

A DE-DUPLICATION TOOL: IMPLEMENTATION OF DUPLICATE DETECTION ALGORITHMS

J.Shana*

Dr.T.Venkatachalam**

ABSTRACT

Data quality is an important issue in any data store especially when it is used for analysis. Poor quality of data would drastically affect the decision making process. Cleansing is an activity performed before data is loaded into the warehouse to enhance the quality and consistency of the data. This paper addresses one of the key data quality problem namely detection of duplicates. A duplicate detection tool is built that implements two algorithms namely Sorted Neighborhood method and Token based cleansing method. The tool was tested using student dataset of 118 records and showed an accuracy of 88.09%. This paper is a primitive analysis of how duplicates can be detected in a dataset.

Keywords- data quality, data cleansing, duplicate detection, data warehouse

* Department of MCA, Coimbatore Institute of Technology, Tamilnadu, India.

** Department of Physics, Coimbatore Institute of Technology, Tamilnadu, India.

INTRODUCTION

Explosion in the amount of data entered is leading to increased corporate attention on improving the quality and accuracy of business data. The data given as input for the data mining process should be of high quality in order for the results to be accurate and reliable. Before it is stored in warehouse or mined the data goes through a process called data cleansing [6][9]. Poor data quality can have a severe impact on the overall effectiveness of information systems. Data cleansing tools help in improving data consistency and accuracy. A TDWI survey published in a research report on data quality showed that customer and product data are the two main types of business entity susceptible to quality problems. Customer data involves variations in names, titles, phone numbers and addresses [10]. Detecting and eliminating duplicated data is an important problem in the broader area of data cleansing and data quality. A duplicate value is the most serious kind of dirty data which leads to wastage of various resources. Many times the same logical real world entity has multiple representations in a relation due to data entry errors, varying conventions and other reasons. Hence a significant amount of time and money is spent on the task of detecting and eliminating duplicates.

1. EXPERIMENTAL EVALUATION

We developed a prototype of a de-duplication tool that implements the two algorithms, Sorted Neighborhood Method and Token Based cleansing algorithm.[5] The dataset used is a student dataset consisting of 118 records. The dataset consist of fields namely Roll Number, Name, Gender, Date of Birth, Father's name, Address, Phone Number.

2.1 Preprocessing

Preprocessing is one of the important steps in the design of data warehouse or prior to data mining. Preprocessing is done to standardize the data before implementing the algorithm to remove errors such as mistyped letters or abbreviations. For example, DOB formats such as 19-Dec-1978 or 19-12-1978 are converted into one standard format 19-12-1978. Abbreviations such as rd, st are converted into road and street respectively. We introduce two types of errors common in data warehouse namely equivalence errors and spelling errors [3]. Exact duplicate tuples are those for which all the fields are the same. If there is a mismatch in any one field then it is inexact duplicate. It is due to typographic errors or format mismatch in data.

2.2 Evaluation metrics

Precision and Recall metrics are used to evaluate the duplicate elimination algorithms [2]. Precision is ratio of identified duplicates to the original duplicates. Recall is the fraction of pairs of tuples an algorithm returns which are truly duplicates.

$$\text{Recall} = \frac{\text{No. of duplicates identified}}{\text{Total no. of duplicates present}} * 100$$

(1)

2. ALGORITHMS USED IN THIS STUDY

3.1 Sorted Neighborhood Method

This method involves the following phases.

Create Key: A key for each record in the list is computed by extracting relevant fields or portions of fields.

Sort Data: The records are sorted by using the key found previously. By sorting, duplicate records will come closer in the list. The effectiveness of this method highly depends upon the comparison key that is selected to sort the records.

Merge: A fixed size window is moved through the sequential list of records in order to limit the comparisons for matching records to those in the window. If the size of the window is 'w' records then every new record that enters that window is compared with the previous 'w-1' records to find matching records. The first record in the window slides out of it.

3.2 Token Based data cleansing method

This algorithm is implemented as follows:

Select Tokens: Select and rank 2 or 3 fields based on their record identifying abilities. Extract smart token for each selected field as follows. Form numeric, alphabetic or alphanumeric tokens after removing stop words and unimportant characters [4]. Table 1 gives the category of tokens used.

Table 1: Tokens for the dataset

Tokens
Numeric
Alphanumeric
Alphabetic

Sort Data: The records in the database are sorted by using the token found in the previous step. To detect the duplicates two ratios, Similarity Match Count (SMC) and Similarity Match Ratio (SMR) are calculated. They are calculated using the equation (2) and (3) as given below.

$$SMC = \frac{\text{no of matching token fields}}{\text{no of token fields used}}$$

(2)

$$SMR = 2 * \frac{\text{no of common characters in the 2 tokens}}{\text{Total no of characters in the 2 tokens}}$$

(3)

SMC and SMR use the token fields and identify all pairs of records as duplicates using Table 2.

Table 2: SMC values for identifying the duplicates

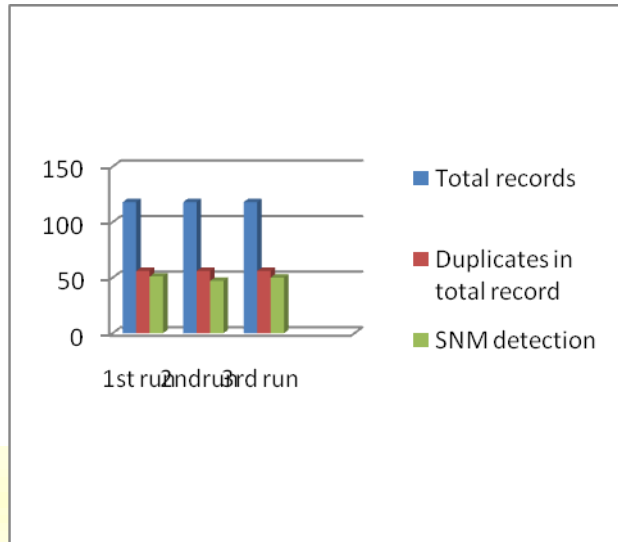
Match type	SMC Values
Perfect match	1.0
Near perfect match	Between 0.67 and 0.99
No Match	Less than 0.33
May be match	Between 0.33 and 0.66

Calculate SMR for may be match and the two tokens are a match if their SMR is greater or equal to 0.67.

3. EXPERIMENTAL RESULTS

4.1 SNM Method

Figure 1 shows the result of three runs made to detect duplicates using SNM Method.



Recall values for different runs of records in Sorted Neighborhood Method.

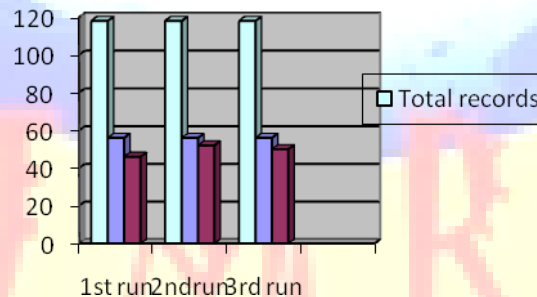
Recall value for 1st run = **91.07 %**

Recall value for 2nd run = **83.92 %**

False Positive error for 3rd run = 89.2%

Average recall value for SNM is **87.5%**

4.2 Token Based Data Cleansing Method



Recall values for different runs of records in Token Based data cleansing Method.

Recall value for 1st run= **82.14 %**

Recall value for 2nd run= **92.86 %**

Recall value for 3rd run= **89.28 %**

Average recall value for Token based data cleansing method is 88.09

Figure 3: Screen shot for SNM Method

Original Records With Duplicates

S.No	Roll No.	Name	Gender	DOB	Father Name
1	07AC01	Aashir F	F	16/01/95	D. Pathanqul
2	07AC01	Aashir F	F	16/01/95	D. Pathanqul
3	07AC02	Aneesh varshu R	F	20/05/90	Pengalath R
4	07AC02	Aneesh varshu R	F	20/05/90	Pengalath R
5	07AC02	Aneesh varshu R	F	20/05/90	Pengalath R
6	07AC04	Ash P	F	23/02/90	Palani S
7	07AC04	Ash P	F	23/02/90	Palani S
8	07AC05	Azad Kalia R	M	29/09/89	Pandeyan C
9	07AC05	Azad Kalia R	M	29/09/89	Palani S
10	07AC05	Azad Kalia R	M	29/09/89	Pandeyan C
11	07AC05	Azusa M	F	15/02/90	Mangappa R
12	07AC05	Azusa M	F	15/02/90	Mangappa R
13	07AC07	Deepthi S	F	13/07/88	Shankupill C D
14	07AC07	Deepthi S	F	13/07/88	Shankupill C D
15	07AC08	Deepika R	F	13/10/89	Pandeyan N
16	07AC08	Deepika R	F	13/10/89	Pandeyan N
17	07AC08	Deepika R	F	13/10/89	Pandeyan N
18	07AC08	Deepika R	F	13/10/89	Pandeyan N
19	07AC08	Deepika R	F	13/10/89	Pandeyan N
20	07AC08	Deepika R	F	13/10/89	Pandeyan N
21	07AC08	Deepika R	F	13/10/89	Pandeyan N
22	07AC08	Deepika R	F	13/10/89	Pandeyan N
23	07AC10	Deepika S	F	10/10/88	Kulasek T S

After Duplicate Elimination

S.No	Roll No.	Name	Gender	DOB	Father Name
1	07AC01	Aashir F	F	16/01/95	D. Pathanqul
2	07AC02	Aneesh varshu R	F	20/05/90	Pengalath R
3	07AC04	Ash P	F	23/02/90	Palani S
4	07AC05	Azad Kalia R	M	29/09/89	Palani S
5	07AC05	Azusa M	F	15/02/90	Mangappa R
6	07AC07	Deepthi S	F	13/07/88	Shankupill C D
7	07AC08	Deepika R	F	13/10/89	Pandeyan N
8	07AC08	Deepika R	F	13/10/89	Pandeyan N
9	07AC10	M Dhanya P	M	23/10/89	Palanathy N
10	07AC10	ANBUSELVAN S	M	20/03/89	SATHIHEL A
11	07AC10	SALAMURIGANAN W	M	04/02/89	THANGAEL N
12	07AC10	CAROLINE J	F	01/01/89	JANARAJ
13	07AC11	Geetha S	M	06/09/87	Sudhakar M P
14	07AC12	Sudhith F	F	19/10/89	Dhayan N
15	07AC14	SIVATHI S	F	08/12/89	SELVARAJAN
16	07AC15	Malathi V	F	05/12/89	Vidyalay V
17	07AC17	Indukathi P	F	27/09/89	Pandeyan M K
18	07AC19	Kalyani F	F	04/09/89	Vidyanan M
19	07AC20	Karthi C	F	21/01/89	Chidambaram S
20	07AC21	Karthi P	F	20/07/89	Palanathy M
21	07AC22	Mahabharth M	F	13/05/89	Mangayandhar R
22	07AC22	Mahesh R	F	26/05/91	Pagalaraj S S
23	07AC28	Nikhilata A	M	19/01/87	Kulasek

Total Records: 118 Total Records: 48

HOME EXIT

Figure 4: Screen shot for Token Based Method for single record.

DEDUCTION OF DUPLICATES FOR EACH RECORD USING TOKEN BASED ALGORITHM

S.No	Roll No.	Name	DOB	Address	SMC	Match	SMR	Match	NO
8	07AC05	Anud Kulkar R	29/08/89	banur Coimbatore-23		1 Perfect Match			1
9	07AC05	Anud Kulkar R	29/08/89	banur Coimbatore-23		1 Perfect Match			2
10	07AC05	Anud Kulkar R	29/08/89	banur Coimbatore-23		1 Perfect Match			3
9	07AC05	Anud Kulkar R	29/08/89	banur Coimbatore-23		1 Perfect Match			2
8	07AC05	Anud Kulkar R	29/08/89	banur Coimbatore-23		1 Perfect Match			1
10	07AC05	Anud Kulkar R	29/08/89	banur Coimbatore-23		1 Perfect Match			3
10	07AC05	Anud Kulkar R	29/08/89	banur Coimbatore-23		1 Perfect Match			3
9	07AC05	Anud Kulkar R	29/08/89	banur Coimbatore-23		1 Perfect Match			2
8	07AC05	Anud Kulkar R	29/08/89	banur Coimbatore-23		1 Perfect Match			1
11	07AC06	Anuna M	15/02/90	1st Singanallur Cba 5		1 Perfect Match			4
12	07AC06	Anuna M	15/02/90	1st Singanallur Cba 5		1 Perfect Match			5
12	07AC06	Anuna M	15/02/90	1st Singanallur Cba 5		1 Perfect Match			5
11	07AC06	Anuna M	15/02/90	1st Singanallur Cba 5		1 Perfect Match			4
63	07AC28	Nikhilata A	28/08/90	anethapuram Cba-45		1 Perfect Match			6
64	07AC28	Nikhilata A	28/08/90	sdagaM road CBE-25		0.67 Near Perfect Match			7
65	07AC28	Nikhilata A	28/08/90	sdagaM road CBE-25		0.67 Near Perfect Match			8

Duplication for each Record

NEXT >>

Figure 5: Screen shot for Token Based Method

Duplicate Records detected using Token Based Method

S.No	Roll No	Name	Gender	DOB	Father Name	Occupation	Address	Mobile
9	07AC26	Arun Kumar.R	Male	29/08/89	Reg. C	Section Supervisor	h 5/Postanur,Coimbatore-23	9789411466
10	07AC26	Arun Kumar.R	male	29/08/89	Rajendran.C	Section Supervisor	h 5/Postanur,Coimbatore-23	9789411466
12	07AC26	Aruna.M	F	15/02/90	Manojagan.R	Milkworker	the Mill rd Singanailur, Che 5	9894208223
65	07AC28	Nalata.A	F	28/08/90	ArunaM.R	business	igar thadegam road, Che 25	9003511620
1	07AC01	Aarathi.R	F	16/03/89	Doctor.RajMaraju.A	Business	J St,Civil Aerodrome, Che-14	9600402389
4	07AC03	Aneeth varshini.R	Female	20/06/90	Rajengrathan	Service/Minister OF T	No 5 Velavan Nagar,Bharat	9790507904
5	07AC03	Aneeth varshini.R	F	20/06/90	Rajengrathan	Service/Minister OF T	No 5 Velavan Nagar,Bharat	9790507904
25	07AC10	Meter Dhanapal.P	Male	23/12/89	Palanisamy.K	Lete	Uppilpalayam Post,Che-15	97906-21722
17	07AC08	Deepika.R	Female	13/10/89	Rajendran.N	Business	No 11, Jothi Nagar,3rd St,Up	9380933736
18	07AC08	Deepika.R	Female	13/10/89	Rajendran.N	Business	No 11, Jothi Nagar,3rd St,Up	9380933736
19	07AC08	Deepika.R	F	13/10/89	Rajendran.N	Business	No 11, Jothi Nagar,3rd St,Up	9380933736
20	07AC08	Deepika.R	f	13/10/89	Rajendran.N	Business	No 11, Jothi Nagar,3rd St,Up	9380933736
21	07AC08	Deepika.R	f	13/10/89	Rajendran.N	Business	No 11, Jothi Nagar,3rd St,Up	9380933736
22	07AC08	Deepika.R	f	13/10/89	Rajendran.N	Business	No 11, Jothi Nagar,3rd St,Up	9380933736
15	07AC08	Deepika.R	Female	13/10/89	Rajendran.N	Business	No 11, Jothi Nagar,3rd St,Up	9380933736
13	07AC07	Deepthi.S	F	13/07/80	ShanMugaM.C.D	Business	agar,Edyerpalayam,Che 25	9852207340

Total No Of Records

No of Perfect match identified

4. CONCLUSION

The data de-duplication tool is implemented successfully using Sorted Neighborhood Method and Token Based Data Cleansing Method. Both the algorithms were tested with the same dataset and the result analyzed. The tool was successfully tested for various runs. The analysis was done based on recall value and false positive error. SNM could detect and eliminate duplicates but it also identified false positives. Token based could detect duplicates with near match and may be match also. Further processing for near perfect match records needs to be done on token based method to detect the duplicates very efficiently.

References

- [1] Ramzi A. Haraty, Ralph Varjabedian Lebanese, ADD: Arabic Duplicate Detector- A Duplicate Detector ,IEEE June 2009
- [2] Surajit Chaudhuri, Venkatesh Ganti, Rajeev Motwani, Robust Identification of Fuzzy Duplicates, IEEE, June 2008.
- [3] M.Bilenko, R.Mooney ,On Evaluation and Training-set Construction for Duplicate Detection , Proceedings of the ACM SIGKDD workshop on Data Cleaning, Record Linkage and Object Identification, August 2003.
- [4] Erhard Rahm and H. Hai Do. Data cleaning: Problems and current approaches. IEEE Data Engineering Bulletin, December 2000
- [5] Timothy Ohanekwu, C.I Ezeife, A Token Based Data Cleansing Technique for Data warehouse System, IJDWM, June 2003.
- [6] Abdelkader Hameurlain, Rosine Cicchetti, Roland Traunmuller, Database and Expert System Applications, Proceedings of International Conference DEXA '97, 1997
- [7] Bing Chen, Beizhan Wang, Analysis and Solution of Data Quality in Data Warehouse of Chinese Materia Medica, Proceedings of International Conference on Computer Science, 2009.
- [8] Vijayshankar Raman , Joseph M. Hellerstein, Potter's Wheel: An Interactive Data Cleaning System, Proceedings of the 27th International Conference on Very Large Data Bases, September , 2001
- [9] R. Colin, Developing Universal Approach to Cleansing Customer and Product Data, BI research white paper from Business Object.
- [10] Philip Russom, Taking Data Quality to the Enterprise through Data Governance, A White Paper.
- [11] Ronald Forina, Data Equality –A Behind the Scene Perspective on Data Cleansing, <http://www.dmreview.com/>, March 2001