving accuracy on entity resolution would greatly improve the overall performance of a data integration task, one of the key topics in this course.

2 One Sentence Summary

We plan to adapt the use of the triplet loss function which has shown success in facial recognition systems to a deep neural network trained on triplets of database records in hopes that such a model achieves better performance than current state of the art techniques for entity resolution, such as random forests.

3 Audience and Needs

This project will be of particular value to the academic community, database managers, and data wranglers. If successful, this project would suggest a new avenue for research in entity resolution. Furthermore, a computationally reasonable model that can resolve entities with high accuracy will lead to more efficient matching, more reliable de-duplication procedures, and more accurate table merging. At its most practical, database managers and data wranglers might be saved substantial amounts of time that would otherwise be spent writing rules and manually examining records.

4 Approach

The problem of entity resolution bears some resemblance to that of facial recognition. In both cases, we are less concerned about what a record is or who a face belongs to and more interested in whether two records or two faces are the same. Recent innovations in facial recognition technology uses neural networks trained on a triplet loss function. In a triplet loss network for facial recognition, rather than predicting the class of an input image X , the network outputs a multi-dimensional vector embedding $f(X)$ of that image with the constraint that the Euclidean

distance between two embeddings of images of two different people is significantly large. That is, let $\mathbb{I}(X)$ be the identity of the person whose face appears in X . If $\mathbb{I}(X) = \mathbb{I}(Y)$ and $\mathbb{I}(X) \neq \mathbb{I}(Z)$, then $\|f(X) - f(Z)\| \gg \|f(X) - f(Y)\|$.

A training instance for such a network consists of a triplet of three images (A, P, N) , an anchor, positive, and negative image, respectively, that satisfies $\mathbb{I}(A) = \mathbb{I}(P)$ and $\mathbb{I}(A) \neq \mathbb{I}(N)$. Given said triplet, the loss that they incur on a network is $\mathcal{L}(A, D, P) = \max(0, \|f(A) - f(P)\|^2 - \|f(A) - f(N)\|^2 + \alpha)$, where α is a threshold hyperparameter defined prior to training.

The above loss function has proven successful for face recognition, and we believe it can be applied similarly to entity resolution. To generate triplets, we will use a *discriminator model* to predict whether pairs of records refer to the same entity. The mispredictions of the discriminator model will either consist of false positives which can be used as the A, N records within a triplet or false negatives which can be used as the A, P records within a triplet. Ideally, we would like to show that given an initial discriminator model, using its mispredictions as training examples for a triplet loss network results in a neural model that outperforms the discriminator.

5 (Best Case) Impact

We can show that the algorithm produces more accurate results than state-of-the-art ML and crowd approaches, including the deepER algorithm, which on average has a 0.6 to 0.8 F-measure on challenging and noisy datasets [6]. In addition, our model should also preserve fast runtimes and not require prohibitive amounts of training data.

6 Milestones

1. Identify datasets tested in previous methodologies that are computationally feasible to train on and that we can use as benchmarks.
2. Develop an adversarial procedure for generating triplets. This will entail developing a tree-based ensemble and corresponding featurization procedures to generate false positives and false negatives. The accuracy of this procedure may also serve as a benchmark.

3. Design a reasonable means of featurizing records as input for a neural network.
4. Implement a neural network based on the triplet loss function and compare against benchmark used to generate "hard" negatives.
5. Develop data augmentation procedures to generate positive pairs and improve the performance of the triplet loss network.
6. Run experiments on various datasets and compare results with current widespread methods.

7 Obstacles

*Major obstacles are situations where we would consider **killing** the project. Minor obstacles are situations that would delay the project or increase the overall cost in energy, time, people, and money.*

7.1 Major obstacles

- If we cannot show that a triplet loss network can at the very least improve on the performance of the initial model used to generate false positives and negatives, it is highly unlikely that said network can compare favorably to the state of the art.

7.2 Minor obstacles

- The original purpose of the triplet loss function was to enable facial recognition systems to compare new images of faces against a database of faces the system has already trained on. In other words, a facial recognition system is responsible for determining whether a new face is the same as one it has been trained to recognize, but it is not responsible for determining whether two new faces are the same. It is possible the triplet loss function is effective for only the former and not the latter, in which case a triplet loss network would require significantly more training data and computational resources to be effective.
- A typical dataset or even pair of datasets with duplicate entities largely consists of negative pairings. While an initial model such as a random forest can be used to pick out negative pairings most effective for a triplet loss network, we

may be constrained by the number of positive pairings, in which case we will need to devise data augmentation methods to generate additional positive pairings to increase the number of training triplets.

- We may discover that the neural network architecture needed to effectively embed records as vectors is prohibitively large for the computational resources available. If so, we should consider training on datasets with fewer attributes or using simpler benchmark models for initially detecting false positives and negatives.

8 Additional Resources

- Some computational time to run our neural network algorithm.
- Access to a GPU-enabled server in order to perform parallel processing in the neural network model.

9 Literature Review

- This study builds on prior work in entity matching [1], particularly the use of supervised learning to differentiate matching and non-matching record pairs [8].
- Some preliminary work on embedded representations of attributes and entire records suggest a new approach towards calculating the similarity between two records [6]. We will build on this direction and borrow the use of the triplet loss function developed for facial recognition [7].
- We are not aware of other studies applying distance metric learning such as triplet loss to the problem of entity resolution. However, the use of word vector embeddings to generate distributed representations of records has outperformed the state of the art [6].
- Hybrid human-machine workflows that combine crowdsourcing and active learning [5, 4] as well as tree-based ensembles for supervised learning [2] are alternate strands of research towards more efficient and more reliable entity matching.

10 Define Success

Developing a model that can exceed the entity resolution accuracy set by state-of-the-art machine learning and crowd-sourcing approaches would be considered a success, particularly if our model shows high accuracy given noisy training data.

References

- [1] A. K. Elmagarmid, P. G. Ipeirotis, and V. S. Verykios, "Duplicate record detection: A survey," *IEEE Transactions on knowledge and data engineering*, vol. 19, no. 1, pp. 1–16, 2007.
- [2] S. Varma, N. Sameer, and C. R. Chowdary, "Relic: Entity profiling by using random forest and trustworthiness of a source-technical report," *arXiv preprint arXiv:1702.00921*, 2017.
- [3] L. Yi, D. Xing-chun, C. Jian-jun, Z. Xing, and S. Yu-ling, "A method for entity resolution in high dimensional data using ensemble classifiers," *Mathematical Problems in Engineering*, vol. 2017, 2017.
- [4] C. Gokhale, S. Das, A. Doan, J. F. Naughton, N. Rampalli, J. Shavlik, and X. Zhu, "Corleone: hands-off crowdsourcing for entity matching," in *Proceedings of the 2014 ACM SIGMOD international conference on Management of data*, pp. 601–612, ACM, 2014.
- [5] J. Wang, T. Kraska, M. J. Franklin, and J. Feng, "Crowder: Crowdsourcing entity resolution," *Proceedings of the VLDB Endowment*, vol. 5, no. 11, pp. 1483–1494, 2012.
- [6] M. Ebraheem, S. Thirumuruganathan, S. Joty, M. Ouzzani, and N. Tang, "Deeper-deep entity resolution," *arXiv preprint arXiv:1710.00597*, 2017.
- [7] F. Schroff, D. Kalenichenko, and J. Philbin, "Facenet: A unified embedding for face recognition and clustering," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 815–823, 2015.
- [8] H. Köpcke, A. Thor, and E. Rahm, "Evaluation of entity resolution approaches on real-world match problems," *Proceedings of the VLDB Endowment*, vol. 3, no. 1-2, pp. 484–493, 2010.