

An Overview of Approaches to Quantify Open Data Catalog Similarity

Jorge Martinez-Gil

*Software Competence Center Hagenberg GmbH
Softwarepark 32a, 4232 Hagenberg, Austria
jorge.martinez-gil@scch.at*

Abstract

As open data initiatives continue to gain importance, the need for effective methods to assess the similarity between different open data catalogs becomes increasingly essential. The task of measuring catalog similarity can be helpful in many processes, such as catalog curation, data discovery, and interconnectivity between various open data repositories. This research provides an overview of existing approaches to quantify the similarity between open data catalogs. We explore various strategies ranging from the use of traditional methods based on comparing triples to advanced semantic-based and hashing methods for specific domain languages. Additionally, we identify key challenges and future research directions in open data catalog similarity measurement.

Keywords: Open Data, Data Catalogs, Semantic Similarity Measurement

1. Introduction

In times when the use of data has exploded exponentially, open data catalogs (ODCs) have emerged as an indispensable source of information that facilitates transparency, innovation, and multidisciplinary collaboration [1]. These catalogs, which often must be manually curated through government agencies, research organizations, or even non-profit organizations, offer invaluable knowledge spanning a wide range of disciplines, industries, and geographies [8]. However, as the volume of ODCs skyrockets, methods and tools are needed to facilitate their proper processing. For example, to assess their quality [2], to quantify their similarity [14], to use them in an effective way [22], and so on.

The concept of catalog similarity goes far beyond the mere comparison of metadata. The reason is that behind this concept lies a significant challenge in interconnecting different data repositories. Therefore, the possibility of automatically determining the similarity of ODC can positively impact a large group of professionals involved in data economics, for example, researchers, policymakers,

scientists, etc. In this way, providing methods and tools that facilitate the effective and efficient comparison of ODCs can facilitate decision-making processes and multidisciplinary collaborations, and it holds great potential for the collective use and analysis of open data [21].

This research is designed to explore the various possibilities for measuring the similarity of ODCs. To this end, we aim to explore the different methodologies and techniques to address this challenge. In fact, we aim to provide a clear roadmap for ODC similarity assessment by surveying existing approaches in this context. Our goal is to identify the options for ODC similarity measurement, allowing the user to choose the most appropriate method for the needs of their specific scenarios. Therefore, we aim to boost the full potential of open data in a wide range of application scenarios.

The remainder of this paper is structured as follows: Section 2 introduces related work concerning similarity between ODCs and challenges that remain pending. Section 3 provides an illustrative overview of existing methods to address the challenge of ODC similarity using different computational approaches. Finally, Section 4 presents the lessons that can be drawn from this research and possible future research directions.

2. Related Works

ODCs are typically used to describe and organize data assets, making them more discoverable and accessible for the stakeholders of an organization. The standard DCAT (Data Catalog Vocabulary) [9] is a widely used specification for describing data catalogs, often used with government and open data portals. When comparing two DCAT data catalogs, it is possible to find similarities in various aspects inherent to them.

The aspects that should be considered in a specific situation are an open discussion. However, some viewpoints suggest assessing metadata elements for similarities, including title, description, keywords, publisher, contact information, access URLs, data formats, and licenses, should work fine in general cases. The organizational structure of the catalogs, whether hierarchical or based on tags and keywords, is also a widely suggested option for cases of this kind.

Other aspects that could be considered are compliance with standards like DCAT and additional metadata elements or support for standards like CKAN¹ or Dublin Core². Also, the data format and licensing model compatibility suggest similarities in ODCs. The search and discovery capabilities, including filters and faceted search [19], and other advanced options could also be

¹<https://ckan.org/>

²<https://www.dublincore.org/>

evaluated, along with methods for accessing data such as direct downloads or APIs. Additionally, the sources and diversity of data, data coverage, versioning support, metadata quality, update frequency, and the use of linked data principles for semantic relationships can also be sources of great value to be investigated.

However, several challenges remain open in the field, and there is still a need for solutions to handle them, including standardized metadata schemas, addressing data quality issues [13], and managing evolving catalogs. Additionally, as open data initiatives evolve, future research should focus on dynamic and real-time similarity assessment, multi-modal data catalogs, and cross-domain catalog comparisons. In this work, we try to shed light on the topic of similarity [10], as we try to identify ways in which similarity could be determined automatically.

3. Similarity Methods for Open Data Catalogs

Similarity methods for ODCs play a crucial role in enhancing the discoverability [16] of vast and diverse datasets available in the public domain. These methods could, for example, employ techniques such as natural language processing [4, 17], metadata analysis through machine learning [6, 7], and other kind of semantic-based techniques [18], or even sophisticated combinations of all of these methods [5], to establish connections and relationships between datasets based on their content, structure, and context. Although this topic is still very incipient, some authors have some works that focus additionally on the interpretability [15] to help stakeholders understand how the similarity value is calculated.

In this work, we focus on calculating similarity scores [11] since this should enable users to find datasets relevant to their specific needs, even when they may not know the exact dataset names or categories. This facilitates more efficient data discovery and encourages data reuse and integration across domains. This can lead to increased collaboration, innovation, and democratizing data-driven insights [12]. As the volume of open data grows, similarity methods remain essential to work with ODCs and address complex challenges and data-driven decision-making processes.

In the context of this work, we are going to focus on an illustrative technique for each of the existing approaches:

1. Considering the two ODCs as two repositories of triples.
2. Considering the two ODCs as two repositories of tokens.
3. Considering the two ODCs as two character sequences.
4. Considering the two ODCs as two documents written in a common purpose-specific language.

We will now go deeper into the technical details of these approaches, exploring their unique attributes, advantages, and potential challenges to gain a comprehensive understanding of their applicability and potential.

3.1. Considering the two ODCs as two repositories of triples

Calculating the similarity between two ODCs, considering them as repositories of triples, involves comparing the triples associated with the datasets in each ODC to determine their similarity. This process can be helpful in various data-related tasks, such as data integration or alignment, where the people work with a granularity at the level of triples.

The idea is simple and based on calculating the similarity between two sets of triples by iterating through each triple in both ODCs and counting the number of triples common to both ODCs. Algorithm 1 shows us how an algorithm could implement this. The significant advantage of this method is its simplicity and ease of understanding.

Algorithm 1 Calculate Similarity Between Two Repositories of Triples

```
1: Initialize an empty RDF graph  $g$ 
2: Parse ODCs into  $g$ 
3: Initialize two empty sets  $triples$  and  $triples2$ 
4: for each triple  $(s, p, o)$  in  $g$  where  $s = dataset\_001$  do
5:   Add  $o$  to  $triples$ 
6: end for
7: for each triple  $(s, p, o)$  in  $g$  where  $s = dataset\_002$  do
8:   Add  $o$  to  $triples2$ 
9: end for
10: Initialize  $similarity$  to 0
11: for each triple  $p1$  in  $triples$  do
12:   for each triple  $p2$  in  $triples2$  do
13:     if  $p1 = p2$  then
14:       Increment  $similarity$  by 1
15:     end if
16:   end for
17: end for
18: Calculate  $similarity$  as  $similarity / \max(\text{length of } triples, \text{length of } triples2)$ 
```

3.2. Considering the two ODCs as two repositories of tokens

Calculating the similarity of two ODCs as two repositories of tokens, whether ordered or unordered, can offer several advantages. For example, tokenizing and comparing data catalogs based on their tokens allows us to identify similarities without loading and processing the entire datasets

quickly. Moreover, the calculations to be performed are highly scalable. The reason is that, as the size of the ODC grows, the computational cost of comparing tokens remains relatively low.

Algorithm 2 encapsulates the *CheckSimilarity* function to calculate the similarity between two ODCs (converted into sets of RDF data) using TF-IDF vectorization and cosine similarity. The idea is first to convert the ODC into TF-IDF feature vectors and then compute the cosine similarity between these vectors, representing the similarity between the two ODCs. The result is scaled to a percentage and returned as the final similarity score. Comparing the similarity is commonly used in many natural language processing tasks.

Algorithm 2 Calculate Similarity Between Two Repositories of Tokens

```

1: function SIMILARITY(odc, odc2)
2:   vectorizer  $\leftarrow$  TfidfVectorizer()
3:   tfidf_matrix  $\leftarrow$  vectorizer.fit_transform([odc, odc2])
4:   cosine_sim  $\leftarrow$  cosine_similarity(tfidf_matrix[0], tfidf_matrix[1])[0][0]
5:   return cosine_sim
6: end function

```

3.3. Considering the two ODCs as two character sequences

Calculating the similarity of two ODCs using the Longest Common Subsequence (LCS) [3] method involves finding the longest sequence of items common to both catalogs. LCS provides a structured approach to aligning sequences of characters. This alignment helps understand how data items in one catalog correspond to those in the other, which can be crucial for data matching and integration. The LCS method is often used in various fields to measure the similarity between sequences, such as information retrieval, text analysis, etc.

Algorithm 3 is intended to construct a matrix dynamically and iteratively comparing their elements. It utilizes dynamic programming to find the maximum common subsequence length and returns this value.

Algorithm 4 calculates the similarity between two ODCs. The idea is first to determine their LCS's length and then normalize it by dividing it by the maximum length of the input sequences. The result is a similarity score between 0 and 1, where 1 indicates complete similarity, and 0 indicates no commonality.

The great advantage of the LCS method is that it primarily measures structural similarity based on the order of items in the catalogs. If the order of items is crucial for the analysis, the LCS method can be suitable. However, in cases where the item order is less important, it is better

Algorithm 3 Calculation of the Longest Common Subsequence

```
1: function LCS(odc, odc2)
2:   m ← length(odc)
3:   n ← length(odc2)
4:   lcs_matrix ← initialize a 2D array of size  $(m + 1) \times (n + 1)$  filled with zeros
5:   for i ← 1 to m do
6:     for j ← 1 to n do
7:       if odc[i − 1] = odc2[j − 1] then
8:         lcs_matrix[i][j] ← lcs_matrix[i − 1][j − 1] + 1
9:       else
10:        lcs_matrix[i][j] ← max(lcs_matrix[i − 1][j], lcs_matrix[i][j − 1])
11:      end if
12:    end for
13:  end for
14:  return lcs_matrix[m][n]
15: end function
```

Algorithm 4 ODC Similarity Calculation using LCS

```
1: function SIMILARITY(odc, odc2)
2:   lcs_length ← LCS(odc, odc2)
3:   return  $\frac{lcs\_length}{\max(\text{length}(odc), \text{length}(odc2))}$ 
4: end function
```

to consider other token-based measures, which focus on item presence or vector representations of the ODCs.

3.4. Considering the two ODCs as two documents written in a common purpose-specific language

These methods are known for their computational efficiency, especially when working with huge volumes of text data. It can quickly process large volumes of data, making it a potentially fast approach for performing calculations. To illustrate this approach, we have chosen the Winnow hashing algorithm [20], which works very well in identifying similarities between specific-purpose languages.

The Algorithm 5 adapts the famous Winnow algorithm that takes a text and an integer k , splitting the text into k -grams, hashing them, and finding the minimum hash value within sliding windows. This process creates a very valuable fingerprint of the text for later processing and comparison.

As a second step, Algorithm 6 compares two fingerprints previously calculated using Algorithm 5 by calculating the Jaccard similarity coefficient between their fingerprints. As with previous cases, a higher Jaccard similarity indicates greater similarity between the texts, making these

Algorithm 5 Winnow Algorithm

```
1: function WINNOW(text, k)
2:   k-grams  $\leftarrow$  emptylist
3:   for i  $\leftarrow$  0 to len(text) - k do
4:     append(k-grams, text[i : i + k])
5:   end for
6:   hashes  $\leftarrow$  emptylist
7:   for k - gram in k-grams do
8:     append(hashes, hash(k - gram))
9:   end for
10:  w  $\leftarrow$  10
11:  min_hashes  $\leftarrow$  emptylist
12:  for i  $\leftarrow$  0 to len(hashes) - w do
13:    min_hash  $\leftarrow$  min(hashes[i : i + w])
14:    append(min_hashes, min_hash)
15:  end for
16:  return min_hashes
17: end function
```

algorithms useful for measuring ODC similarity efficiently.

Algorithm 6 ODC Similarity Calculation using Winnow Hashing

```
1: function SIMILARITY(odc, odc2)
2:   k  $\leftarrow$  5
3:   fingerprint1  $\leftarrow$  set(Winnow(odc, k))
4:   fingerprint2  $\leftarrow$  set(Winnow(odc2, k))
5:   return  $\frac{\text{len}}{(\text{fingerprint1} \cap \text{fingerprint2})} \text{len}(\text{fingerprint1} \cup \text{fingerprint2})$ 
6: end function
```

4. Conclusion

This research provides an overview of various approaches to quantifying the similarity between ODCs, shedding light on the diverse methodologies and techniques employed in this domain. Our analysis has revealed multiple strategies for measuring catalog similarity, ranging from traditional methods based on comparing triples to advanced semantic-based methods for specific domain languages. Each approach brings its strengths and limitations, making it essential to choose the most appropriate method based on the particular goals and characteristics of the catalogs.

Identifying these varied approaches stresses the importance of selecting the most appropriate method based on the specific needs and characteristics of the catalogs in question. Catalogs spanning different domains may require tailored similarity metrics. Additionally, the scale of

catalogs, the granularity of data, and the available computational resources must all be considered when choosing the right approach.

In future work, it is necessary to keep researching in this direction since solutions are needed in this context. As open data initiatives grow in importance, catalog similarity measurement will play a vital role in facilitating data interoperability, discovery, and utilization across diverse domains and sectors.

References

- [1] Albertoni, R., Browning, D., Cox, S., González-Beltrán, A. N., Perego, A., & Winstanley, P. (2023). The W3C data catalog vocabulary, version 2: Rationale, design principles, and uptake. *CoRR*, *abs/2303.08883*. doi:10.48550/arXiv.2303.08883. arXiv:2303.08883.
- [2] Albertoni, R., & Isaac, A. (2021). Introducing the data quality vocabulary (DQV). *Semantic Web*, *12*, 81–97. doi:10.3233/SW-200382.
- [3] Bergroth, L., Hakonen, H., & Raita, T. (2000). A survey of longest common subsequence algorithms. In *Proceedings Seventh International Symposium on String Processing and Information Retrieval. SPIRE 2000* (pp. 39–48). IEEE.
- [4] Devlin, J., Chang, M., Lee, K., & Toutanova, K. (2019). BERT: pre-training of deep bidirectional transformers for language understanding. In J. Burstein, C. Doran, & T. Solorio (Eds.), *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers)* (pp. 4171–4186). Association for Computational Linguistics. doi:10.18653/v1/n19-1423.
- [5] Dzeroski, S., & Zenko, B. (2004). Is combining classifiers with stacking better than selecting the best one? *Mach. Learn.*, *54*, 255–273. doi:10.1023/B:MACH.0000015881.36452.6e.
- [6] English, T. M., & Gotesman, M. (1995). Stacked generalization and fitness ranking in evolutionary algorithms. In J. R. McDonnell, R. G. Reynolds, & D. B. Fogel (Eds.), *Proceedings of the Fourth Annual Conference on Evolutionary Programming, EP 1995, San Diego, CA, USA, March 1-3, 1995* (pp. 205–218). A Bradford Book, MIT Press. Cambridge, Massachusetts.
- [7] Goldberg, Y. (2017). Neural network methods for natural language processing. *Synthesis lectures on human language technologies*, *10*, 1–309.

- [8] Lakomaa, E., & Kallberg, J. (2013). Open data as a foundation for innovation: The enabling effect of free public sector information for entrepreneurs. *IEEE Access*, 1, 558–563. doi:10.1109/ACCESS.2013.2279164.
- [9] Maali, F., Erickson, J., & Archer, P. (2014). Data catalog vocabulary (dcat). w3c recommendation. *World Wide Web Consortium*, (pp. 29–126).
- [10] Martinez-Gil, J. (2019). Semantic similarity aggregators for very short textual expressions: a case study on landmarks and points of interest. *J. Intell. Inf. Syst.*, 53, 361–380. doi:10.1007/s10844-019-00561-0.
- [11] Martinez-Gil, J. (2022). A comprehensive review of stacking methods for semantic similarity measurement. *Machine Learning with Applications*, 10, 100423.
- [12] Martinez-Gil, J. (2023). A comparative study of ensemble techniques based on genetic programming: A case study in semantic similarity assessment. *Int. J. Softw. Eng. Knowl. Eng.*, 33, 289–312. doi:10.1142/S0218194022500772.
- [13] Martinez-Gil, J. (2023). Framework to automatically determine the quality of open data catalogs. *CoRR*, abs/2307.15464. URL: <https://doi.org/10.48550/arXiv.2307.15464>. doi:10.48550/arXiv.2307.15464. arXiv:2307.15464.
- [14] Martinez-Gil, J., & Chaves-Gonzalez, J. M. (2019). Automatic design of semantic similarity controllers based on fuzzy logics. *Expert Syst. Appl.*, 131, 45–59. doi:10.1016/j.eswa.2019.04.046.
- [15] Martinez-Gil, J., & Chaves-Gonzalez, J. M. (2021). Semantic similarity controllers: On the trade-off between accuracy and interpretability. *Knowl. Based Syst.*, 234, 107609. doi:10.1016/j.knosys.2021.107609.
- [16] Martinez-Gil, J., & Chaves-Gonzalez, J. M. (2022). Sustainable semantic similarity assessment. *Journal of Intelligent & Fuzzy Systems*, 43, 6163–6174. doi:10.3233/JIFS-220137.
- [17] Mikolov, T., Sutskever, I., Chen, K., Corrado, G. S., & Dean, J. (2013). Distributed representations of words and phrases and their compositionality. In *Advances in Neural Information Processing Systems 26: 27th Annual Conference on Neural Information Processing Systems 2013. Proceedings of a meeting held December 5-8, 2013, Lake Tahoe, Nevada, United States.* (pp. 3111–3119).

- [18] Navigli, R., & Martelli, F. (2019). An overview of word and sense similarity. *Nat. Lang. Eng.*, 25, 693–714. doi:10.1017/S1351324919000305.
- [19] Paoletti, A. L., Martinez-Gil, J., & Schewe, K. (2016). Top-k matching queries for filter-based profile matching in knowledge bases. In S. Hartmann, & H. Ma (Eds.), *Database and Expert Systems Applications - 27th International Conference, DEXA 2016, Porto, Portugal, September 5-8, 2016, Proceedings, Part II* (pp. 295–302). Springer volume 9828 of *Lecture Notes in Computer Science*. URL: https://doi.org/10.1007/978-3-319-44406-2_23. doi:10.1007/978-3-319-44406-2_23.
- [20] Schleimer, S., Wilkerson, D. S., & Aiken, A. (2003). Winnowing: local algorithms for document fingerprinting. In *Proceedings of the 2003 ACM SIGMOD international conference on Management of data* (pp. 76–85).
- [21] Skoda, P., Bernhauer, D., Necaský, M., Klímek, J., & Skopal, T. (2020). Evaluation framework for search methods focused on dataset findability in open data catalogs. In M. Indrawan-Santiago, E. Pardede, I. L. Salvadori, M. Steinbauer, I. Khalil, & G. Kotsis (Eds.), *iWAS '20: The 22nd International Conference on Information Integration and Web-based Applications & Services, Virtual Event / Chiang Mai, Thailand, November 30 - December 2, 2020* (pp. 200–209). ACM. doi:10.1145/3428757.3429973.
- [22] Subramaniam, P., Ma, Y., Li, C., Mohanty, I., & Fernandez, R. C. (2021). Comprehensive and comprehensible data catalogs: The what, who, where, when, why, and how of metadata management. *CoRR*, abs/2103.07532. arXiv:2103.07532.