

A New Approach for Data Cleaning to Improve Quality of Data Warehouse

Ali Shahnava¹, Mehdi Afzali^{2*}, Shima Rahimzadeh³

1. Assistant Professor, Department of Mathematics and Statistics, Zanzan Branch, Islamic Azad University, Zanzan, Iran. shahnava_ali2000@yahoo.com
2. Assistant Professor, Department of IT Engineering, Zanzan Branch, Islamic Azad University, Zanzan, Iran. (Corresponding Author) afzali@iauz.ac.ir
3. MSc Student, Zanzan University of Medical Sciences, Zanzan, Iran.

Abstract

Introduction: Business intelligence is defined as the process of converting data into information and then knowledge. Business intelligence represents a broad application and technology space for collecting, storing, analyzing and accessing information to improve high-quality business process modeling. Business intelligence consists of several steps. In the first step, data sources are collected. These sources can include data from various types of databases or information from existing software. The information collected during the "ETL" process is loaded into the analytical database or data warehouse. There have been several works addressing data cleaning, including correcting spelling and phonetic errors, unifying the format of the date field, correcting quantitative data, deleting duplicate records, etc. But based on the research done, nothing has been done to resolve the contradiction in interrelated fields (e.g. geographical locations), and since this field is available in the organization's database, its resolution is a problem for improving the organization's performance. It is important and necessary that it should be addressed.

Method: It is necessary to use the standard data of provinces and cities of Iran, and for this purpose, a database containing the information of the provinces and cities of Iran and the first three digits of the national code for each city was created in the Microsoft Access 2010, and the connection between the related tables was established. The health certificate database of Zanzan University of Medical Sciences students who graduated in 1992 and 1993 was received from the information technology department of Zanzan University of Medical Sciences. These data were based on different semesters for two consecutive years. The proposed program is written in the Visual Studio 2013 environment and has two databases, one is the standard culture of provinces and cities and the three digits of their codes, and the other is a student database for testing and resolving data contradictions.

Finding: Through the implementation of this approach we have been able to detect dirty data and then by using students' national codes, the correction process has been applied to them. Based on the achieved results, the amount of dirty data decreased from 25.79% to 4.97%.

Conclusion: In this article, an attempt was made to provide a new method for data cleaning. In order to achieve the desired goal, especially the data of provinces and cities were cleaned using the prepared standard data. By implementing the proposed method, we were able to identify the contaminated data in this characteristic in the examination of the students' national code and then implement the data correction process on them.

Keywords: Data Preparation, Data warehouse, Data Quality, Data Cleaning, Data Warehouse, Dirty Data, Data Management

ارائه روشی جدید برای پاکسازی داده‌ها جهت بهبود کیفیت انبار داده

دوره اول، زمستان ۱۳۹۹
شماره دوم، صص: ۳۳ - ۴۱

تاریخ دریافت: ۱۳۹۹/۰۸/۱۱
تاریخ پذیرش: ۱۳۹۹/۱۰/۲۳

علی شهناز^۱، مهدی افصلی^{۲*}، شیما رحیم‌زاده^۳

۱. استادیار، گروه ریاضی و آمار، واحد زنجان، دانشگاه آزاد اسلامی، زنجان، ایران. shahnavaz_ali2000@yahoo.com
۲. استادیار، گروه مهندسی فن‌آوری اطلاعات، واحد زنجان، دانشگاه آزاد اسلامی، زنجان، ایران. (نویسنده مسئول) afzali@iauz.ac.ir
۳. کارشناسی ارشد، دانشگاه علوم پزشکی و خدمات درمانی و بهداشتی زنجان، زنجان، ایران.

چکیده: مهمترین مسئله در مدیریت داده‌ها، موضوع کیفیت داده است. کیفیت داده می‌تواند پاکسازی داده‌ها را قبل از بارگذاری به انبار داده‌ها تضمین کند. پاکسازی داده فعالیتی است شامل فرآیند تشخیص و اصلاح اشتباهات و تناقضات در انبار داده‌ها. به دلیل وجود اطلاعات زیاد در بانک‌های اطلاعاتی مشکلات و تناقضات فراوانی در آن‌ها به وجود آمده است. هدف اصلی ما ارائه روشی برای رفع تناقضات موجود در بانک‌های اطلاعاتی برای پاکسازی داده‌های آلوده می‌باشد. با هدف بهبود کیفیت انبار داده برای تصمیم‌گیری‌های صحیح، روش جدیدی ارائه شده است و برای آزمایش روش پیشنهادی، از بانک اطلاعاتی شناسنامه سلامت دانشجویان دانشگاه علوم پزشکی زنجان ورودی سال‌های ۹۲ و ۹۳، شامل ۸۴۵ نفر که در حال حاضر همه آن‌ها فارغ‌التحصیل شده‌اند به عنوان داده‌های مورد بررسی استفاده شده است. برنامه پیشنهادی با زبان برنامه‌نویسی سی‌شارپ پیاده‌سازی و اجرا شده است. برنامه یا اپلیکیشن ما در چهار لایه و به صورت ویندوز اپلیکیشن نوشته شده است. از طریق اجرای روش پیشنهادی توانستیم با بررسی کدملی دانشجویان، داده‌های آلوده در این مشخصه را تشخیص داده و سپس فرآیند اصلاح داده را روی آن‌ها اعمال نماییم. براساس نتایج به دست آمده، میزان داده آلوده در انبار داده تولیدشده از ۲۵٫۷۹ درصد به ۴٫۹۷ درصد کاهش یافت.

واژه‌های کلیدی: مدیریت داده، آماده‌سازی، انبار داده، داده‌کاو، داده‌های آلوده، پاک‌سازی.

۱. مقدمه

فن آوری‌های نوین با سرعتی سرسام‌آور در حال پیشرفت هستند، آن گونه که جوامع به صورت عام و بازار به صورت خاص با شتابی وصف‌ناپذیر به دنبال ترندهایی می‌گردند که بقایشان را در این عرصه آشفته و متلاطم تضمین کنند. سازمان‌ها باید بپذیرند که فلسفه حیاتشان تغییر کرده است و دیگر زنده‌بودن به معنای رسیدن به وضعیت سوددهی مداوم نمی‌تواند باشد و باید به دنبال رقابت و ابزار آن بود. یک سازمان در طول حیاتش، داده ایجاد می‌کند و این حقیقت سازمان‌ها را ملزم به جستجوی ابزارهایی برای تسهیل فرآیند کسب اثربخش داده‌ها، پردازش و تحلیل وسیع آن‌ها کرده است تا براساس آن پایه‌ای را برای کشف دانش جدید بنا نهند. بنابراین تسلط بر فن آوری‌های جدیدی مانند هوش تجاری^۱ در کسب-وکارها یک الزام و ضرورت اجتناب‌ناپذیر تلقی می‌شود.

هوش تجاری به‌عنوان فرآیند تبدیل داده به اطلاعات و سپس دانش تعریف می‌شود [1]. هوش تجاری نشان‌دهنده فضای گسترده کاربرد و فن آوری برای جمع‌آوری، ذخیره‌سازی، تجزیه و تحلیل و دسترسی به اطلاعات برای بهبود مدل‌سازی فرآیند تجاری باکیفیت است.

مراحل هوش تجاری شامل مراحل اول، مرحله اول، منابع داده جمع‌آوری می‌شوند. این منابع می‌تواند داده‌های انواع پایگاه داده یا اطلاعات نرم‌افزارهای موجود را دربرگیرد. اطلاعات جمع‌آوری شده طی فرآیند "ای تی ال"^۲ در پایگاه داده تحلیلی یا همان انبارداده بارگذاری می‌شود.

رالف کیمبال انبار داده‌ها را به‌عنوان "یک کپی از سیستم‌های معاملاتی، به‌ویژه ساختار برای پرس‌وجو و تجزیه و تحلیل" تعریف می‌کند. داده در پایگاه داده تحلیلی در بخش‌های مجزایی به نام انبارک^۳ قرار می‌گیرد. در انبار داده‌ها کل داده‌های جمع‌آوری شده وجود دارد و انبارک‌ها شامل تنها یک زیرمجموعه از حجم داده‌های سازمان، مختص به یک گروه خاصی از کاربران، در موضوعی خاص محدود شده است [2]. انبارداده‌ها برای اهداف تحلیلی یعنی پردازش تحلیلی آنلاین "اولپ"^۴ استفاده می‌شوند، انبار داده موضوع‌گرا، یکپارچه، غیرفرآر و مغایر زمان جمع‌آوری داده‌ها در حمایت از مدیریت تصمیم‌گیری است [3]. انبار داده‌ها به همراه "ای تی ال" و ابزار گزارش، یک محیط یکپارچه برای پردازش کسب‌وکار ایجاد می‌کند [4]. در مرحله بعد هوش تجاری وارد عمل شده و اطلاعات طبقه‌بندی شده را تجزیه و تحلیل می‌کند. در نهایت اطلاعات جهت انتشار به ابزارهای سطح بالا تحویل داده می‌شود [2].

بیش از ۷۰٪ از پروژه‌های انبار داده‌ها در بخش «ای تی ال» صرف می‌شود [5]. امروزه سناریوی ابزار «ای تی ال» قسمت مهمی از مسئولیت نرم‌افزار برای یکپارچه‌سازی اطلاعات ناهمگون از منابع مختلف است [6,7]. انجام فرآیند «ای تی ال» بالقوه پیچیده، سخت و وقت‌گیر است.

شوران (۲۰۱۴) کل مشکلات موجود در داده‌های انبار داده را براساس کیفیت داده دسته‌بندی کرده است [13]. چودھاری در سال ۲۰۱۴ مشکلات و روش‌های پاکسازی داده‌ها را مطرح کرد. [7] میلانی

و دکتر گوپتا (۲۰۱۵) نیز مشکلات کیفیت داده‌ها در مرحله "ای تی ال" را نام‌بردند [9]. این پژوهش‌ها نشان از اهمیت مشکلات موجود است و همچنان تحقیقات روی دسته‌بندی و معرفی این مشکلات ادامه دارد.

تاکنون کارهای فراوانی جهت پاک‌سازی داده‌ها صورت گرفته است، اعم از رفع غلط املایی و آوایی، یکی کردن فرمت فیلد تاریخ، اصلاح داده‌های کمی، حذف رکورد تکراری و غیره. ولی براساس پژوهش انجام‌گرفته، برای رفع تناقض در فیلدهای وابسته به هم (مانند: محل‌های جغرافیایی) کاری انجام‌نشده است و چون این فیلد در بانک اطلاعاتی سازمان‌ها موجود است، بنابراین رفع آن برای بهبود عملکرد سازمان‌ها مسئله‌ای مهم و ضروری است که باید به آن پرداخته شود. کارهای انجام‌گرفته همیشه از طریق دو فیلد تشخیص و اصلاح خطا انجام شده است در این پژوهش برآنیم که از طریق سه فیلد وابسته به هم، این کار را انجام دهیم. هدف اصلی ارائه روش جدیدی جهت بالابردن کیفیت داده از طریق رفع تناقض در فیلدهای وابسته به هم و پرکردن فیلدهای خالی -در صورت وجود- و حذف رکوردهای تکراری در بانک اطلاعاتی است.

۲. پیشینه پژوهش

وارول و همکاران سال (۲۰۱۰) الگوریتم PNR^۵ را ارائه دادند. این الگوریتم برای تصحیح خطاهای متنی و آوایی، از فرهنگ لغت جهانی استفاده کرده و به پاک‌سازی داده می‌پردازد. مشکل روش یادشده این منحصربودن فرهنگ لغت جهانی به زبان انگلیسی است [14].

لی و همکاران (۲۰۱۰) الگوریتم بستر متعددی^۶ را ارائه دادند، که در آن برای حذف رکوردهای تکراری و پرکردن مقادیر فیلدهای خالی از کلید برای تطبیق استفاده می‌شود اما مشکل این است که تنها یک کلید تطبیق استفاده می‌شود و به‌صورت کامل هم خودکار نیست [15].

هماد و جیهاد (۲۰۱۱) تکنیک‌های پیشرفته^۷ برای پاک‌سازی داده را ارائه دادند که جهت تصحیح داده‌های کمی استفاده می‌شود این تکنیک در اصلاح داده‌هایی مثل سن و تاریخ تولد به‌کار می‌رود [16].

پائول و همکاران در سال ۲۰۱۲ روش ترکیبی HADCLEAN را برای پاک‌سازی داده در انبار داده‌ها ارائه دادند و در آن از ترکیب اصلاح-شده نسخه‌های الگوریتم PNR^۵ و الگوریتم بستر متعددی استفاده کردند. در روش مذکور به جای فرهنگ لغت جهانی از فرهنگ اصطلاحات مخصوص سازمان استفاده و هجده رکورد بیشتر از الگوریتم اصلی PNR^۵ اصلاح شده است [17].

پوروال و وورا در سال ۲۰۱۳ دو الگوریتم قوانین انجمنی^۸، HADCLEAN و روش‌های آن‌ها برای کیفیت داده‌ها را شرح داده‌اند. این تحقیق همچنین شامل مقایسه عوامل مختلف و جنبه‌های مشترک در این دو الگوریتم است. الگوریتم قوانین انجمنی با استفاده از محاسبات سنگین ریاضی خطای دورویی از نوع رشته داده (فیلد نام) را می‌یابد [12].

کولکارتی و باکال در سال ۲۰۱۴ هم روش ترکیبی HADCLEAN را برای پاک‌سازی داده در انبار داده‌ها دوباره مطرح کرده و شرح داده‌اند

[18]. آشوبنی و همکاران نیز (۲۰۱۴) از تکنیک‌های ترکیبی برای پاک‌سازی داده‌ها استفاده کرده و با ترکیب سه الگوریتم PNRs، تکنیک‌های پیشرفته، و بستر متعددی توانسته‌اند چندین فیلد را با هم پاک‌سازی کنند [19].

باتاچارجی و همکاران نیز سال ۲۰۱۴ روی پاک‌سازی بانک اطلاعاتی براساس «ای تی ال» تحقیق کردند. آن‌ها با اعمال الگوریتم‌های نسبت به موفق، حذف رکوردهای تکراری، ادغام جداول به ترتیب باعث کاهش خطا در بانک اطلاعاتی شدند که در نهایت بعد از اعمال الگوریتم، ۳٪ از کل خطا، باقی‌ماند [8].

تانجا و همکاران در سال ۲۰۱۴ الگوریتم نوین DFT^9 (تبدیل فیلد تاریخ) را برای پاک‌سازی داده‌ها ارائه کردند. این الگوریتم که با استفاده از جاوا پیاده‌سازی شده است و مشکلات ذکر شده در فیلد تاریخ را اصلاح می‌کند [10].

دکتر کالیا و دوی (۲۰۱۵) به معرفی و مقایسه ابزارهای پاک‌سازی داده اعم از: RAPIDMINOR-WINPURE CLEAN & MATCH - MS EXCEL DATA CLEANER پرداختند [11].

۳. روش‌شناسی پژوهش

استفاده از فرهنگ استاندارد داده‌های استان‌ها و شهرستان‌های ایران ضروری است و برای این کار ابتدا یک بانک اطلاعاتی شامل اطلاعات استان‌ها و شهرستان‌های ایران و سه رقم اول کد ملی مخصوص هر شهرستان، در محیط اکسس ۲۰۱۰ ساخته شد و ارتباط میان جداول مربوط برقرار شد. سپس بانک اطلاعاتی حاصل از محیط اکسس به اس کیو ال سرور ۲۰۱۴ انتقال داده شد (این کار به صورت پویا^{۱۱} انجام شده و قابلیت اضافه کردن اطلاعات جدید و یا اصلاح اطلاعات موجود و همچنین امکان حذف از بانک اطلاعاتی در مواقع خاص در آن وجود دارد).

بانک اطلاعاتی شناسنامه سلامت دانشجویان دانشگاه علوم پزشکی زنجان ورودی سال‌های ۹۲ و ۹۳ که فارغ‌التحصیل شده‌اند به عنوان داده‌های مورد بررسی از بخش فن‌آوری اطلاعات دانشگاه علوم پزشکی زنجان دریافت شد. این داده‌ها براساس ترم‌های مختلف برای دو سال متوالی بود که پس از دریافت، در محیط اکسل ۲۰۱۰ یکپارچه گردید و مجموعاً شامل ۸۴۵ رکورد شد. سپس بانک اطلاعاتی دانشجویان از محیط اکسل به اس کیو ال سرور ۲۰۱۴ انتقال داده شد. لازم به ذکر است که برنامه پیشنهادی با زبان برنامه‌نویسی سی شارپ پیاده‌سازی و اجرا شده است (دلیل استفاده از سی شارپ این است که دارای زبان مفهومی تری برای کاربر است و در آن از دات نت فریم ورک^{۱۱} که زیربنای زبان‌های برنامه‌نویسی است، استفاده شده است) (دات نت فریم ورک تمام لایه‌های پیاده‌سازی نرم‌افزار را از سطح سیستم‌عامل به بالا، پوشش می‌دهد و فریم ورک مجموعه‌ای از فایل‌های مورد نیاز سیستم‌عامل است که برای اجرای برنامه نوشته شده تحت دات نت ضروری است). برنامه پیشنهادی در محیط ویژوال استودیو^{۱۲} ۲۰۱۳ نوشته شده و دارای دو بانک اطلاعاتی یکی فرهنگ استاندارد استان‌ها و شهرها و سه رقم کد

مربوط به آن‌ها و دیگری بانک اطلاعاتی دانشجویی برای انجام آزمایش و رفع تناقضات داده‌هاست، هر دو بانک اطلاعاتی در اس کیو ال سرور ۲۰۱۴ هستند (دلیل استفاده از اس کیو ال سرور هم این است که این نسخه تقریباً با بانک اطلاعاتی اوراکل برابری می‌کند و در میلیاردها رکورد، بسیار سریع پاسخ می‌دهد). سپس ارتباط برنامه با این دو بانک اطلاعاتی از طریق دیتاست انجام می‌شود، یعنی برای ارتباط از کلاس دیتاست استفاده شده است. برنامه یا اپلیکیشن ما در چهار لایه و به صورت ویندوز اپلیکیشن^{۱۳} نوشته شده است:

(۱) لایه ارتباط با بانک که جهت ارتباط اپلیکیشن با بانک اطلاعاتی استفاده می‌شود.

(۲) لایه دستورات جستجو و غیره، اسکریپت‌های مختلف جستجو در این لایه نوشته می‌شود. اسکریپت‌ها می‌توانند مستقیماً در محیط برنامه نوشته یا به صورت store procedure از بانک اطلاعاتی فراخوانی شوند.

(۳) لایه ارتباط لایه اول و دوم؛

(۴) لایه دستورات زبان برنامه‌نویسی، که بخش اصلی برنامه بوده و طراحی اینترفیس^{۱۴} برنامه هم در این لایه انجام می‌گیرد.

در برنامه پیشنهادی از کامپوننت Janus Winforms Controls v4 جهت زیباسازی خروجی برنامه استفاده شده است.

۱.۳. شرح مراحل برنامه

ابتدا کلیه اطلاعات فرهنگ استاندارد استان‌ها و شهرستان‌ها و سه رقم کد ملی هر شهرستان خوانده و در یک دیتاست ذخیره می‌شود. سپس کلیه اطلاعات بانک اطلاعاتی دانشجویان خوانده می‌شود و در دیتاست دیگری ذخیره می‌شود که حاصل کار اطلاعات زیر است: استان، شهرستان و کد ملی.

در مرحله بعد، برنامه می‌سنجد که آیا فیلد کد ملی دارای مقدار می‌باشد یا مقدار فیلد کد ملی null و یا blank است؟ اگر خیر آن سطر از جدول دانشجویان را در جدول داده‌های فاقد کد ملی ذخیره می‌کند و در خروجی نمایش می‌دهد. اگر بله به مرحله بعدی می‌رود و ده تایی بودن تعداد کاراکترهای فیلد کد ملی را می‌سنجد. اگر بله بود، ایده به-کاررفته در این مرحله استفاده از جدولی به نام temp است. جدول موقتی که در اس کیو ال سرور ساخته شده داده‌های آلوده کشف شده در این بخش را از جدول اصلی دانشجویان می‌خواند و آن‌ها را در جدول temp کپی می‌کند و در خروجی نمایش می‌دهد. این جدول موقت در هر مرحله و هر بار که فراخوانی می‌شود، ابتدا همه داده‌های خود را پاک می‌کند و از اول آن را بسته به دستور خواسته شده از داده‌های آلوده پرمی‌کند و نمایش می‌دهد. حال اگر تعداد کاراکتر کد ملی مطابق انتظار ده بود و درست بود به مرحله بعدی می‌رود و می‌سنجد که فرمت فیلد کد ملی درست باشد. اگر چنین بود، باز هم از جدول temp استفاده می‌کند و اطلاعات قبلی را پاک کرده و داده‌های آلوده جدید را در آن می‌ریزد و نمایش می‌دهد.

اما خطای کد ملی چگونه تشخیص داده می‌شود؟ کد ملی شماره‌ای است ده رقمی که از سمت چپ سه رقم کد شهرستان محل صدور

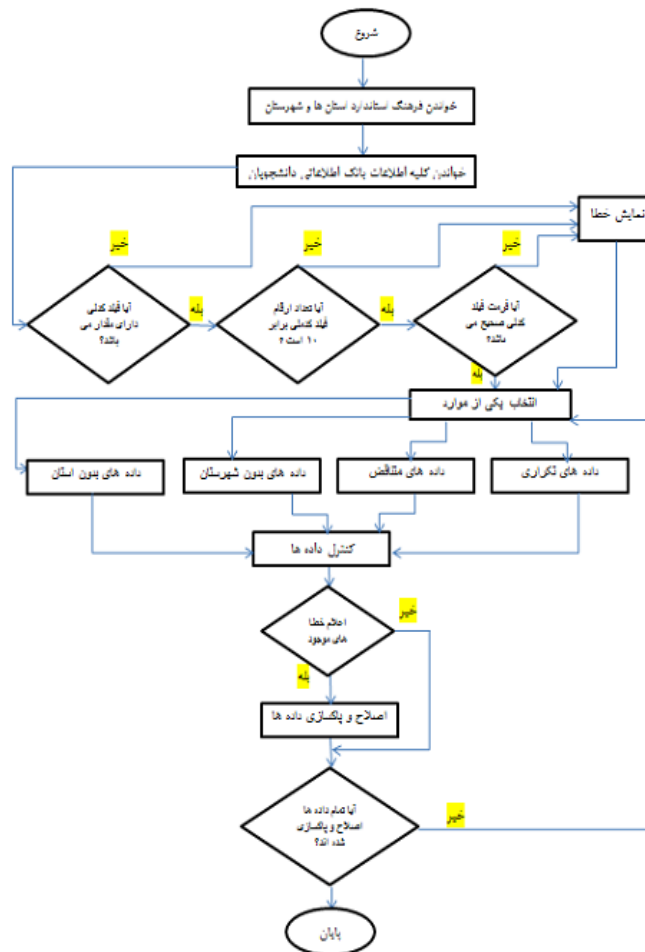
شناسنامه، شش رقم بعدی کد منحصر به فرد برای فرد دارنده شناسنامه در شهرستان محل صدور و رقم آخر هم یک رقم کنترل است که از روی نه رقم سمت چپ به دست می‌آید. برای بررسی کنترل کد کافی است مجدداً از روی نه رقم سمت چپ رقم کنترل را محاسبه کنیم. برای محاسبه رقم کنترل از روی سایر ارقام:

- هر رقم را در موقعیت آن ضرب کرده و حاصل را با هم جمع می‌کنیم.
- مجموع به دست آمده از مرحله یک را بر عدد یازده تقسیم می‌کنیم.
- اگر باقیمانده کمتر از عدد دو باشد، رقم کنترل باید برابر باقیمانده باشد در غیر این صورت، رقم کنترل باید برابر با یازده منهای باقیمانده باشد.

این فرمول کلی در برنامه برای تک تک کدملی‌ها محاسبه می‌گردد و به وسیله آن داده‌های آلوده این بخش مشخص می‌گردد.

در مرحله بعد داده‌های تکراری، داده‌های بدون استان، داده‌های بدون شهرستان و داده‌هایی که استان و شهرستان آن‌ها با هم تناقض دارند، کنترل و محاسبه می‌شوند.

ایده به کار رفته برای کشف تکرار چنین است: کدملی‌هایی که کاراکتر به کاراکتر عین هم باشند را می‌یابد و باز هم با استفاده از جدول temp آن‌ها را نمایش می‌دهد. برای سه قسمت بعدی در قسمت view اس کیو ال سرور یک جدول جدید از ترکیب جداول فرهنگ استاندارد استان‌ها و شهرستان‌ها ساخته‌ایم که نام استان‌ها و شهرستان‌ها و سه رقم کد مربوط به آن شهرستان‌هاست و داده‌های آلوده استان و شهرستان با توجه به این جدول ترکیبی ساخته شده استخراج می‌شود، جدول موقت temp با این خطاها پر شده و در خروجی نشان داده می‌شود.



شکل ۱: فلوچارت برنامه پیشنهادی

پس از کشف این داده‌های آلوده وارد مرحله بعد می‌شویم و اصلاح و پاک‌سازی داده‌ها آغاز می‌شود. تا قبل از این مرحله خطاها فقط تشخیص داده می‌شدند، پس از این مرحله وارد فاز پاک‌سازی داده‌ها می‌شویم، اصلاح داده‌های آلوده با استفاده از فرهنگ استاندارد استان‌ها و شهرستان‌ها آغاز می‌شود و داده‌های تکراری نیز از بین می‌روند و فقط یک نسخه از آن‌ها باقی می‌ماند. در این برنامه همچنین اگر داده‌ها دارای غلط املائی باشند به عنوان داده آلوده شناسایی و با توجه به فرهنگ استاندارد تهیه شده، اصلاح می‌شوند. از آنجاکه فرهنگ استاندارد طراحی و ساخت خود ماست، عاری از هرگونه غلط املائی است و می‌تواند این نوع از خطاها را هم تشخیص داده و اصلاح کند. لازم به ذکر است که در هر مرحله خطاها با ذکر تعداد در خروجی نمایش داده می‌شود.

غلط املائی باشند به عنوان داده آلوده شناسایی و با توجه به فرهنگ استاندارد تهیه شده، اصلاح می‌شوند. از آنجاکه فرهنگ استاندارد طراحی و ساخت خود ماست، عاری از هرگونه غلط املائی است و می‌تواند این نوع از خطاها را هم تشخیص داده و اصلاح کند. لازم به ذکر است که در هر مرحله خطاها با ذکر تعداد در خروجی نمایش داده می‌شود.

برنامه پس از اتمام مرحله پاک‌سازی داده، در مرحله بعد می‌سند که همه داده‌ها از جدول خوانده شده باشند که در این صورت، کار به اتمام می‌رسد و داده‌های اصلاح شده نمایش داده می‌شوند.

۲.۳. شرح روند اجرایی و خروجی‌های برنامه

برنامه پیشنهادی یک صفحه اصلی و سه منو دارد: (۱) بانک استان‌ها (۲) بانک افراد (۳) تحلیل داده‌ها.

(۱) بانک استان‌ها: در این بخش، بانک اطلاعاتی فرهنگ استاندارد استان‌ها و شهرها به نمایش درآمده است که در فرم ابتدایی آن نام تمامی سی و دو استان ایران ثبت شده است و اگر وارد ثبت شهرستان هر کدام از استان‌ها شویم، شهرستان‌های آن استان را در فرم بعدی به صورت کامل نمایش می‌دهد. تعداد کل شهرستان‌های ثبت شده ۵۱۴ شهرستان است. اگر وارد ثبت کدهای شهرستان هر استان شویم، سه رقم اول کد ملی مخصوص آن شهرستان را در فرم بعدی می‌بینیم. هر شهرستان ممکن است دارای یک یا چندین کد شهرستان باشد، به عنوان مثال شهرستان اسلامشهر دارای یک کد شهرستان است ولی تهران دارای پانزده کد شهرستان می‌باشد. که ایده به کار رفته در کدنویسی برنامه برای جدا کردن کدهای شهرستان استفاده از کاراکتر * (ستاره) است. تعداد کل کد شهرستان‌های ثبت شده ۶۵۵ مورد است. لازم به ذکر است که فرهنگ استان و شهرستان تهیه شده براساس آخرین تقسیم‌بندی جغرافیایی کشور ایران است.

در این فرم‌ها هر کدام از استان‌ها و شهرستان‌ها و کدهای آن‌ها دارای قسمت‌های حذف و ویرایش و اضافه کردن مورد است یعنی به صورت پویا، اگر زمانی تقسیم‌بندی کشور دچار هرگونه تغییری شود، امکان اعمال آن در فرهنگ استاندارد مذکور تعبیه شده باشد و نیز اگر به هر دلیلی بانک اطلاعاتی دچار اختلال شد، اضافه، ویرایش و حذف رکوردهای آن امکان پذیر باشد.

استانی فاقد شهرستان باشد، برای اینکه پیغام خطا ندهد و از اجرای برنامه خارج نشود، از حلقه try, catch استفاده شده است. وقتی در بخش‌های انتخابی استان و شهرستانی را انتخاب می‌کنیم، تمام داده‌هایی را که استان و شهرستان آن‌ها درست است، طبق انتخاب ما از بانک دانشجویان می‌خواند و در جدولی با مشخص کردن تعداد رکوردهای موجود نمایش می‌دهد. این بخش برنامه هم دارای گزینه ویرایش و حذف می‌باشد.

کد ملی	نام شهرستان	حسبیت	فرص	تاریخ	مکان
۰۰۰۰۰۰۰۰	تهران	دقیق	تهران	۱۳۹۹/۰۱/۰۱	تهران
۰۰۰۰۰۰۰۰	تهران	دقیق	تهران	۱۳۹۹/۰۱/۰۱	تهران
۰۰۰۰۰۰۰۰	تهران	دقیق	تهران	۱۳۹۹/۰۱/۰۱	تهران

شکل ۳: خروجی برنامه در بخش بانک افراد

(۳) تحلیل داده‌ها: این قسمت دارای هشت قسمت انتخابی می‌باشد شامل نمایش کلیه داده‌ها، نمایش داده‌های تکراری، داده‌های بدون کد ملی، داده‌های با فرمت نادرست کد ملی، داده‌های با تعداد ارقام نادرست کد ملی، داده‌های بدون استان، داده‌های بدون شهرستان و داده‌های دارای استان و شهرستان متناقض، که براساس انتخاب ما در خروجی داده‌های مربوطه را به همراه تعدادشان نمایش می‌دهد. پس از کشف داده‌های آلوده از طریق دکمه اصلاح، می‌توان خطاها را اصلاح نمود.

کد ملی	نام شهرستان	حسبیت	فرص	تاریخ	مکان
۰۰۰۰۰۰۰۰	تهران	دقیق	تهران	۱۳۹۹/۰۱/۰۱	تهران
۰۰۰۰۰۰۰۰	تهران	دقیق	تهران	۱۳۹۹/۰۱/۰۱	تهران
۰۰۰۰۰۰۰۰	تهران	دقیق	تهران	۱۳۹۹/۰۱/۰۱	تهران

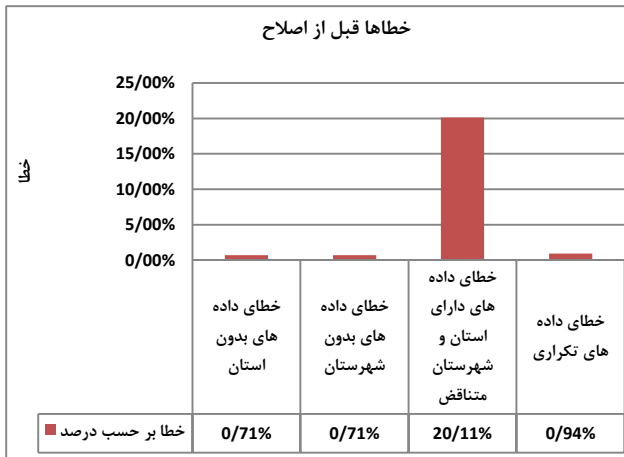
شکل ۴: خروجی برنامه در قسمت تحلیل داده‌ها



شکل ۲: خروجی برنامه در بخش بانک استان‌ها

۴. پیاده‌سازی روش پیشنهادی

(۲) بانک افراد: قسمت بانک افراد دارای دو بخش انتخابی برای استان و شهرستان است. ایده به کار رفته در این بخش چنین است که اگر زمانی



شکل شماره ۶: نمودار خطای کشف شده در برنامه قبل از اصلاح

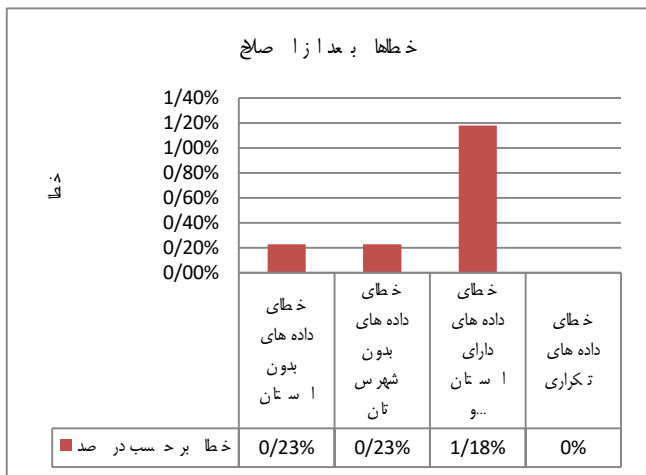
۳.۴. مرحله پاکسازی داده‌ها

پس از اتمام فاز تشخیص خطا در برنامه نوبت به پاکسازی این داده‌های آلوده می‌رسد. کل فرآیند پاکسازی داده‌ها در برنامه پیشنهادی ده ثانیه زمان می‌برد. و تمامی داده‌های آلوده شناسایی شده در فاز دوم تشخیص، به‌طور کامل اصلاح می‌شوند و پس از آن داده‌های اصلاح شده را به ما نمایش می‌دهد.

چون روش پیشنهادی ما شامل اصلاح خطاهای مربوط به صحت کد ملی نیست، به‌همین دلیل پس از مرحله اصلاح تمام خطاها به صفر درصد نمی‌رسند. و ما به‌طور دقیق تاثیر این خطاها بر درصد اصلاح را محاسبه نموده‌ایم، که در داده‌های بدون استان و شهرستان دو رکورد دارای خطای صحت کد ملی هستند، پس خطا از ۰,۷۱ درصد به ۰,۳۵ درصد کاهش می‌یابد.

در داده‌هایی که دارای استان و شهرستان متناقض هستند، تعداد ده خطای صحت کد ملی داریم که با این حساب خطا از ۲۰,۱۱ درصد به ۱,۱۸ درصد می‌رسد.

در داده‌های تکراری چون خطای صحت کد ملی وجود ندارد، خطا از ۰,۹۴ درصد به صفر درصد می‌رسد.



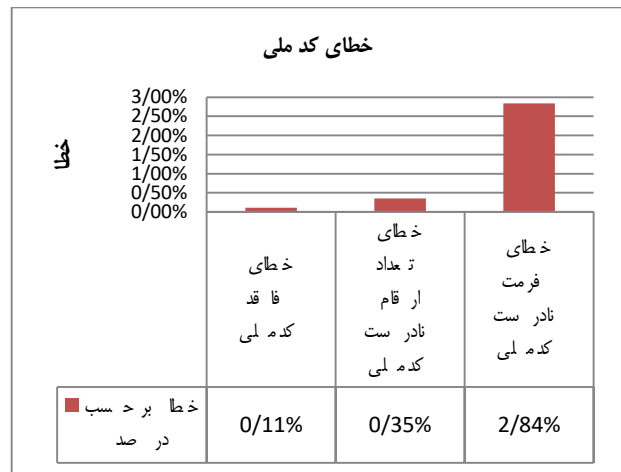
شکل ۷: نمودار خطای باقی‌مانده پس از اصلاح و پاکسازی داده‌ها

در این مرحله پیاده‌سازی برنامه آغاز می‌شود و همان‌طور که قبلاً گفته شد روش پیشنهادی را روی جامعه آماری یعنی بانک اطلاعاتی مربوط به شناسنامه سلامت دانشجویان، اعمال می‌کنیم. داده‌های آلوده به‌طور کامل شناسایی می‌شوند و تعداد داده‌های آلوده مربوط به هر بخش در برنامه مشخص می‌شود.

خطاهای فاز اول پیاده‌سازی برنامه شامل این سه نوع خطای زیر است:

۱.۴. خطاهای مرتبط با کد ملی

بانک اطلاعاتی دانشجویان تنها دارای یک رکورد فاقد کد ملی می‌باشد، بر این اساس ۰,۱۱ درصد از خطا را به خود اختصاص داده‌است. تعداد سه رکورد وجود دارد که دارای تعداد ارقام نادرست فیلد کد ملی می‌باشند، لذا ۰,۳۵ درصد از خطا مربوط به این نوع از داده‌های آلوده می‌باشد. تعداد ۲۴ رکورد دارای فرمت نادرست کد ملی برطبق فرمول‌بندی تشخیص صحت کد ملی هستند که ۲,۴۸ درصد از خطا را به خود اختصاص داده‌اند.



شکل ۵: نمودار خطاهای مربوط به کد ملی

۲.۴. خطاهای مربوط به داده‌های استان‌ها و شهرستان‌ها

داده‌هایی که فاقد استان می‌باشند، شش مورد کشف شده‌اند و ۰,۷۱ درصد از خطا مربوط به این نوع از داده‌هاست. داده‌هایی که فاقد شهرستان می‌باشند، هم شش مورد است که باز هم ۰,۷۱ درصد از خطا را به خود اختصاص داده‌اند. داده‌هایی که استان و شهرستان آن‌ها با هم در تناقض‌اند، تعداد ۱۷۰ داده است و ۲۰,۱۱ درصد خطا مربوط به این نوع از داده‌ها می‌باشد. داده‌های تکراری هشت مورد است و ۰,۹۴ درصد از خطا را به خود اختصاص داده‌است.

به این ترتیب فاز دوم تشخیص خطا هم به پایان می‌رسد.

۵. نتیجه‌گیری

در این مقاله تلاش بر این بود تا روش جدیدی برای پاکسازی داده‌ها ارائه شود. در راستای دستیابی به هدف موردنظر، به‌طور خاص پاکسازی روی داده‌های استان‌ها و شهرستان‌های با به‌کارگیری فرهنگ استاندارد تهیه شده انجام شد. برای آنکه نتایج واقعی باشد از اطلاعات دانشجویان دانشگاه علوم پزشکی زنجان استفاده و فرآیند پاکسازی روی آن اعمال گردید. در این پژوهش سه فیلد وابسته به هم را در نظر گرفته و همچنین صحت فیلد کد ملی را بررسی کرده‌ایم و در پایان خطای حاصل از آن را محاسبه کرده‌ایم. روش ارائه شده از دقت مناسبی برخوردار است، و به میزان قابل توجهی میزان داده‌های آلوده را کاهش داده است. هدف اصلی بالابردن کیفیت داده از طریق رفع تناقض در فیلدهای وابسته به هم، پرکردن فیلدهای خالی احتمالی و حذف رکوردهای تکراری در بانک اطلاعاتی بود و روش پیشنهادی توانست ما را به هدف اصلی برساند و سطح کیفیت مطلوبی که در ابتدای کار از ارائه این روش انتظار داشتیم کاملاً تضمین کند.

با اجرای روش پیشنهادی توانستیم در بررسی کد ملی دانشجویان، داده‌های آلوده در این مشخصه را تشخیص داده و سپس فرآیند اصلاح داده را روی آن‌ها پیاده کنیم. با اعمال روش پیشنهادی و براساس نتایج به دست آمده، میزان داده آلوده در انبار داده تولیدشده از ۲۵٫۷۹٪ به ۴٫۹۷٪ کاهش یافت.

مراجع

- [1] Golfarelli, Matteo "New Trends in Business Intelligence". Proceedings of the 28th International Convention MIPRO (BIS&DE&ISS). MIPRO (May 30-June /2005), Opatija, Croatia. PP: (15-20)
- [2] NEDELICU, Bogdan. "Business Intelligence Systems". Database Systems Journal. Vol IV. no.4(2013). PP: (12-20)
- [3] Inmon, William H. "Building the Data Warehouse". Wiley Publishing. 4th Edition. (2005). P:32
- [4] Ghosh, Ranak; Halder, Sujay; Sen, Soumya. "An Integrated Approach to Deploy Data Warehouse in Business Intelligence Environment". IEEE Third International Conference. (7-8 Feb- 2015), Hooghly). PP: (1 – 4)
- [5] Talebzadeh, Hossein. A Service-Based Framework for ETL Process Based on Metadata. Journal of Basic and Applied Scientific Research. (2/1/ 2012). PP:(54-59)
- [6] Gill, Rupali; Singh, Jaiteg. "A Review of Contemporary Data Quality Issues in Data Warehouse ETL Environment". Journal on Today's Ideas Tomorrow's Technologies. Vol. 2. No. 2. (19 December 2014). PP: (153_160)
- [7] Choudhary, Nidhi. "A Study over Problems and Approaches of Data Cleansing/Cleaning", International Journal of Advanced Research in Computer Science and Software Engineering. Vol. 4, Issue 2. (February 2014). PP: (774 _779)
- [8] Bhattacharjee, Arup Kumar ; Chatterjee, Partha ; Prasad Shaw, Mukesh ; Chakraborty, Manomoy. "ETL based Cleaning on Database", International Journal of Computer Applications. Vol.105, No. 8. (November 2014). PP: (34– 40)
- [9] Miglani, Sakshi; Dr. Gupta, Neha. "An Overview On Evocations Of Data Quality at ETL Stage", International Journal of Advanced Technology in Engineering and

Science. Vol. No. 03. Special Issue No. 01. (March 2015). PP: (1429- 1436)

- [10] Taneja, Shweta; Ashri, Ishita; Gupta, Shipra; Sharma, Mehak. DFT: "A Novel Algorithm for Data Cleansing", International Journal of Computer Science and Information Technologies. Vol..5 (2). (2014). PP: (2297- 2301)
- [11] Devi, Sapna; Dr. Kalia, Arvind. "Study of Data Cleaning & Comparison of Data Cleaning Tools", International Journal of Computer Science and Mobile Computing. Vol. 4. Issue. 3. (March 2015). PP: (360 – 370)
- [12] Porwal, Sonal; Vora, Deepali. "A Comparative Analysis of Data Cleaning Approaches to Dirty Data". International Journal of Computer Applications. Vol. 62, No.17. (January 2013). PP: (30- 34)
- [13] Sheoran, Jyoti. "Issues of Data Quality in Data Warehouses", International Conference on Advances in Computer Engineering & Applications (ICACEA-2014 at IMSEC.GZB). PP: (6 – 8)
- [14] Varol, Cihan; Bayrak, Coskun; Wagner, Rick; Goff, Dana. "Application of the Near Miss Strategy and Edit Distance to handle Dirty Data", International Series in Operations Research & Management Science. Springer US. Vol. 32, (2010). PP: (91 -101)
- [15] Ning Li, Wing; Bheemavaram, Roopa; Zhang, Johnson. "Transitive Closure of Data Records", Application and Computation. International Series in Operations Research & Management Science. Springer US. vol. 132. (2010). PP: (39-75)
- [16] Dr. M. Hamad, Mortadha; Jihad, Alaa Abdulkar. "An Enhanced Technique to clean Data in the Data Warehouse". Developments in E-system Engineering (DeSE). IEEE International Conference on (6-8 Dec. 2011). Dubai. PP: (306-311)
- [17] Paul, Arindam; Ganesan, Varuni; Challa, Jagat Sesh; Sharma, Yashvardhan. "HADCLEAN: A Hybrid Approach to Data Cleaning in Data Warehouses", Information Retrieval & Knowledge Management (CAMP). IEEE International Conference (13-15 March 2012). Kuala Lumpur. PP: (136- 142)
- [18] Kulkarni, Prerna S; Bakal, J.W. "Hybrid Approaches for Data Cleaning in Data Warehouse". International Journal of Computer Applications. Vol 88 , No.18. (February 2014). PP: (7-10)
- [19] M. Save, Ashwini; Kolkur, Seema. "Hybrid Technique for Data Cleaning". National Conference on Role of Engineers (in Nation Building. 2014. NCRENB-14). PP: (4-8)

پی‌نوشت

1. Business Intelligence
2. ETL (Extract, Transform, Load)
3. Data Mart
4. OLAP (On-Line Analytical Processing)
5. Personal Name Recognizing Strategy
6. Transitive Closure Algorithm
7. Enhanced Technique
8. Alliance rules algorithms
9. Date Field Transformation Algorithm
10. dynamic
11. .Net Framework
12. Visual Studio
13. Windows Application
14. interface