

Name Disambiguation Method Based on Attribute Match and Link Analysis*

Yu-Feng Yao

Changshu Institute of Technology, Computer Science and Engineering College, Suzhou, China.
Email: fengsingal81@gmail.com

Received October 30th, 2011; revised November 31st, 2011; accepted December 13th, 2011

ABSTRACT

A name disambiguation method is proposed based on attribute match and link analysis applying in the field of insurance. Aiming at the former name disambiguation methods such as text clustering method needs to be considered in a lot of useless words, a new name disambiguation method is advanced. Firstly, the same attribute matching is applied, merging the identity of a successful match, secondly, the link analysis is used, structural analysis of customers network is analyzed, Finally, the same cooperating information is merged. Experiment results show that the proposed method can realize name disambiguation successfully.

Keywords: Name Disambiguation; Data Mining; Attribute Match; Link Analysis

1. Introduction

Data Mining [1] is searching large amounts of data, reveals the hidden laws in it, and further models it to the advanced and effective method according to the established business objectives. It is the process of extracting the potentially useful information and knowledge in a lot of not complete, noisy, fuzzy, random data. Name disambiguation is to divide the different people using the same name to be represented according to the context or chapter information. The current systems can manage name disambiguation are such as SnakeT, Vivisimo and Apex, etc., the drawback is only the name of people is managed as an ordinary word. The clustering result labels are the vocabulary related with people name, can not distinguish the disambiguation results. The literature [2] proposed the character name disambiguation method based on social network. Bollegala D *et al.* [3] achieve clustering for the final result by calculating the vector similarity. The above methods [2-5] all considered a lot of useless words, and the digestion process has the strong reliance with the information extraction.

A new name disambiguation method is proposed in this paper, the method concludes three procedures: Firstly, the attribute match can be done, the attribute math algorithm is used to calculate the similarity between the customers; then the link analysis is followed, using the relationships between the customers to recognize the customer disambiguation problem. The experiment shows

the method is this paper can optimize the match processing of the different people having the same name.

2. The Respective Work

The original data of Field insurance records the contract number, insured name, age of the insured, the insured employment, insured unit, occupation, age and other information. The following are the definitions of several variables:

Identifier: The cooperative network node extracted from the original nodes data, is the carrier of customer information, with the letter "r" to represent;

Entities: The true customer related with the identifier, and it construct the relation of one to many, using the letter "e" to represent;

Attribute: Describe the author attribute such as address, profession etc. using the character "a" to represent, $r_i \cdot a_j$ represent the attribute of identifier a_j ; $r_i \equiv r_j$ represent r_i and r_j are the identifier of the same entity.

Relation: The relations between the two entities are called the relation, relation is a wide notion.

3. Name Disambiguation Method

3.1. Name Disambiguation Flow

Name disambiguation procedure can be mainly concluded as three steps: the first step is attribute match, to match the same identifier customer information, and combine the attribute with the successful matched identifier. The

*Supported by the research startup fund for young teachers of the Changshu Institute of Technology (Grant No. QZ0912).

Second step is the link analysis, through the structural analysis for the customer cooperative net to find the customers with the same identifier is cooperated with the same agent. entity, therefore, the customers having the same identifier is the same entity, and combine and analyze the entity with the same cooperative information. The steps are showed as **Figure 1**.

3.2. Attribute Match

For the customers having insured in some insurance company, the attribute has the name, customer name, purchased product name, and the customer information. Firstly, divide the two attributes c of the input entity, and match the first part of it showed as Equation (1):

$$\text{SubStringMatch}(r_i \cdot \text{Attr} \cdot \text{FirstPart}, r_j \cdot \text{Attr} \cdot \text{FirstPart}, R_l) \quad (1)$$

In Equation (1), R_l is the threshold value between 0 and 1, given the relative short string s in c , the relative long character string is l , and given the maximum long string of l in s is m , so the Equation (2) can be concluded.

$$\text{SubStringMatch} = \begin{cases} \text{true}, s \cdot \text{LENGTH} \leq m \cdot \text{LENGTH} \times R_l \\ \text{false}, s \cdot \text{LENGTH} > m \cdot \text{LENGTH} \times R_l \end{cases} \quad (2)$$

When the Equation (1) returns value is true, then the second part is matched by Equation (3):

$$\text{CountMatch}(r_i \cdot \text{Attr}, r_j \cdot \text{Attr}, R_c) \quad (3)$$

In Equation (3), R_c is also the threshold between 0 and 1, given s is the relative short string between $r_i \cdot \text{Attr}$ and $r_j \cdot \text{Attr}$, and the long string is l , and given the appearance times of s in string l is n , so the Equation (4) can be conclude:

$$\text{CountMatch} = \begin{cases} \text{true}, n \geq s \cdot \text{LENGTH} \times R_c \\ \text{false}, n < s \cdot \text{LENGTH} \times R_c \end{cases} \quad (4)$$

The attribute match function showed as Equation (5) can be obtained by combined Equations (1) and (3):

$$\begin{aligned} \text{Match}_2(r_i, r_j, R_l, R_c) = \\ \text{SubStringMatch} \\ (r_i \cdot \text{DEPT} \cdot \text{FirstPart}, r_j \cdot \text{DEPT} \cdot \text{FirstPart}, R_l) \\ \wedge \text{CountMatch}(r_i \cdot \text{DEPT}, r_j \cdot \text{DEPT}, R_c) \end{aligned} \quad (5)$$

Considering the below case: Identifier A is the same with identifier B, and they have the insurance cooperative relation with the identifier C, then identifier A and the identifier B correspond with the same customer entity, showed as the Equation (6):

$$\begin{aligned} \text{If } \exists r_x, (r_i, r_x) \in S_c \wedge (r_j, r_x) \in S_c \\ \wedge r_i \cdot \text{NAME} = r_j \cdot \text{NAME}, \text{ then } r_i \equiv r_j. \end{aligned} \quad (6)$$

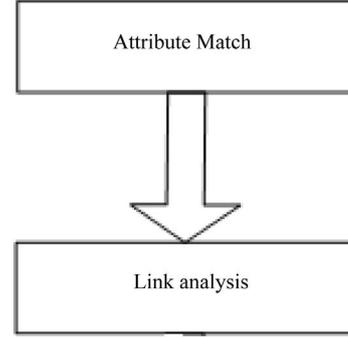


Figure 1. Name disambiguation flow.

4. Simulation Experiment

4.1. Experiment Parameters

The data set used in the experiment is from the core business system of some big insurance company, mainly contains integrated business manage system, universal system and pension system etc. The insured information contains customer number, social insurance number, certificate type, certificate number, occupation categories, subcategories, small categories, health state, smoke years, the smoked number, marriage status, the relationship with the insured; the information of the insured contains customer number, certificate type, certificate number and name etc. Insurance information contains insured number, sign date, effective date, paid premium and policy duration. Insurance category contains insurance code, insurance name, insurance type, duration and the agent-related information. The data set is showed as **Figure 2**.

4.2. Attribute Match Simulation

It is a match experiment showed in **Table 1** after applied Equation (5) in the dataset.

4.3. Link Analysis Simulation

Using the node to represent the insurance relation person such as insurer, insured, and agent, the node size represents the degree of the node, the edge of connecting node represents the cooperative relations between authors. The customer name is labeled at the edge. For the researched sample object, the unit information of the customer can be labeled. For the link analysis, the perfect manage result is: the same name customer not in the same cooperative customer entities are different customer entities, and the same name author in one cooperative entities are attribute to the customer entity. Before the link analysis simulation, the net of three people called Zhao aixiang is showed as **Figure 3**, and the resulted result after link analysis is showed as **Figure 4**.

For the distance between two same name authors called

Table 1. Attribute match.

	Id	Gender	Name	birthday	age	Body state	Order date
Before Match	1315	F	Wang Si	1953.10.23	55	1	2008-5-7
	1328	F	Wang Si	1953.10.23	55	0	2008-5-7
After Match	1328	F	Wang Si	1953.10.23	55	rep	2008-5-7

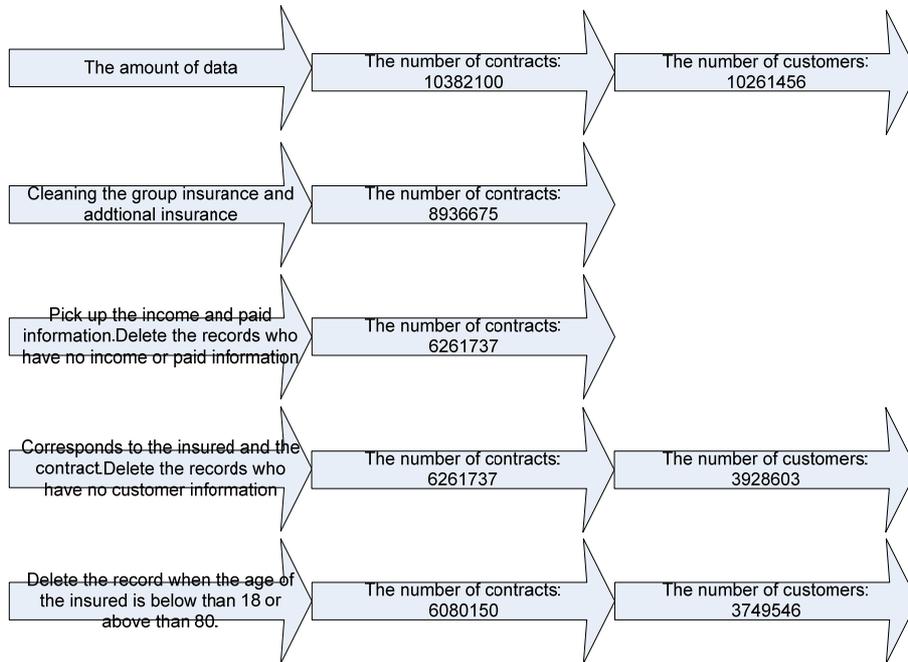


Figure 2. Data set of some insurance company.

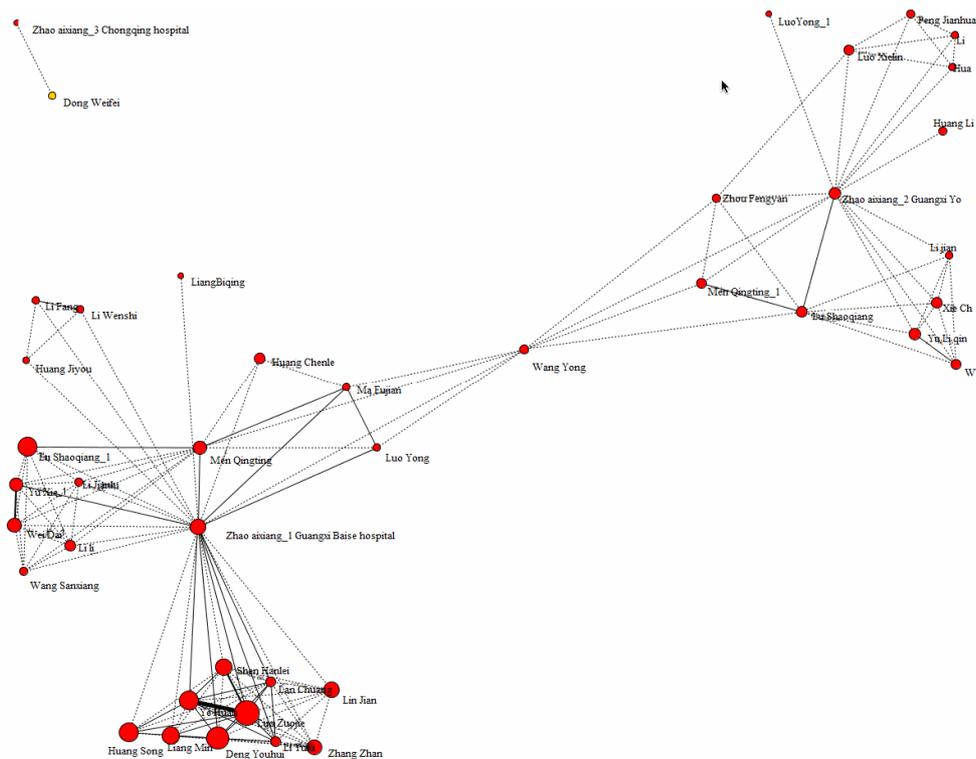


Figure 3. Customer net of “Zhao aixiang” before link analysis.

