

[USPTO PATENT FULL-TEXT AND IMAGE DATABASE](#)

(1 of 1)

United States Patent
Cassidy , et al.

10,558,627
February 11, 2020

Method and system for cleansing and de-duplicating data

Abstract

Method and system for cleansing and de-duplicating data in database are provided. The method includes filtering garbage records from a plurality of records based on data fields, and applying cleansing rules to create a cleansed database. A similarity vector is generated, where each vector corresponds to pairwise comparison of distinct data entries in cleansed database. Matching rules are applied to label each vector as one of matched, unmatched and unclassified. The method analyzes the vectors labeled as matched and unmatched to train a machine learning model to identify duplicates in the cleansed database. Unclassified vectors in the cleansed database are labeled as matched or unmatched by applying machine learning model on unclassified vectors. Thereafter, the method processes all the vectors labeled as matched to create clusters of records that are duplicates of each other. Further, records in each cluster are merged to obtain de-duplicated cleansed database using predefined consolidated rules.

Inventors: **Cassidy; Hugh** (Campbell, CA), **DeMarco; Sofia** (Mountain View, CA), **Lakshmikanthan; Jayant** (San Jose, CA)

Applicant: **Name** **City** **State** **Country** **Type**

LeanTaas, Inc. Santa Clara CA US

Assignee: *LeanTaas, Inc.* (Santa Clara, CA)

Family ID: 60089645

Appl. No.: 15/488,388

Filed: April 14, 2017

Prior Publication Data**Document Identifier**

US 20170308557 A1

Publication Date

Oct 26, 2017

Related U.S. Patent Documents**Application Number**

62325968

Filing Date

Apr 21, 2016

Patent Number**Issue Date**

Current U.S. Class:

1/1

Current CPC Class:

G06N 20/00 (20190101); G06F 12/0253 (20130101); G06F 16/24556 (20190101); G06F 16/215 (20190101)

Current International Class:

G06F 16/215 (20190101); G06F 12/02 (20060101); G06N 20/00 (20190101); G06F 16/2455 (20190101)

References Cited [\[Referenced By\]](#)**U.S. Patent Documents**

6961721	November 2005	Chaudhuri
7672942	March 2010	Weinberg
8688603	April 2014	Kurup
8793201	July 2014	Wang
8914366	December 2014	Li
9471609	October 2016	Kienzle
2004/0107203	June 2004	Burdick
2004/0158562	August 2004	Caulfield
2009/0245573	October 2009	Saptharishi
2013/0054541	February 2013	Kaldas
2013/0238623	September 2013	Wyllie
2015/0269494	September 2015	Kardes
2016/0092494	March 2016	Kabra et al.
2016/0180245	June 2016	Tereshkov
2017/0083825	March 2017	Battersby
2017/0242891	August 2017	Doan
2019/0080247	March 2019	Dubey

Primary Examiner: Richardson; James E
Attorney, Agent or Firm: Fenwick & West LLP

Claims

What is claimed is:

1. A computer-implemented method for cleansing and de-duplicating data in a database, the computer-implemented method comprising: identifying garbage records from a plurality of records in the database based on distinct data entries of data fields from each of the plurality of records; applying cleansing rules to create a cleansed database by removing the garbage records and standardizing the distinct data entries; for each pair of records in the cleansed database: generating a similarity vector using one or more string matching algorithms, wherein the similarity vector corresponds to a pairwise comparison of the pair of records and each dimension of the similarity vector represents a similarity score between distinct data entries of one data field of the distinct records; and labeling the similarity vector as a matched vector, an unmatched vector, or an unclassified vector by applying matching rules to the distinct data entries of the pair of records; training a machine learning model to identify matched records in the cleansed database by analyzing the labelled similarity vectors labeled as matched vector and unmatched vectors; applying the machine learning model to the unclassified vectors to label each of the unclassified vectors in the cleansed database as matched vector or unmatched vector; processing all of the matched vectors in the cleansed database to create clusters of matched records; and merging the matched records in each cluster to obtain a de-duplicated cleansed database using one or more predefined consolidated rules.

2. The computer implemented method of claim 1 further comprising prior to the step of filtering the garbage records, extracting the data fields based on a profiling of data sets in the database, wherein the profiling of

Over time, as data entry and merging of records from different sources occur, duplicate copies may begin to creep into the database. Such occurrence of duplicate copies is referred to as data duplication. Storing duplicate data in a database is inefficient for several reasons, for example, duplicate data could make pricing analysis almost impossible, duplicate vendor data could make any vendor rationalization difficult, duplicate data may lead to memory constraints, and the like. Therefore, identifying and eliminating duplicate data is one of the major problems in the area of data cleaning and data quality. Several approaches have been implemented to counter the problem of data duplication. However, none of the approaches are effective specifically in large-scales.

SUMMARY

Various methods, systems and computer readable mediums for cleansing and de-duplicating data in a database disclosed. In an embodiment, a computer implemented method for cleansing and de-duplicating data in a database is disclosed. The computer-implemented method includes filtering unnecessary records from a plurality of records in the database based on data fields. A cleansed database is then created by applying cleansing rules to remove the unnecessary records. The computer-implemented method then generates similarity vectors, wherein each vector corresponds to a pairwise comparison of distinct data entries in the cleansed database. Further, matching rules are applied to find matched data and unmatched data in the cleansed database for labeling each vector in the similarity vectors as one of matched, unmatched and unclassified. The computer-implemented method thereafter analyses the vectors labeled as matched and unmatched to train a machine learning model to identify duplicates in the cleansed database. The unclassified vectors in the cleansed database are then labeled as matched or unmatched by applying the machine learning model. The computer-implemented method then processes all the vectors labeled as match to create clusters of records that are duplicates of each other in the cleansed database. Further, the records in each cluster are merged to obtain a de-duplicated cleansed database using predefined consolidated rules.

In another embodiment, a system for cleansing and de-duplication of data is disclosed. The system includes a memory and a processor. The memory stores instructions for cleansing and de-duplicating data. The processor is operatively coupled with the memory to fetch instructions from the memory for cleansing and de-duplicating data. The processor is configured to filter unnecessary records from a plurality of records in the database based on data fields. The processor then creates a cleansed database by applying cleansing rules to remove the unnecessary records. Thereafter, the processor generates similarity vectors, wherein each vector corresponds to a pairwise comparison of distinct data entries in the cleansed database. The processor further applies matching rules to find matched data and unmatched data in the cleansed database for labeling each vector in the similarity vectors as one of matched, unmatched and unclassified. The processor also analyses the vectors labeled as matched and unmatched to train a machine learning model to identify duplicates in the cleansed database. The unclassified vectors are then labeled by the processor as matched or unmatched by applying the machine learning model. Further, all the vectors labeled as match are processed by the processor to create clusters of records that are duplicates of each other in the cleansed database. Subsequently, the processor merges records in each cluster to obtain a de-duplicated cleansed database using predefined consolidated rules.

In yet another embodiment, a computer system for cleansing and de-duplication of data is disclosed. The computer system includes a processor and an application program. The application program is executed by the processor. The application program filters unnecessary records from a plurality of records based on data fields. A cleansed database is then created by applying cleansing rules to remove the unnecessary records. The application program then generates similarity vectors, wherein each vector corresponds to a pairwise comparison of distinct data entries in the cleansed database. Further, matching rules are applied to find matched data and unmatched data in the cleansed database for labeling each vector in the similarity vectors as one of matched, unmatched and unclassified. The application program thereafter analyses the vectors labeled as matched and unmatched to train a machine learning model to identify duplicates in the cleansed database. The unclassified vectors are then labeled in the cleansed database as matched or unmatched by applying the machine learning model on the unclassified vectors. All the vectors labeled as match is then processed by the application program to create clusters of records that are duplicates of each other in the cleansed database. Subsequently, the application program merges the records in each cluster obtain a de-duplicated cleansed database using predefined consolidated rules.

Other aspects and example embodiments are provided in the drawings and the detailed description that follows.

BRIEF DESCRIPTION OF THE FIGURES

For a more complete understanding of example embodiments of the present technology, reference is now made to the following descriptions taken in connection with the accompanying drawings in which:

FIG. 1 illustrates a schematic representation of an example environment 100, where at least some example embodiments of the present invention can be implemented;

FIG. 2 illustrates a block diagram depicting different modules included in a data cleansing and de-duplication system, in accordance with an example embodiment;

FIG. 3 is a flow chart describing a method for cleansing and de-duplicating data, in accordance with an example embodiment;

FIG. 4 illustrates a block diagram of a data cleansing and de-duplicating system, in accordance with an example embodiment;

FIG. 5 illustrates a pictorial representation of records removed at each stage while performing the method for cleansing and de-duplicating data, in accordance with an example embodiment; and

FIG. 6 is a block diagram of a machine in the example form of a computing device within which instructions for causing the machine to perform any one or more of the methodologies discussed herein may be executed, in accordance with an example embodiment.

The drawings referred to in this description are not to be understood as being drawn to scale except if specifically noted, and such drawings are only example in nature.

DETAILED DESCRIPTION

In the following description, for purposes of explanation, numerous specific details are set forth in order to provide a thorough understanding of the present disclosure. It will be apparent, however, to one skilled in the art that the present disclosure can be practiced without these specific details. In other instances, apparatuses and methods are shown in block diagram form only in order to avoid obscuring the present disclosure.

Reference in this specification to "one embodiment" or "an embodiment" means that a particular feature, structure, or characteristic described in connection with the embodiment is included in at least one embodiment of the present disclosure. The appearance of the phrase "in one embodiment" in various places in the specification are not necessarily all referring to the same embodiment, nor are separate or alternative embodiments mutually exclusive of other embodiments.

Moreover, although the following description contains many specifics for the purposes of illustration, anyone skilled in the art will appreciate that many variations and/or alterations to said details are within the scope of the present disclosure. Similarly, although many of the features of the present disclosure are described in terms of each other, or in conjunction with each other, one skilled in the art will appreciate that many of these features can be provided independently of other features. Accordingly, this description of the present disclosure is set forth without any loss of generality to, and without imposing limitations upon, the present disclosure.

The term data used throughout the description represents a plurality of records stored in a database. Further, a machine learning algorithm is used to build a machine learning model which is then to classify vectors in the proposed method. Herein, unless the context suggests otherwise, the terms `machine learning algorithm` and `machine learning model` are used interchangeably for the purposes of this description.

FIG. 1 illustrates a schematic representation of an example environment 100, where at least some example embodiments of the present invention can be implemented. The environment 100 includes a computing

device 102, a network 104 and a database 106.

The computing device 102 is a portable electronic or a desktop device configured with a user interface (not shown in FIG. 1) to interact with a user of the computing device 102. Examples of the computing device 102 include, but are not limited to, a personal computer (PC), a mobile phone, a tablet device, a personal digital assistant (PDA), a smart phone and a laptop. Examples of the user interface include, but are not limited to, display screen, keyboard, mouse, light pen, appearance of a desktop, illuminated characters and help messages.

The computing device 102 is also associated with a data cleansing and de-duplication system 108. The data cleansing and de-duplication system 108 is configured with a non-transitory computer-readable medium (referred as "Prime software"), the content of which causes to perform the method disclosed herein. For the sake of clarity and for the purpose of this description the "prime software" is referred as software. Typically, the software uses a combination of user-defined rules and machine learning algorithms to traverse multiple combinations of records rapidly to cleanse and isolate potential duplicates in the database 106.

The data cleansing and de-duplication system 108 which includes the software, comprising the invention, may be directly associated with the computing device 102. In an embodiment, the data cleansing and de-duplication system 108 is associated with different computing devices (not shown in FIG. 1) and may be accessible by the computing device 102 over the network 104. In another embodiment, the data cleansing and de-duplication system 108 may be associated with one or more servers (not shown in FIG. 1) and is accessible by the computing device 102 over the network 104.

The network 104 is a group of points or nodes connected by communication paths. The communication paths may be connected by wires, or they may be wirelessly connected or may use any such combination. The network 104 as defined herein can interconnect with other networks and may contain subnetworks. Examples of the network 104 includes, but are not limited to, a local area network (LAN), a personal area network (PAN), a metropolitan area network (MAN), and a wide area network (WAN).

The data cleansing and de-duplication system 108 may include a plurality of modules for performing different steps and functions associated with cleansing and de-duplicating data. The plurality of modules is explained in conjunction with FIG. 2. In an embodiment, the steps and function performed by different modules may be combined into a single module, or into other combinations of modules.

In the environment 100, the computing device 102 is associated with the database 106 through the network 104. In an embodiment, the database 106 may be directly associated with the computing device 102. The database 106 is configured to store data after cleansing and de-duplication. In an embodiment, the database 106 also stores data that includes garbage and duplicated data. In the embodiment, the data contains a plurality of records that are duplicates and may include garbage data.

In the environment 100, the data, including the garbage data and duplicate records, is first diagnosed to extract fields that can be used to filter out unwanted and unnecessary records (also referred as `junk` or `garbage records`). In an embodiment, a unique identity, for example a signature, is created to be assigned to each record. The garbage records and values are removed from the database 106 with the aid of cleansing rules. Consequently, a cleansed database that is ready for de-duplication is obtained.

Rules are then applied to the data in the cleansed database to identify matches and non-matches. This will help train a machine learning model to identify patterns in the data. The trained machine learning model then extracts patterns and identifies duplicates in the entire data. Pair-wise comparisons are made in the cleansed data to label the pairs as a match or a non-match. All match-pairs are then processed to create clusters. The clusters are then merged using predefined consolidated rules. The merged records could reside in the database 106, in an external database or in other source systems.

It should be appreciated by those of ordinary skill in the art that FIG. 1 depicts the computing device in an oversimplified manner and a practical embodiment may include additional components and suitably configured processing logic to support known or conventional operating features that are not described in detail herein.

FIG. 5 illustrates a pictorial representation of records removed at each stage while performing the method for cleansing and de-duplicating data, in accordance with an example embodiment.

At 502, a bar represents a total number of records in a database. For example, it shows `250,000` records. At 504, number of garbage records (junk records) is filtered. In an embodiment, the garbage records are filtered out by the data diagnostic module 202 and the data cleansing module 204. In another embodiment, the garbage records are filtered out by the processor 404. The numbers of junk records are identified as `80,000`. Thus, the number of records for de-duplication, after removing junk records, is identified by excluding the total number junk records (`80,000`) from the total number of records in the database (`250,000`). At 506, the number of records considered for de-duplication is depicted. For example, the number of records considered for de-duplication are `170,000` (i.e. `250,000` - `80,000`).

Thereafter, a de-duplication method described in the present description is applied on the records, thereby de-duplicating the duplicate records. At 508, number of records that are de-duplicated is depicted as `15,000`. Thus, the remaining records in the database could be identified by subtracting the number of records that are de-duplicated from the number of records considered for de-duplication. At 510, the remaining records after de-duplication are depicted as `155,000` (i.e. `170,000` - `15,000`).

Thereafter, a predefined consolidation rule is applied. For example, inactive records are removed from the remaining records, for example the records left after de-duplication. At 512, inactive records with no transactions in the past five (5) years are identified, for example such records are 50,000. At 514, the inactive records are removed from the de-duplicated records and a cleansed database is prepared. Thus, at 514, final set of records in the cleansed and de-duplicated database is depicted as `105,000` (i.e. `155,000` - `50,000`).

Machine Execution

FIG. 6 is a block diagram of a machine in the example form of a computer system within which instructions for causing the machine to perform any one or more of the methodologies discussed herein may be executed. In alternative embodiments, the machine operates as a standalone device or may be connected (e.g., networked) to other machines. In a networked deployment, the machine may operate in the capacity of a server or a client machine in a server-client network environment, or as a peer machine in a peer-to-peer (or distributed) network environment. The machine may be a personal computer (PC), a tablet PC, a set-top box (STB), a Personal Digital Assistant (PDA), cellular telephone, a web appliance, a network router, switch or bridge, or any machine capable of executing instructions (sequential or otherwise) that specify actions to be taken by that machine. Further, while only a single machine is illustrated, the term "machine" shall also be taken to include any collection of machines that individually or jointly execute a set (or multiple sets) of instructions to perform any one or more of the methodologies discussed herein.

The example computer system 600 includes a processor 602 (e.g., a central processing unit (CPU), a graphics processing unit (GPU), or both), a main memory 604, and a static memory 606, which communicate with each other via a bus 608. The computer system 600 may further include a video display unit 610 (e.g., a liquid crystal display (LCD) or a cathode ray tube (CRT)). The computer system 600 also includes an alphanumeric input device 612 (e.g., a keyboard), a user interface (UI) navigation device 614 (e.g., a mouse), a disk drive unit 616, a signal generation device 618 (e.g., a speaker), and a network interface device 620. The computer system 600 may also include an environmental input device 628 that may provide a number of inputs describing the environment in which the computer system 600 or another device exists, including, but not limited to, any of a Global Positioning Sensing (GPS) receiver, a temperature sensor, a light sensor, a still photo or video camera, an audio sensor (e.g., a microphone), a velocity sensor, a gyroscope, an accelerometer, and a compass.

The disk drive unit 616 includes a machine-readable medium 622 on which is stored one or more sets of data structures and instructions 624 (e.g., software) embodying or utilized by any one or more of the methodologies or functions described herein. The instructions 624 may also reside, completely or at least partially, within the main memory 604 and/or within the processor 602 during execution thereof by the computer system 600, the main memory 604 and the processor 602 also constituting machine-readable media. In an embodiment, an application program 626 is also executed by the processor 602.

The application program 626 is software (also referred as "Prime Software") that may reside in the memory (604 or 606) or it may be a separate component in the computer system 600. In an embodiment, the application program 626 is a part of different computer system and is associated with the computing system 600 through a computer network, for example the computer network 650. The application program 626 includes instructions that are capable of cleansing and de-duplication of data in a database. The application program 626 is executed by the processor 602. The database is either associated with the computer system 600 or is a part of different computer system. In an embodiment, the application program 626 performs all the function of performed by different modules included in the data cleansing and de-duplicating system 108 (explained in conjunction with FIG. 2).

The application program 626 first filters garbage records from a plurality of records based on data fields. A cleansed database is then created by applying cleansing rules to remove the garbage records. The application program 626 then generates similarity vectors, wherein each vector corresponds to a pairwise comparison of distinct data entries in the cleansed database. Further, the application program 626 applies matching rules to find matched data and unmatched data in the cleansed database for labeling each vector in the similarity vectors as one of matched, unmatched and unclassified.

The application program 626 thereafter analyses the vectors labeled as matched and unmatched to train a machine learning model to identify duplicates in the cleansed database. Further, the application program 626 labels the unclassified vectors in the cleansed database as matched or unmatched by applying the machine learning model on the unclassified vectors. Thereafter, all the vectors labeled as match is processed to create clusters of records that are duplicates of each other in the cleansed database. Subsequently, the application program 626 merges the records in each cluster are to obtain a de-duplicated cleansed database using predefined consolidated rules.

While the machine-readable medium 622 is shown in an example embodiment to be a single medium, the term "machine-readable medium" may include a single medium or multiple media (e.g., a centralized or distributed database, and/or associated caches and servers) that store the one or more instructions 624 or data structures. The term "non-transitory machine-readable medium" shall also be taken to include any tangible medium that is capable of storing, encoding, or carrying instructions for execution by the machine and that cause the machine to perform any one or more of the methodologies of the present subject matter, or that is capable of storing, encoding, or carrying data structures utilized by or associated with such instructions. The term "non-transitory machine-readable medium" shall accordingly be taken to include, but not be limited to, solid-state memories, and optical and magnetic media. Specific examples of non-transitory machine-readable media include, but are not limited to, non-volatile memory, including by way of example, semiconductor memory devices (e.g., Erasable Programmable Read-Only Memory (EPROM), Electrically Erasable Programmable Read-Only Memory (EEPROM), and flash memory devices), magnetic disks such as internal hard disks and removable disks, magneto-optical disks, and CD-ROM and DVD-ROM disks.

The instructions 624 may further be transmitted or received over a computer network 650 using a transmission medium. The instructions 624 may be transmitted using the network interface device 620 and any one of a number of well-known transfer protocols (e.g., HTTP). Examples of communication networks include a local area network (LAN), a wide area network (WAN), the Internet, mobile telephone networks, Plain Old Telephone Service (POTS) networks, and wireless data networks (e.g., Wi-Fi and WiMAX networks). The term "transmission medium" shall be taken to include any intangible medium that is capable of storing, encoding, or carrying instructions for execution by the machine, and includes digital or analog communications signals or other intangible media to facilitate communication of such software.

As described herein, computer software products can be written in any of various suitable programming languages, such as C, C++, C#, Pascal, Fortran, Perl, Matlab (from MathWorks), SAS, SPSS, JavaScript, AJAX, and Java. The computer software product can be an independent application with data input and data display modules. Alternatively, the computer software products can be classes that can be instantiated as distributed objects. The computer software products can also be component software, for example Java Beans or Enterprise Java Beans. Much functionality described herein can be implemented in computer software, computer hardware, or a combination.

Furthermore, a computer that is running the previously mentioned computer software can be connected to a network and can interface to other computers using the network. The network can be an intranet, internet, or

the Internet, among others. The network can be a wired network (for example, using copper), telephone network, packet network, an optical network (for example, using optical fiber), or a wireless network, or a combination of such networks. For example, data and other information can be passed between the computer and components (or steps) of a system using a wireless network based on a protocol, for example Wi-Fi (IEEE standard 802.11 including its sub-standards a, b, e, g, h, i, n, et al.). In one example, signals from the computer can be transferred, at least in part, wirelessly to components or other computers.

It is to be understood that although various components are illustrated herein as separate entities, each illustrated component represents a collection of functionalities which can be implemented as software, hardware, firmware or any combination of these. Where a component is implemented as software, it can be implemented as a standalone program, but can also be implemented in other ways, for example as part of a larger program, as a plurality of separate programs, as a kernel loadable module, as one or more device drivers or as one or more statically or dynamically linked libraries.

The present disclosure is described above with reference to block diagrams and flowchart illustrations of method and device embodying the present disclosure. It will be understood that various block of the block diagram and flowchart illustrations, and combinations of blocks in the block diagrams and flowchart illustrations, respectively, may be implemented by a set of computer program instructions. These set of instructions may be loaded onto a general purpose computer, special purpose computer, or other programmable data processing apparatus to cause a device, such that the set of instructions when executed on the computer or other programmable data processing apparatus create a means for implementing the functions specified in the flowchart block or blocks. Although other means for implementing the functions including various combinations of hardware, firmware and software as described herein may also be employed.

Various embodiments described above may be implemented in software, hardware, application logic or a combination of software, hardware and application logic. The software, application logic and/or hardware may reside on at least one memory, at least one processor, an apparatus or, a non-transitory computer program product. In an example embodiment, the application logic, software or an instruction set is maintained on any one of various conventional computer-readable media. In the context of this document, a "computer-readable medium" may be any non-transitory media or means that can contain, store, communicate, propagate or transport the instructions for use by or in connection with an instruction execution system, apparatus, or device, such as a computer, with one example of a system described and depicted in FIG. 6. A computer-readable medium may comprise a computer-readable storage medium that may be any media or means that can contain or store the instructions for use by or in connection with an instruction execution system, apparatus, or device, such as a computer.

The foregoing descriptions of specific embodiments of the present disclosure have been presented for purposes of illustration and description. They are not intended to be exhaustive or to limit the present disclosure to the precise forms disclosed, and obviously many modifications and variations are possible in light of the above teaching. The embodiments were chosen and described in order to best explain the principles of the present disclosure and its practical application, to thereby enable others skilled in the art to best utilize the present disclosure and various embodiments with various modifications as are suited to the particular use contemplated. It is understood that various omissions and substitutions of equivalents are contemplated as circumstance may suggest or render expedient, but such are intended to cover the application or implementation without departing from the spirit or scope of the claims of the present disclosure.

* * * * *



