

بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ



دانشگاه تربیت دبیر شهید رجائی

## پاکسازی داده‌ها

رفع ناسازگاری، تکرار، مقدار از دست رفته و اغتشاش

تألیف:

دکتر نگین دانشپور

عضو هیأت علمی دانشگاه تربیت دبیر شهید رجائی

مهدیه عطایان

سرشناسه	: دانشپور، نگین، ۱۳۵۵-
عنوان و نام پدیدآور	: پاک‌سازی داده‌ها: رفع ناسازگاری، تکرار، مقدار از دست‌رفته و اغتشاش / تالیف نگین دانشپور، مهدیه عطاییان؛ ویراستار ادبی ساغر سلمانی‌نژاد.
مشخصات نشر	: تهران: دانشگاه تربیت دبیر شهید رجایی، ۱۴۰۰.
مشخصات ظاهری	: ۲۲۰ ص: جدول.
شابک	: ۹۷۸-۶۲۲-۶۵۸۹-۲۰-۸
وضعیت فهرست نویسی	: فیپا
یادداشت	: واژه‌نامه.
یادداشت	: کتابنامه: ص. [۱۸۶] - ۱۸۸.
یادداشت	: نمایه.
موضوع	: داده‌پردازی
موضوع	: Electronic data processing
موضوع	: داده‌کاوی
موضوع	: Data mining
شناسه افزوده	: عطاییان، مهدیه، ۱۳۷۱-
شناسه افزوده	: دانشگاه تربیت دبیر شهید رجایی
شناسه افزوده	: Shahid Rajaei Teacher Training University
رده بندی کنگره	: QA۷۶
رده بندی دیویی	: ۰۰۴
شماره کتابشناسی ملی	: ۸۵۴۱۹۸۸
اطلاعات رکورد کتابشناسی	: فیپا



عنوان	: پاک‌سازی داده‌ها: رفع ناسازگاری، تکرار، مقدار از دست‌رفته و اغتشاش
تألیف	: دکتر نگین دانشپور، عضو هیأت علمی دانشگاه تربیت دبیر شهید رجایی / مهدیه عطاییان
ویراستار ادبی	: دکتر ساغر سلمانی‌نژاد
نوبت چاپ	: اول - زمستان ۱۴۰۰
انتشارات	: دانشگاه تربیت دبیر شهید رجایی
لیتوگرافی، چاپ	: رجاء نقشینه، شریف
طراح جلد	: محمد معتمدی‌نژاد
ناظر چاپ	: محمد معتمدی‌نژاد
صفحه‌آرا	: نیره فیروزی
شمارگان	: ۱۰۰ جلد
قیمت	: ۶۰۰۰۰۰ ریال
شابک	: ۹۷۸-۶۲۲-۶۵۸۹-۲۰-۸
ISBN: 978-622-6589-20-8	

کلیه حقوق این اثر برای مؤلفان و دانشگاه تربیت دبیر شهیدرجایی محفوظ است.

نشانی: تهران، لویزان، کد پستی ۱۶۷۸۸-۱۵۸۱۱، صندوق پستی ۱۶۳ - ۱۶۷۸۵، تلفن: ۹ (۲۶۳۲) - ۲۲۹۷۰۰۶۰،  
 تلفکس: ۲۲۹۷۰۰۴۲، پست الکترونیکی: [publish@sru.ac.ir](mailto:publish@sru.ac.ir)، وب سایت: <http://publish.sru.ac.ir>

## پیش‌گفتار

داده‌های جمع‌آوری شده از منابع مختلف و توزیع شده برای اتخاذ تصمیمات مفید و سودمند باید به اطلاعات و دانش تبدیل شوند. به‌طور سنتی استخراج دانش را تحلیل‌گران داده انجام می‌دهند اما با توجه به رشد روزافزون داده‌ها، نیازمند روش‌هایی مبتنی بر رایانه هستیم. فرآیندهای مبتنی بر رایانه برای کشف دانش را «داده‌کاوی» می‌گویند. کیفیت داده‌های مورد بررسی در هنگام استخراج دانش، اهمیت بسزایی دارد. صحت، کامل بودن و سازگاری از جمله معیارهای مورد بررسی در بحث کیفیت داده‌ها هستند. داده‌های دنیای واقعی ممکن است به دلایل مختلفی دچار مشکل عدم کیفیت شوند. مقادیر جاافتاده، داده‌های مغشوش، ناسازگاری و مقادیر تکراری از جمله مشکلات عمده داده‌ها هستند. از این‌رو پیش‌پردازش که اغلب به منظور رفع این مشکلات انجام می‌شود، یکی از مهم‌ترین مراحل کشف دانش است.

پیش‌پردازش فرآیندی زمان‌بر است اما با توجه به هزینه‌هایی که داده‌های فاقد کیفیت برای سازمان‌ها ایجاد می‌کنند، فرآیندی اجتناب‌ناپذیر است. در صورتی که دانش مورد نیاز برای اتخاذ تصمیمات، از داده‌های فاقد کیفیت استخراج شود، منجر به شکست تصمیمات می‌شود. پاک‌سازی و یکپارچه‌سازی از مراحل پیش‌پردازش هستند. در این کتاب تمرکز بر روی روش‌های جدید پاک‌سازی داده‌ها است. روش‌های مورد بررسی بر روی داده‌های عددی، ترتیبی و اسمی متمرکز هستند. در سال‌های اخیر پاک‌سازی داده‌ها از اهمیت ویژه‌ای برخوردار شده است، زیرا الگوریتم‌های داده‌کاوی برای عملکرد صحیح و ارائه دانش مفید، پیش‌بینی‌ها یا توصیف‌ها، نیاز به داده‌های صحیح، کامل و سازگار دارند. در فرآیند پاک‌سازی داده‌ها هدف، تشخیص داده‌های مغشوش، تکراری، پرت و مقادیر جاافتاده است.

در این کتاب به معرفی مفاهیم اولیه پیش‌پردازش داده‌ها و ابعاد مختلف آن می‌پردازیم. همچنین مشکلات کیفیت داده‌ها به صورت اجمالی بیان می‌شوند و راهکارهای ابتدایی برای پاک‌سازی داده‌ها به‌طور مختصر معرفی می‌شوند. تمرکز اصلی در این کتاب بر روی شرح الگوریتم‌های پیشنهادی نویسندگان در حوزه پاک‌سازی داده‌ها است. کارایی این الگوریتم‌ها به دقت ارزیابی شده و برتری آن‌ها در مقایسه با کارهای معتبر جدید و قابل‌قیاس در آن حوزه تأیید شده است.

به این منظور بر روی این الگوریتم‌ها آزمایش‌های متعدد در شرایط مختلف انجام داده‌ایم. بر اساس آزمایش‌ها و ارزیابی‌های انجام‌شده، این نتیجه حاصل شد که این الگوریتم‌ها نسبت به الگوریتم‌های مشابه عملکرد بهتری دارند و بر این اساس در این کتاب ارائه کردیم.

با توجه به این که پاک‌سازی داده‌ها امری اجتناب‌ناپذیر در تحلیل داده‌ها است و در سال‌های اخیر این امر از اهمیت ویژه‌ای برخوردار شده است، هدف اصلی از ویرایش این کتاب ارائه یک نسخه فارسی متمرکز در این حوزه است. اگرچه منابع انگلیسی متعددی در زمینه پیش‌پردازش و کیفیت داده‌ها موجود است، اما منبع فارسی متمرکزی در حوزه پاک‌سازی داده‌ها موجود نیست. این کتاب می‌تواند برای محققان یا تحلیلگران داده در دانشگاه‌ها و مراکز تحقیقاتی، مفید واقع شود.

## فهرست مطالب

۱	پیشگفتار
۱	<b>فصل ۱</b>
۱	<b>پیش‌پردازش داده‌ها</b>
۴	۱.۱. کیفیت داده‌ها و ابعاد آن
۵	۱.۱.۱. صحت
۷	۲.۱.۱. کامل بودن
۸	۳.۱.۱. سازگاری
۱۰	۴.۱.۱. ابعاد مرتبط با زمان
۱۲	۵.۱.۱. تفسیرپذیری و قابل اعتماد بودن
۱۲	۲.۱. پاکسازی داده‌ها
۱۲	۱.۲.۱. مقادیر ازدست‌رفته
۱۵	۲.۲.۱. تصحیح داده‌های مغشوش
۱۸	۳.۲.۱. شناسایی داده‌های پرت
۲۷	۴.۲.۱. شناسایی داده‌های تکراری
۳۲	۵.۲.۱. شناسایی تناقضات محدودیت‌های یکپارچگی
۳۴	۶.۲.۱. معیارهای ارزیابی الگوریتم‌های پاکسازی داده
۳۷	<b>فصل ۲</b>
۳۷	<b>شناسایی تناقضات محدودیت‌های یکپارچگی و تعمیر داده‌ها</b>
۴۰	۱.۲. تصحیح خودکار مبتنی بر وابستگی تابعی و سیستم یادگیری مرکب
۴۱	۱.۱.۲. مرحله اول: شناسایی تناقضات وابستگی‌های تابعی با رکوردها

۴۴	۲.۱.۲. مرحله دوم: استخراج مقادیر تکرارشونده
۴۵	۳.۱.۲. مرحله سوم: شناسایی خطا
۵۰	۴.۱.۲. مرحله دوم: تصحیح خطا
۵۲	۲.۲. تصحیح خودکار داده و وابستگی تابعی با الگوریتم اکتشافی
۵۴	۱.۲.۲. تعمیر داده‌ها
۵۷	۲.۲.۲. تعمیر محدودیت‌ها
۵۹	۳.۲.۲. الگوریتم تجربی انتخاب محدودیت‌ها
۶۳	۳.۲. خلاصه
<b>۶۵</b>	<b>فصل ۳</b>
<b>۶۵</b>	<b>جایگذاری مقادیر جافتاده</b>
۶۹	۱.۱.۳. جایگذاری مقادیر جافتاده مبتنی بر خوشه‌بندی و رویکرد ترکیبی
۶۹	۱.۱.۳. مرحله اول: خوشه‌بندی
۷۱	۲.۱.۳. مرحله دوم: جایگذاری مقادیر جافتاده مبتنی بر رگرسیون خطی و نزدیکترین همسایگی
۷۴	۲.۳. جایگذاری مقادیر جافتاده با استفاده از روش‌های مبتنی بر همبستگی
۷۵	۱.۲.۳. مرحله اول: محاسبه اولویت صفات جافتاده
۷۵	۲.۲.۳. مرحله دوم: انتخاب مجموعه داده پایه
۷۷	۳.۲.۳. مرحله سوم: حداکثر نمودن همبستگی بین ویژگی‌ها با بسط دادن مجموعه پایه
۸۱	۴.۲.۳. مرحله چهارم: تخمین مقادیر جافتاده با استفاده از زیرمجموعه‌های یافته شده
۸۳	۵.۲.۳. مرحله پنجم: ترکیب گام‌های حداکثرسازی و تخمین
<b>۸۴</b>	<b>۳.۳. جایگذاری مقادیر جافتاده با استفاده از خوشه‌بندی فازی و معیار اطلاعات متقابل</b>
۸۵	۱.۳.۳. مرحله اول: محاسبه اولویت صفات جافتاده
۸۶	۲.۳.۳. مرحله دوم: خوشه‌بندی
۹۱	۳.۳.۳. مرحله سوم: انتخاب ویژگی در هر خوشه انتخاب شده
۹۲	۴.۳.۳. مرحله چهارم: اعمال مدل رگرسیون بر هریک از خوشه‌ها و صفات انتخاب شده
۹۳	۵.۳.۳. مرحله پنجم: انجام نظرخواهی وزن دار برای تخمین نهایی
<b>۹۳</b>	<b>۴.۳. جایگذاری مقادیر جافتاده با استفاده از درون‌یابی معکوس فاصله وزن دار و رویکردی ترکیبی</b>

۹۴	۱.۴.۳. مرحله اول: تعیین شعاع جستجو با استفاده از خوشه‌بندی $k$ -means
۱۰۶	۲.۴.۳. مرحله دوم: تعیین توان تأثیر همسایگان با الگوریتم جستجوی فاخته
۱۰۸	۳.۴.۳. مرحله سوم: تخمین مقادیر جافتاده
	۵.۳. جایگذاری مقادیر جافتاده با استفاده از شبکه عصبی مدلسازی گروهی داده‌ها و الگوریتم هوش
۱۱۰	هیجانی
۱۱۱	۱.۵.۳. شبکه عصبی مدلسازی گروهی داده‌ها
۱۱۵	۲.۵.۳. الگوریتم هوش هیجانی
۱۱۷	۳.۵.۳. ترکیب شبکه عصبی مدلسازی گروهی داده‌ها و الگوریتم هوش هیجانی
۱۲۰	۶.۳. خلاصه
۱۲۳	<b>فصل ۴</b>
۱۲۳	<b>شناسایی رکوردهای تکراری</b>
۱۲۶	۱.۴. شناسایی رکوردهای تکراری مبتنی بر خوشه‌بندی و نشانه‌سازی
۱۲۶	۱.۱.۴. مرحله اول: انتخاب ویژگی و یکسان‌سازی داده‌ها
۱۲۷	۲.۱.۴. مرحله دوم: محاسبه شباهت نمونه‌ها مبتنی بر کد کردن و خوشه‌بندی داده‌ها
۱۲۹	۳.۱.۴. مرحله سوم: محاسبه فاصله نمونه‌ها مبتنی بر نشانه‌سازی
۱۳۰	۴.۱.۴. مرحله چهارم: شناسایی رکوردهای مشابه
۱۳۲	۱.۲.۴. مرحله اول: انتخاب ویژگی
۱۳۳	۲.۲.۴. مرحله دوم: خوشه‌بندی سلسله مراتبی
۱۳۷	۳.۲.۴. مرحله سوم: مقایسه رکوردها
۱۴۱	۳.۴. خلاصه
۱۴۳	<b>فصل ۵</b>
۱۴۳	<b>تشخیص داده‌های مغشوش</b>
۱۴۶	۱.۵. تشخیص اغتشاش با استفاده از فیلترها
۱۴۷	۱.۱.۵. فیلتر ترکیبی
۱۴۷	۲.۱.۵. فیلتر بخش‌بندی تکرارشونده
۱۴۸	۲.۵. تشخیص اغتشاش مبتنی بر خوشه‌بندی و نزدیک‌ترین همسایگی

۱۴۹	۱.۲.۵. مرحله اول: خوشه‌بندی
۱۴۹	۲.۲.۵. مرحله دوم: مقایسه
۱۵۲	۳.۲.۵. مرحله سوم: تشخیص ویژگی‌های مغشوش
۱۵۲	۴.۲.۵. مرحله چهارم: کاهش خطا
۱۵۲	۳.۵. تشخیص اغتشاش مبتنی بر طبقه‌بند بیز و نزدیکترین همسایگی
۱۵۳	۱.۳.۵. مرحله اول: شناسایی ویژگی‌های مغشوش
۱۵۷	۲.۳.۵. مرحله دوم: تصحیح ویژگی‌های مغشوش
۱۵۸	۴.۵. خلاصه
<b>۱۵۹</b>	<b>فصل ۶</b>
<b>۱۵۹</b>	<b>تشخیص داده‌های پرت</b>
۱۶۲	۱.۶. تشخیص داده پرت مبتنی بر مجاورت بدون پارامتر
۱۶۸	۲.۶. تشخیص داده پرت با الگوریتم ROCF
۱۷۳	۳.۶. تشخیص داده پرت با الگوریتم LDBSCAN
۱۷۹	۴.۶. تشخیص داده پرت با ترکیب الگوریتم ROCF و LDBSCAN
۱۸۳	۵.۶. خلاصه
<b>۱۸۵</b>	<b>منابع و مآخذ</b>
<b>۱۸۹</b>	<b>واژه‌نامه‌ها</b>
۱۹۱	واژه‌نامه انگلیسی به فارسی
۱۹۹	واژه‌نامه فارسی به انگلیسی
<b>۲۰۰</b>	<b>نمایه</b>



# فصل ۱

## پیش‌پردازش داده‌ها



## مقدمه

رشد روزافزون داده‌ها موجب شده تا سازمان‌ها با داده‌های حجیم، منابع ناهمگون<sup>۱</sup> و توزیع‌شده<sup>۲</sup> روبرو شوند. استخراج دانش از داده‌ها برای ارائه خدمات و هدایت فرایند تصمیم‌گیری، به یک وظیفه اصلی در مدیریت داده‌ها تبدیل شده است [۱]. سازمان‌ها برای استخراج دانش و تصمیم‌گیری مبتنی بر داده‌های ذخیره‌شده یا گذرای<sup>۳</sup> خود، نیازمند روش‌ها و ابزارهای خودکار هستند. وظیفه اصلی این ابزارها تبدیل حجم عظیم داده‌های خام به دانش و اطلاعات مفید است. این امر منجر به شکوفایی و رشد سریع دانش داده‌کاوی و کاربردهای مختلف آن در علوم رایانه در سال‌های اخیر شده است.

جمع‌آوری این حجم از داده‌ها و تبدیل آن‌ها به دانش، موجب ایجاد موضوعی به نام «کیفیت داده‌ها»<sup>۴</sup> برای سازمان‌ها می‌شود. زمانی که داده‌ها از منابع و سامانه‌های مختلف به یک سامانه دیگر منتقل می‌شوند، ممکن است دچار مشکلاتی از قبیل قالب و دامنه ناهمگون، ناسازگاری<sup>۵</sup> و مقادیر از دست‌رفته<sup>۶</sup> شوند. تصمیم‌گیری بر مبنای این داده‌های فاقد کیفیت، علاوه بر این که به ساختار سازمان آسیب می‌زند، باعث می‌شود هزینه‌های زیادی به سازمان وارد شود. بر اساس مطالعات انجام‌شده بیش از ۳۰٪ داده‌ها فاقد کیفیت هستند. این‌گونه داده‌ها موجب شده سالانه ۳ تریلیون دلار به دولت آمریکا خسارت وارد شود [۳]؛ لذا تامین کیفیت داده‌ها در منابع داده از اهمیت بالایی برخوردار است.

کیفیت داده‌ها شامل شاخص‌های صحت<sup>۷</sup>، کامل بودن<sup>۸</sup>، سازگاری<sup>۹</sup>، به هنگام بودن<sup>۱۰</sup>، قابل اعتماد بودن<sup>۱۱</sup> و تفسیرپذیری<sup>۱۲</sup> است [۱۰]. داده‌ها در دنیای واقعی به دلایلی از جمله معیوب بودن تجهیزات جمع‌آوری اطلاعات، خطای سامانه و یا انسانی در هنگام وارد کردن داده و قالب‌های ناهمگون در هنگام جمع‌آوری اطلاعات از منابع توزیع‌شده دچار مشکل عدم کیفیت می‌شوند. از

---

<sup>1</sup> Heterogeneous

<sup>2</sup> Distributed

<sup>3</sup> Transient

<sup>4</sup> Data quality

<sup>5</sup> Inconsistency

<sup>6</sup> Missing value

<sup>7</sup> Accuracy

<sup>8</sup> Completeness

<sup>9</sup> Consistency

<sup>10</sup> Timeliness

<sup>11</sup> Believability

<sup>12</sup> Interpretability

جمله متداول‌ترین مشکلات کیفیت داده‌ها می‌توان به مقادیر ازدست‌رفته، ناسازگاری، مقادیر ناصحیح<sup>۱</sup> و تکراری<sup>۲</sup> اشاره کرد.

پیش‌پردازش داده‌ها یکی از مهم‌ترین و ابتدایی‌ترین مراحل کشف دانش است که نقش اساسی در حصول اطمینان از کیفیت داده‌ها دارد. پاک‌سازی<sup>۳</sup> داده‌ها و یکپارچه‌سازی<sup>۴</sup> آن‌ها از جمله فرایندهایی هستند که در پیش‌پردازش انجام می‌گیرد. فراهم کردن داده‌هایی با کیفیت بالا کاری زمان‌بر و پرهزینه است؛ تا جایی که برآورد شده است که ۳۰٪-۸۰٪ هزینه و زمان ساخت پایگاه داده تحلیلی<sup>۵</sup> مربوط به پاک‌سازی داده است؛ اما صرف این هزینه و زمان زیاد در قبال هزینه‌های شکست بر اثر داده‌های فاقد کیفیت بسیار ارزشمند است [۳]. تضمین کیفیت داده‌ها در ساخت پایگاه داده تحلیلی بسیار حیاتی است.

بر این اساس در این فصل، در بخش ۱.۱ به معرفی کیفیت داده‌ها و ابعاد آن و در بخش ۲.۱ به بررسی مشکلات کیفیت داده و راهکارهای پاکسازی آن‌ها می‌پردازیم.

## ۱.۱. کیفیت داده‌ها و ابعاد آن

تضمین کیفیت داده‌ها از اهمیت بالایی برای سازمان‌ها برخوردار است؛ بالاخص زمانی که تصمیمات آینده سازمان، مبتنی بر این داده‌ها انجام پذیرد. منابع داده‌ای سازمان‌ها ممکن است به صورت ساخت‌یافته<sup>۶</sup>، نیمه‌ساخت‌یافته<sup>۷</sup> و غیرساخت‌یافته<sup>۸</sup> باشند. منابع ساخت‌یافته دارای موجودیتی<sup>۹</sup> با ساختار ثابت و مشخص هستند که مجموعه داده‌های رابطه‌ای<sup>۱۰</sup> و پایگاه داده‌های تحلیلی از رایج‌ترین انواع آن هستند. در منابع نیمه‌ساخت‌یافته موجودیت‌های مربوط به یک کلاس دارای ویژگی‌های متفاوتی هستند که رایج‌ترین نوع این منابع، داده‌های ارائه‌شده در قالب XML است. داده‌های منابع غیرساخت‌یافته دارای هیچ ساختار یا سازماندهی از پیش تعیین شده‌ای نیستند. فایل‌های متنی<sup>۱۱</sup> نمونه متداول از این منابع هستند [۱۳].

یک منبع داده در صورتی دارای کیفیت است که طرح<sup>۱۲</sup> و داده‌های آن دارای کیفیت باشند. کیفیت یک منبع، بر اساس میزان نزدیکی آن به نمونه‌هایش در دنیای واقعی سنجیده می‌شود. برای اندازه‌گیری این نزدیکی، از شاخص‌های کیفیت استفاده می‌شود. چنانچه طراحی اولیه طرح

<sup>1</sup> Incorrect value

<sup>2</sup> Duplicate

<sup>3</sup> Cleaning

<sup>4</sup> Integration

<sup>5</sup> Data warehouse

<sup>6</sup> Structured

<sup>7</sup> Semi-structured

<sup>8</sup> Unstructured

<sup>9</sup> Entity

<sup>10</sup> Relational

<sup>11</sup> Text

<sup>12</sup> Schema