92

93

94

95

96

97

98

99

100

101

102

103

104

105

106

107

108

109

110

111

112

113

114

115

116

59

60

Author Name Disabiguation using Markov Chain Monte Carlo

Dhruv Singh Chandel dhruv.singh.chandel@fit.fraunhofer.de Fraunhofer FIT Sankt Augustin, Germany

Zeyd Boukhers zeyd.boukhers@fit.fraunhofer.de Fraunhofer FIT Sankt Augustin, Germany

Abstract

1

2

3

4

5

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

26

27

28

29

30

31

32

33

34

35

36

37

38

39

40

41

42

43

44

45

46

47

48

49

50

51

52

53

54

This paper presents a novel approach to the Incorrect Name Detection (IND) task as part of the KDD Cup 2024 Open Academic Graph Challenge (OAG-Challenge). We propose Author Name Disambiguation using Markov Chain Monte Carlo (AND-MCMC) algorithm to identify incorrectly assigned papers within author profiles in the WhoIsWho dataset [9]. Our method constructs graph structures or "graphlets" for each author and employs an iterative refinement process that prioritizes split actions over merge actions. The approach aims to effectively separate anomalous papers from those correctly attributed to the predominant author. Leveraging the dataset's structure that includes correctly and incorrectly assigned publications, the algorithm employed in this work processes one author's file at a time. We evaluate our method using a weighted Area Under the Receiver Operating Characteristic Curve (AUC) metric^[2], which accounts for varying error distributions across authors. This work contributes to academic graph mining by addressing the challenges associated with detecting incorrect paper attributions in large-scale scholarly databases.

Keywords

Incorrect Name Detection using Author Name Disambiguation(IND-AND), Author Name Disambiguation, Graph Clustering, AND-MCMC Algorithm, Anomaly Detection, WhoIsWho Dataset, Weighted AUC, Graphlet Analysis, Scholarly Databases, Iterative Graph Refinement, OAG-Challenge, KDD Cup 2024

ACM Reference Format:

Dhruv Singh Chandel, Nagaraj Bahubali Asundi, Zeyd Boukhers, and Sulayman K. Sowe. 2024. Author Name Disabiguation using Markov Chain Monte Carlo. In *Proceedings of (KDD Cup 2024)*. ACM, New York, NY, USA, 4 pages. https://doi.org/10.1145/nnnnnnnnnn Nagaraj Bahubali Asundi nagaraj.bahubali.asundi@fit.fraunhofer.de Fraunhofer FIT Sankt Augustin, Germany

> Sulayman K. Sowe sowe@bdis.rwth-aachen.de RWTH Aachen University Aachen, Germany

1 Introduction

The exponential growth of academic literature in recent years has led to an increasing need for accurate and efficient methods of organizing scholarly information[12]. One of the critical challenges in this domain is the correct attribution of academic papers to their respective authors, a task complicated by factors such as name ambiguity, inconsistent name spellings, and errors in database entries. The IND task, a crucial component of academic graph mining [3, 9], aims to address this challenge by identifying and rectifying instances where papers have been erroneously attributed to authors.

In the context of the OAG-Challenge at the KDD Cup 2024¹, this paper presents a novel approach to the IND task, leveraging the WhoIsWho dataset provided by AMiner.cn². Our work builds upon the foundation laid by previous research in author name disambiguation and graph-based approaches to scholarly data analysis [11]. We propose the AND-MCMC algorithm, specifically tailored to detect incorrectly assigned papers within individual author profiles.

Our approach innovates by introducing an iterative refinement process that prioritizes the separation of anomalous papers from those correctly attributed to the predominant author. This is achieved through the construction and manipulation of graph structures, or "graphlets," for each author profile. *By favoring split actions over merge actions in our algorithm*, we aim to effectively isolate incorrectly assigned papers, even in cases where they represent a small fraction of an author's overall publication record.

In the following sections, we provide a comprehensive overview of our methodology, including the data preprocessing steps, the AND-MCMC algorithm, and our evaluation framework. We then present and discuss our results, achieving a weighted AUC score of 0.499933, which indicates a moderate ability to distinguish between correct and incorrect paper assignments. Finally, we conclude with insights into the strengths and limitations of our approach, as well as potential avenues for future research in this critical area of academic graph mining.

This work contributes to the broader field of scholarly data analysis by addressing the specific challenge of incorrect paper attribution, a problem that has significant implications for the accuracy and reliability of academic databases and bibliometric studies. Our findings not only offer a novel solution to the Incorrect Name Detection

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

⁵⁵ KDD Cup 2024, March 20th, 2024 - June 14th, 2024, Tsinghua University, China

^{56 © 2024} Copyright held by the owner/author(s). Publication rights licensed to ACM. ACM ISBN 978-x-xxxx-x/YY/MM

⁵⁷ https://doi.org/10.1145/nnnnnnnnnn

⁵⁸

¹https://www.biendata.xyz/kdd2024/

²https://www.aminer.cn/

(IND) task but also provide valuable insights into the complexities of author name disambiguation in large-scale academic datasets.

2 Background and Related Work

117

118

119

120

121

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

140

154

155

156

157

158

159

160

161

162

163

164

165

166

174

Author Name Disambiguation (AND) is a critical challenge in the field of academic graph mining, particularly as the volume of scholarly literature continues to grow exponentially[7] [1]. The task involves correctly attributing academic papers to their respective authors, a process complicated by factors such as name ambiguity, inconsistent name spellings, and errors in database entries. Accurate author name disambiguation is crucial for ensuring the reliability of bibliometric analyses, researcher evaluations, and literature discovery systems[10].

The AND task has been approached through various methodologies, broadly categorized into supervised, unsupervised, and hybrid approaches[4]. While supervised methods can achieve high accuracy, they often face challenges related to scalability and the need for substantial manually labeled data[6]. Unsupervised methods, 139 on the other hand, offer better scalability but may sacrifice some accuracy compared to supervised approaches[8].

In recent years, graph-based approaches have gained significant 141 142 traction in the field of AND. These methods leverage the complex relationships between authors, publications, and other metadata to 143 improve disambiguation accuracy[1]. The use of co-authorship net-144 145 works, in particular, has shown promise in distinguishing between authors with similar names but different collaboration patterns. The 146 incorporation of advanced natural language processing techniques 147 148 has also been explored in the context of AND. Word embeddings 149 and topic modeling approaches have been employed to capture 150 semantic similarities between paper titles and abstracts, providing 151 valuable information for disambiguation purposes[5]. The use of 152 these techniques allows for a more nuanced understanding of an author's research interests and how they evolve over time. 153

The concept of publication patterns over time has been identified as a valuable feature for disambiguation, especially in cases where authors have similar research interests but different career trajectories[6]. By modeling how an author's publications are distributed across topics and years, it becomes possible to distinguish between authors who may appear similar based on other metrics.

Despite these advancements, the AND problem remains challenging due to the dynamic nature of academic publishing, the increasing interdisciplinarity of research, and the global diversity of author names[8][1]. There is a continued need for methods that can effectively integrate multiple sources of information, handle large-scale datasets, and adapt to the evolving landscape of scholarly communication.

The AND-MCMC method proposed in this work aims to address 167 168 these challenges by combining multiple features within a Markov Chain Monte Carlo framework. By leveraging domain similarity, 169 co-authorship networks, publication patterns, and affiliation infor-170 mation, the method seeks to provide a comprehensive approach to 171 172 author name disambiguation that can handle the complexities of 173 modern academic databases.

Methodology 3

The data used for the IND task is collected from AMiner.cn and provided in the WhoIsWho dataset. The proposed approach transforms the dataset into a graph with nodes representing various features, and iteratively merges or splits these groups of nodes until clusters of homogeneous papers are generated. The entire approach used is shown in Figure 1.



Figure 1: Architecture of the proposed approach.

For this task, we propose the AND-MCMC algorithm to detect the incorrectly assigned papers within each author's file. The key steps of the approach are as follows:

- Data Loading and Transformation: The process begins by loading the author's data from the provided dataset (e.g., train_author.json), which includes the author's ID, name, and paper IDs (both correctly and incorrectly assigned). This data is then transformed into an atomic name file format, creating a structured representation that facilitates subsequent graph construction and manipulation.
- Graph Construction and Paper Sampling: A graph structure, or "graphlet," is created with all paper nodes initially connected to a single atomic node representing the author. From this graphlet, a paper node is randomly selected for potential merge or split operations in the subsequent step.
- Action Selection and Prioritization: The algorithm selects an action to perform on the sampled paper node, either a merge or split operation. In the first iteration, the split action is favoured, as there are no other graphlets to merge with. Furthermore, throughout the iterative process, the split action is generally prioritised over the merge action to facilitate the separation of incorrectly assigned papers. an exaple of the split and merge actions can be seen in Figure 2.

175

176

Author Name Disabiguation using Markov Chain Monte Carlo



Figure 2: Illustrations of merge and split actions.

- Action Execution and Iterative Refinement: The selected merge or split action is applied to the sampled paper node, either combining it with an existing graphlet or creating a new one. This process is repeated iteratively until a stable graph topology is achieved, progressively refining the graphlets and separating incorrectly assigned papers from the predominant author's works.
- Anomaly Identification: After convergence, the graphlets with the highest number of paper nodes are considered to represent the correct author's papers. Conversely, the remaining graphlets with fewer paper nodes are flagged as anomalies or incorrectly assigned papers.

Acceptance Criterion

The algorithm evaluates proposed actions, such as merging or splitting, by first calculating an acceptance criterion. If this criterion exceeds a certain threshold, the proposed action is performed. The acceptance criterion is based on four factors.

 Domain Similarity (α): We assess the action by comparing the domain similarity of paper titles in the current graph state with the future state after the possible merge or split. For merging, a *target graphlet* is first sampled along with a *merging graphlet*. The decision to merge these graphlets is made based on the ratio of the similarity between the paper groups of the *target graphlet* and the *merging graphlet* to the similarity of the subgroups of papers within the *target graphlet*.

Let $\mathcal{P}(x)$ be a function that returns all papers in a graphlet x, $\operatorname{Avg}(P)$ a function that returns the average of all embeddings of the paper set P, and $\operatorname{Sim}(emb_1, emb_2)$ a function that returns the cosine similarity between the embedding vectors emb_1 and emb_2 . Let g and h represent the target graphlet and the merging graphlet, respectively. Let g' and g'' be the interim splits of g such that $\mathcal{P}(g') \cup \mathcal{P}(g'') = \mathcal{P}(g)$. The ratio of domain similarity for the merge (α_{merge}) is given by:

$$\alpha_{merge} = \frac{\alpha^{(t+1)}}{\alpha^{(t)}} = \frac{\operatorname{Sim}(\operatorname{Avg}(\mathcal{P}(g)), \operatorname{Avg}(\mathcal{P}(h)))}{\operatorname{Sim}(\operatorname{Avg}(\mathcal{P}(g')), \operatorname{Avg}(\mathcal{P}(g'')))}$$
(1)

KDD Cup 2024, March 20th, 2024 - June 14th, 2024, Tsinghua University, China

The interim splits g' and g'' are obtained by applying kmeans clustering over the embeddings of all papers in g. The intuition behind computing α_{merge} is that we want to compare the similarity of the papers in graph state t with that in state t + 1. Here, state t represents the status quo or the current state of the graph, and state t + 1 represents the future state in which the target and merging graphlets are combined. Based on the ratio obtained, we decide whether it is worth merging the graphlets or not.

Unlike merging, we only need a target graphlet for splitting. Therefore, a target graphlet g is selected first. Then, a paper \hat{p} is selected from all papers within g such that \hat{p} is distant from the rest of the papers in terms of the domain. Now 3 interim splits g', g'', and g''' are obtained by applying k-means clustering to the papers of g such that $\mathcal{P}(g') \cup \mathcal{P}(g'') \cup \mathcal{P}(g''') = \mathcal{P}(g)$. Here g''' contains the paper \hat{p} . The ratio of domain similarity for the split (α_{split}) is given by:

$$\alpha_{split} = \frac{\alpha^{(t+1)}}{\alpha^{(t)}} = \frac{\operatorname{Sim}(\operatorname{Avg}(\mathcal{P}(g')), \operatorname{Avg}(\mathcal{P}(g'')))}{\operatorname{Sim}(\operatorname{Avg}(\mathcal{P}(g') \cup \mathcal{P}(g'')), \operatorname{Avg}(\mathcal{P}(g''')))}$$
(2)

- (2) Co-authorship Overlap (β): Authors who have already published together are more likely to continue their collaboration and publish additional papers. Therefore, co-author networks are well-suited for author identification. In this work, the Jaccard index is used to measure the similarity between two graphlets.
- (3) **Publication Pattern** (γ): The publication pattern models the domain of an author's papers over time. This is achieved by feeding the paper titles into a Gaussian Latent Dirichlet Allocation (LDA) model to extract latent topics. Once we have the topical distribution of papers, we determine the topical distribution over the duration of the author's publication period. For this, all the papers in a graphlet are collected, and the publication patterns are determined per topic. We model the distribution of the *i*th topic over the years independently of the *j*th topic, given there are *n* latent topics with $i \le n$, $j \le n$, and $i \ne j$. Determining the distribution per topic allows easier comparison of two sets of papers from independent graphlets.
- (4) Affiliation Overlap (κ): While co-authorship overlap can help determine if both sets of co-authors belong to the same target author, there is still a possibility of ambiguity, as different authors may share the same co-author names. To address this, we also consider the affiliations of these co-authors. By checking for overlap in affiliations, we can better resolve co-author ambiguity.

The algorithm begins by sampling a target graphlet g. Then an appropriate action for g is selected. If the action is *merge*, another graphlet h is selected that could be a potential candidate to be merged with g. Also, two interim splits of g are obtained, namely g' and g''. The interim splits are just imaginary splits of g and have nothing to do with the *split* operation. The decision to merge g with h is based on the acceptance criterion A_{merge} , which is defined as follows:

349 350 351

352

353

354

355

356

357

358

359

360

361

362

363

364

365

366

367

368

369

370

371

372

373

374

375

376

377

378

379

380

381

382

383

384

385

386

387

388

389

390

391

392

393

394

395

396

397

398

399

400

401

402

403

404

405

406

 $A_{merge} = \alpha_{merge} \cdot \beta_{merge} \cdot \gamma_{merge} \cdot \kappa_{merge}$

(3)

Once A_{merge} is calculated a value u is sampled from a uniform distribution $\mathcal{U}(0, 1)$. The proposal to merge *g* with *h* is accepted if $A_{merge} > u$. When the action is *split*, a paper \hat{p} is selected within g that is farthest from the remaining papers with respect to the domain. Also, three interim splits of g are obtained, namely g', g'', and $g^{\prime\prime\prime}$. Here $g^{\prime\prime\prime}$ contains the paper \hat{p} . The decision to separate g''' from g is based on the acceptance criterion A_{split} which is calculated similar to A_{merge} .

By favouring the split action and iteratively refining the graph structure, the proposed approach aims to separate the incorrectly assigned papers from the predominant author's papers within each file. The graphlets with the highest number of papers are considered to represent the correct author, while the remaining graphlets are flagged as anomalies or incorrect assignments.

This approach leverages the structure of the WhoIsWho dataset, where each author's file contains both correctly assigned papers (normal data) and incorrectly assigned papers (outliers). By processing one author's file at a time and favouring the split action, the algorithm can effectively identify the anomalous papers that do not belong to the predominant author.

4 Evaluation Metrics

The performance of the proposed approach is evaluated using the Area Under the Receiver Operating Characteristic (ROC) Curve (AUC) metric

Here's a condensed, academic-style version of the information:

Author-wise AUC Calculation 4.1

The proposed approach is applied to each author's dataset, comprising both correctly and incorrectly assigned papers. Resultant graphlets are size-ranked, with the largest presumed to represent the author's legitimate works. Using ground truth labels, true positive and false positive rates are computed across various graphlet size thresholds. The ROC curve is then plotted, and the area under this curve (AUC) is calculated, yielding the author-specific AUC score.

4.2 Weighted AUC Calculation

To account for the varying importance of different authors in the dataset, a weighted AUC score is calculated as follows:

$$Weight_i = \frac{\# Errors of the Author_i}{\# Total Errors}$$
(4)

The weighted AUC is then calculated as the sum of the AUC scores for each author, multiplied by their respective weights.

Weighted AUC =
$$\sum_{i=1}^{M} (AUC_i \times weight_i)$$
 (5)

Where *M* is the total number of authors in the dataset.

The weighted AUC score accounts for each author's contribution to the dataset's total errors, emphasizing authors with higher incorrectly assigned paper counts. The proposed approach achieved a weighted AUC of 0.499933, indicating moderate success in distinguishing correct from incorrect assignments while highlighting potential for improvement. This metric offers a comprehensive evaluation across all authors, considering varied error distributions and appropriately weighting authors with higher error rates in the overall assessment.

Conclusion 5

This paper introduces the AND-MCMC algorithm, a novel approach to the Incorrect Name Detection (IND) task within the KDD Cup 2024 Open Academic Graph Challenge. Our method addresses the complexity of distinguishing correctly and incorrectly attributed papers in large-scale scholarly databases, achieving a weighted AUC score of 0.499933. While demonstrating potential, this result also indicates opportunities for further refinement. Our research contributes to creating reliable public benchmarks for academic graph mining, a critical gap in the field. The methodology serves as a foundation for future developments in author name disambiguation systems and showcases the potential of graph-based methods in scholarly data analysis. This work aims to stimulate further research in academic graph mining, addressing the limitation of suitable public benchmarks. While our approach shows promise, there remains significant scope for improvement. We anticipate this contribution will foster continued innovation in managing and analyzing scholarly data at scale, benefiting the broader KDD Cup 2024 audience and the academic community.

References

- [1] Michele De Bonis, F. Falchi, and Paolo Manghi. 2023. Graph-based methods for Author Name Disambiguation: a survey. PeerJ Computer Science 9 (2023). https://api.semanticscholar.org/CorpusID:261793272
- [2] Andrew P. Bradley. 1997. The use of the area under the ROC curve in the evaluation of machine learning algorithms. Pattern Recognition 30, 7 (1997), 1145-1159. https://doi.org/10.1016/S0031-3203(96)00142-2
- [3] Başak Buluz and Burcu Yilmaz. 2017. Graph mining approach for modeling academic success. In 2017 25th Signal Processing and Communications Applications Conference (SIU). 1-4. https://doi.org/10.1109/SIU.2017.7960621
- [4] Ijaz Hussain and Sohail Asghar. 2017. A survey of author name disambiguation techniques: 2010-2016. The Knowledge Engineering Review 32 (2017). https: //api.semanticscholar.org/CorpusID:52175217
- Ayesha Manzoor, Sohail Asghar, and Tehmina Amjad. 2022. Toward a New [5] Paradigm for Author Name Disambiguation. IEEE Access 10 (2022), 76055-76068. https://api.semanticscholar.org/CorpusID:250518707
- [6] Andreas Rehs. 2021. A supervised machine learning approach to author disambiguation in the Web of Science, J. Informetrics 15 (2021), 101166. https:// //api.semanticscholar.org/CorpusID:236236516
- Neil R. Smalheiser and Vetle I. Torvik, 2009. Author name disambiguation. Annu. [7] Rev. Inf. Sci. Technol. 43 (2009), 1-43. https://api.semanticscholar.org/CorpusID: 205418775
- Alexander Tekles and Lutz Bornmann. 2019. Author name disambiguation of bib-[8] liometric data: A comparison of several unsupervised approaches1. *Quantitative* Science Studies 1 (2019), 1510-1528. https://api.semanticscholar.org/CorpusID: 139102212
- Fanjin Zhang, Shijie Shi, Yifan Zhu, Bo Chen, Yukuo Cen, Jifan Yu, Yelin Chen, [9] Lulu Wang, Qingfei Zhao, Yuqing Cheng, Tianyi Han, Yuwei An, Dan Zhang, Weng Lam Tam, Kun Cao, Yunhe Pang, Xinyu Guan, Huihui Yuan, Jian Song, Xiaoyan Li, Yuxiao Dong, and Jie Tang. 2024. OAG-Bench: A Human-Curated Benchmark for Academic Graph Mining. In Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining.
- Wenjing Zhang, Zhongmin Yan, and Yongqing Zheng. 2019. Author Name Dis-[10] ambiguation Using Graph Node Embedding Method. 2019 IEEE 23rd International Conference on Computer Supported Cooperative Work in Design (CSCWD) (2019), 410-415. https://api.semanticscholar.org/CorpusID:199509475
- [11] Yutao Zhang, Fanjin Zhang, Peiran Yao, and Jie Tang. 2018. Name Disambiguation in AMiner: Clustering, Maintenance, and Human in the Loop. Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining (2018). https://api.semanticscholar.org/CorpusID:207579405
- [12] Ivan Zupic and Tomaž Čater. 2015. Bibliometric methods in management and organization. Organizational research methods 18, 3 (2015), 429-472

451

452

453

454

455

456

457

458

459

460

461

462

463

464

407

408

409

410