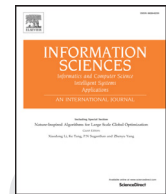


Contents lists available at [ScienceDirect](#)

Information Sciences

journal homepage: www.elsevier.com/locate/ins

Missing value imputation for the analysis of incomplete traffic accident data

Rupam Deb*, Alan Wee-Chung Liew

School of Information and Communication Technology, Gold Coast Campus, Griffith University, QLD 4222, Australia

ARTICLE INFO

Article history:

Received 24 April 2015
Revised 30 December 2015
Accepted 2 January 2016
Available online xxx

Keywords:

Data preprocessing
Decision tree
Missing value imputation
Categorical data
Traffic accident

ABSTRACT

Death, injury and disability resulting from road traffic crashes continue to be a major global public health problem. Recent data suggest that the number of fatalities from traffic crashes is in excess of 1.25 million people each year with non-fatal injuries affecting a further 20–50 million people. It is predicted that by 2030 road traffic accidents will have progressed to be the 5th leading cause of death and that the number of people who will die annually from traffic accidents will have doubled from current levels. Both developed and developing countries suffer from the consequences of increase in human population, and therefore, vehicle population. Therefore, methods to reduce accident severity are of great interest to traffic agencies and the public at large. To analyse traffic accident factors effectively we need a complete traffic accident historical database. Any missing data in the database could prevent the discovery of important environmental and road accident factors and lead to invalid conclusions. In this paper, we present a novel imputation method that exploits the within-record and between-record correlations to impute missing data of numerical or categorical values. In addition, our algorithm accounts for uncertainty in real world data by sampling from a list of potential imputed values according to their affinity degree. We evaluated our algorithm using four publicly available traffic accident databases from the United States, the first of which is the largest open federal database (explore.data.gov) in the United States, and the second is based on the National Incident Based Reporting System (NIBRS) of the city and county of Denver (data.opencolorado.org). The other two are from New York's open data portal (Motor Vehicle Crashes-case information: 2011 and Motor Vehicle Crashes-individual information: 2011, data.ny.gov). We compare our algorithm with four state-of-the-art imputation methods using missing value imputation accuracy and RMSE. Our results indicate that the proposed method performs significantly better than the existing algorithms we compared.

© 2016 Published by Elsevier Inc.

1. Introduction

The high growth of the number of vehicles has led to roads with higher traffic density. The immediate effect of this situation is the dramatic increase of traffic accidents on the road, which has become a serious problem in many countries. For example, 2478 people died on Spanish roads in 2010, which means one death for every 18,551 inhabitants [7,8,11]. In the United States (according to the Department of Transportation, United States), in 2012, 33,561 people died in motor vehicle

* Corresponding author. Tel.: +61 426932694.

E-mail addresses: rupam.deb@griffithuni.edu.au, cse.rupam@gmail.com (R. Deb), a.liew@griffith.edu.au (A.W.-C. Liew).

6 traffic crashes. According to the Australian Bureau of Statistics, the majority of transport related deaths (72% in 2009) in
7 Australia is associated with motor vehicles driven on public roads.

8 A nation's socio-economic development is highly dependent on the health status of its citizens. Road safety, which is
9 mainly affected by road accidents is said to be one of the major health concerns. The burden of road accident causalities and
10 damage is a major headache for both developed and developing countries. To reduce the number of road traffic accidents, it
11 is necessary to characterize the causes of the accidents and also determine the factors which significantly affect the severity
12 of injuries in road crashes. The global economic cost of road traffic accidents has been estimated at US\$518b and has been
13 calculated to account for 0.3–4% of the gross national product of many countries [15].

14 Digital data acquisition methods and storage technology have resulted in the growth of a huge amount of traffic data be-
15 ing stored in different types of databases. Large amounts of traffic accident data have been collected with the advancement
16 in sensor technologies [34]. Using data mining technology, such as classification and clustering, we can uncover patterns of
17 traffic activities and factors which lead to accident. The major reason that data mining has attracted a great deal of attention
18 in information industry is due to the wide availability of huge amounts of data and the imminent need for turning such
19 data into useful information and knowledge [28]. Data mining is the exploration and analysis of large dataset in order to
20 discover knowledge and rules. Data mining is typically conceptualized as a three part process: preprocessing, learning, and
21 post-processing.

22 Appropriate data preprocessing involves transforming the raw data into suitable format for subsequent analysis. To run
23 the classification and clustering algorithms, there is a strong need for good data preprocessing to ensure the data are of
24 good quality. Data preprocessing takes almost 80% of the total data mining effort [35]. It is also known that good results can
25 be achieved by using data mining algorithms only if there is a good quality dataset [19].

26 Real-life data are frequently imperfect: erroneous, incomplete, uncertain and vague. Datasets often have missing values
27 due to various reasons including equipment malfunctioning, human errors, and faulty data transmission. If an organization
28 does not take extreme care during data collection, then approximately 5% or more missing/corrupt data may be introduced
29 into the datasets [10,20,33]. Data preprocessing includes imputation of missing values, smoothing out noisy data, identifica-
30 tion of incorrect data, and correction of inconsistent data. In this work, we propose a novel imputation algorithm to address
31 the problem of missing attribute values.

32 Specifically, we propose a new decision tree and sampling based missing value imputation algorithm called DSMI for
33 missing value imputation. Our algorithm measures the correlation after horizontally dividing the large dataset depending
34 on the missing record(s). We divide the whole dataset on different horizontal segments based on non-missing attributes of
35 the missing records and impute the missing values based on sampling from the distributions induced by the correlation.
36 Our experiments show that the proposed algorithm has better imputation accuracy compared with several other existing
37 algorithms. For validation, we compare our algorithm with some well-known missing value imputation techniques using
38 imputation accuracy and RMSE.

39 In our study, we use two truck crash datasets, the accident dataset of Denver County, and two crash datasets from New
40 York's open data portal. Most of the columns (attributes) of these records are categorical. The domain of the categorical data
41 is represented by nominal, ordinal and interval based variables. Nominal variables have values that have no natural ordering
42 (e.g. airbag conditions: Ruptured, Cut, Torn); Ordinal variables do have a natural order (e.g. day of the week); and Interval
43 variables are created from intervals on a contiguous scale (e.g. age group 13–19).

44 This paper is organized as follows: in Section 2 we present a literature review of related work. Our proposed technique
45 is described in Section 3. Experimental results are discussed in Section 4. Finally Section 5 draws the concluding remarks.

46 2. Related work

47 Imputation of missing values is an important data mining task for improving the quality of the data mining result. Many
48 missing value imputation algorithms have been proposed for various applications [1,3,4,14,16–18,23–27,30]. Some of these
49 methods are: expectation maximization imputation (EMI) [27], decision tree based missing value imputation (DMI) [26],
50 similarity based missing value imputation (SiMI) [26], k -decision tree based missing value imputation (k DMI) [25], k -nearest
51 neighbour based imputation (k NNI) [1], local weighted linear approximation imputation (LWLA) [18], and framework for
52 imputing missing values using co-appearance, correlation and similarity analysis (FIMUS) [24].

53 To impute numerical missing values, the EMI algorithm [27] relies on estimating the mean and covariance matrix of
54 the dataset. The EMI algorithm starts with an initial estimate of the mean and the covariance matrix, and iterates until
55 the imputed values and the estimates of mean and covariance matrix stop changing appreciably from the current iteration
56 to the next iteration [9,27]. The EMI algorithm is only applicable to datasets in which the missing values are missing at
57 random. The main drawback of this method is that for imputing the missing value, EMI algorithm uses information from
58 the whole dataset and therefore is suitable only for datasets that exhibit strong correlations between attributes.

59 Instead of using information from the whole dataset, k NNI method [1] imputes missing values using k number of similar
60 records. This method first finds user-defined k number of records from the total dataset by using the Euclidean distance
61 measure. To impute a numerical missing value, the method utilizes the mean value of the specific attribute within the k
62 most similar records of the entire dataset. If the missing attribute is categorical, then the method utilizes the most frequent
63 value of the attribute within the k most similar records. k NNI is a simple method that performs well on dataset that has
64 strong local correlation structure. However, the method can be expensive for large dataset since for each record with missing

Download English Version:

<https://daneshyari.com/en/article/6857386>

Download Persian Version:

<https://daneshyari.com/article/6857386>

[Daneshyari.com](https://daneshyari.com)