

داده کاوی

(مفاهیم و تکنیک‌ها)

نویسنده:

ژیاوی هان - میشلین کمبر - ژان پی

مترجم:

دکتر مهدی اسماعیلی

نیاز دانش

پیش‌گفتار

مترجم

اگر بر جای من غیری کزیند دوست حاکم اوست

حرامم باد اگر من جان بجای دوست بگزینم

داده‌ها یکی از مولفه‌های ارزشمند دنیای امروزی تلقی می‌شود و بدون شک جمع‌آوری و تحلیل داده‌ها برای سازمان‌ها و شرکت‌ها و حتی زندگی روزمره‌ی افراد سودمند است. پیشرفت‌هایی که در رسانه‌ها و ابزارهای ذخیره‌سازی داده‌ها صورت گرفته است، دغدغه‌های ما را در زمینه‌ی جمع‌آوری و نگهداری داده‌ها تا حدودی حل نموده است، اما مشکل چگونگی تحلیل این داده‌ها است. امروزه شما با حجم انبوهی از داده‌ها روبرو هستید و تحلیل حجم بالایی از داده‌ها نیز یک ضرورت است. یکی از علومی که امروزه در دنیا از اهمیت ویژه‌ای برخوردار است و بی‌شک می‌توان پیشرفت روزافزون آن را در بسیاری از کاربردهای علمی و صنعتی مشاهده نمود، علم داده‌کاوی است. داده‌کاوی همچون هر کاوش و تحلیل دیگری به دنبال گنجی از اطلاعات در میان اقیانوسی از داده‌ها است و رشد بی‌رویه‌ی داده‌ها اهمیت آن را دو چندان کرده است.

صاحب‌نظران، بویژه منبع‌شناسان در زمینه‌ی داده‌کاوی بهتر می‌دانند که کتاب پیش روی شما (منبع اصلی) یکی از ارزنده‌ترین کتاب‌هایی است که در این حوزه به رشته‌ی تحریر درآمده است. از اینرو از میان کتاب‌های فراوانی که در این حوزه وجود دارد، به انتخاب و ترجمه‌ی کتاب مزبور همت گماشتیم. هدف این کتاب، بیان مفاهیم و تکنیک‌های داده‌کاوی است. موضوعات به گونه‌ای در کتاب تنظیم شده‌اند که هم برای مبتدیان و هم برای حرفه‌ای‌ها مناسب هستند، چون ابتدا مفاهیم و پایه‌های نظری بررسی می‌شوند و در ادامه شما می‌توانید روش‌های پیشرفته‌ی داده‌کاوی را مطالعه کنید.

مطالب کتاب در سیزده فصل گردآوری و تهیه شده است. مقدمه‌ای بر داده‌کاوی و همچنین بیانی اجمالی از مطالب کتاب را می‌توانید در فصل اول مطالعه بفرمایید. گونه‌های متفاوت داده‌ها و خصوصیات هر یک از آنها در فصل دوم بررسی می‌شوند و در فصل سوم به سراغ روش‌های پیش‌پردازش داده‌ها می‌رویم. فصل‌های چهارم و پنجم به موضوع انبارش داده‌ها تخصیص داده شده است. کاوش الگوهای مکرر، مشارکت‌ها و همبستگی‌ها موضوع فصل‌های ششم و هفتم را تشکیل می‌دهند. مفاهیم پایه‌ای دسته‌بندی در فصل هشتم بیان می‌شوند و به دنبال آن روش‌های پیشرفته‌ی آن در فصل نهم شرح داده شده‌اند. تحلیل خوشه‌بندی و روش‌های پیشرفته‌ی آن موضوع فصل‌های دهم و یازدهم است. فصلی مجزا برای انواع داده‌های پرت و روش‌های شناسایی آن در نظر گرفته شده است. در فصل دوازدهم بیشتر درباره‌ی داده‌های پرت صحبت می‌کنیم. موضوع آخرین فصل کتاب هم روندها و مرزهای تحقیق در حوزه‌ی داده‌کاوی است.

در ترجمه‌ی این کتاب تلاش شده است ضمن اصالت جمله‌بندی‌های کتاب، مطالب به گونه‌ای ترجمه شوند تا خوانندگان محترم بتوانند مفاهیم آن را به راحتی درک کنند. واژه‌های معادل در این متن پیشنهادی هستند و چه بسا برای برخی از اصطلاحات برابری دیگر (و احياناً بهتر) بتوان یافت. آنچه مسلم است این است که اگر

چه گاهی تبدیل واژه‌ها آنچنانکه باید و شاید انجام نشده است اما در ادای جملات و بیان موضوع امانت مراعات و تلاش فراوان شده است تا خوانندگان گرامی با متنی مبهم و گیج‌کننده روبرو نشوند. بخش پایانی هر فصل در کتاب اصلی به معرفی منابع می‌پردازد که در ترجمه‌ی کتاب وجود ندارد. نه اینکه این بخش اهمیت نداشته باشد که به زعم مترجم یکی از بخش‌های ارزشمند کتاب به خصوص برای دوستان محقق و پژوهشگر محسوب می‌شود. اما از آنجا که کتاب دارای صفحات بالایی است و ما علاقه‌مند بودیم کتاب را در یک جلد و آنهم با مبلغی مناسب در اختیار خوانندگان محترم قرار دهیم، تصمیم گرفته شد تا این بخش در ترجمه قرار نداشته باشد. اما شما به سهولت می‌توانید به مطالب این بخش‌ها دسترسی داشته باشید.

در اینجا لازم می‌دانم از همه‌ی اساتید و دانشجویان به خاطر راهنمایی‌های ارزشمندشان در حین ترجمه‌ی این کتاب سپاسگزاری کنم. همچنین از مدیریت محترم انتشارات نیاز دانش نیز به خاطر آماده‌سازی، چاپ و پخش این کتاب تشکر می‌کنم. رهین محبت بی‌دریغ خانواده‌ام هستم که با فراهم‌سازی محیطی مناسب مرا یاری نمودند. با وجود همه‌ی سعی و تلاشی که در تمام مراحل آماده‌سازی این کتاب انجام گرفته است، یقین دارم که عاری از اشتباه نیست، چرا که تنها مکتوب بی‌نقص همان معجزه‌ی جاوید قرآن کریم است. در آخر ضمن سپاسگزاری از همه کسانی که مرا یاری داده‌اند و با پذیرش مسئولیت هرگونه کاستی احتمالی، امیدوارم که این اندک مفید افتد.

مهدی اسماعیلی

۱۳۹۳

فهرست مطالب

فصل ۱ مقدمه ۱۱

- ۱-۱ چرا داده کاوی؟ ۱۱
۱-۱-۱ حرکت به سوی عصراطلاعات ۱۱
۲-۱-۱ داده کاوی به عنوان سیر تکاملی فن آوری اطلاعات ۱۲
۲-۱ داده کاوی چیست؟ ۱۵
۳-۱ چه نوعی از داده‌ها کاوش می‌شوند؟ ۱۷
۱-۳-۱ داده‌های پایگاه داده‌ها ۱۸
۲-۳-۱ انبارهای داده ۱۹
۳-۳-۱ داده‌های تراکنشی ۲۱
۴-۳-۱ گونه‌های دیگر داده‌ها ۲۲
۴-۱ چه نوعی از الگوها کاوش می‌شوند؟ ۲۳
۱-۴-۱ توصیف و تفکیک ۲۴
۲-۴-۱ کاوش الگوهای مکرر، مشارکت‌ها و وابستگی‌ها ۲۵
۳-۴-۱ دسته‌بندی و رگرسیون برای تحلیل پیشگویی ۲۶
۴-۴-۱ تحلیل خوشه ۲۸
۵-۴-۱ تحلیل داده‌های پرت ۲۹
۶-۴-۱ آیا کلیه الگوها جالب هستند؟ ۲۹
۵-۱ کدام فن آوری‌ها استفاده می‌شوند؟ ۳۱
- ۱-۵-۱ آمار ۳۱
۲-۵-۱ یادگیری ماشین ۳۲
۳-۵-۱ سیستم‌های پایگاه داده‌ها و انبارهای داده ۳۳
۴-۵-۱ بازیابی اطلاعات ۳۴
۶-۱ هدف داده کاوی کدام نوع از برنامه‌های کاربردی است؟ ۳۴
۱-۶-۱ هوش تجاری ۳۵
۲-۶-۱ موتورهای جستجوی وب ۳۵
۷-۱ موضوعات عمده در داده کاوی ۳۶
۱-۷-۱ متدولوژی کاوش ۳۶
۲-۷-۱ تعامل کاربر(دخالت کاربر) ۳۸
۳-۷-۱ کارآمد بودن و قابلیت مقیاس پذیری ۳۸
۴-۷-۱ تنوع گونه‌های پایگاه داده‌ها ۳۹
۵-۷-۱ داده کاوی و جامعه ۳۹
۸-۱ خلاصه ۴۱
۹-۱ تمرین‌ها ۴۳

فصل ۲ شناخت داده‌ها ۴۵

- ۱-۲ نمونه‌ها و انواع صفات خاصه ۴۶
۱-۱-۲ یک صفت خاصه چیست؟ ۴۶
۲-۱-۲ صفات خاصه‌ی اسمی ۴۷
۳-۱-۲ صفات خاصه‌ی دودویی ۴۷
۴-۱-۲ صفات خاصه‌ی ترتیبی ۴۸
۵-۱-۲ صفات خاصه‌ی عددی ۴۸
۶-۱-۲ صفات خاصه‌ی گسسته در برابر پیوسته ۴۹
۲-۲ آمار توصیفی ۵۰
۱-۲-۲ محاسبه‌ی شاخص‌های مرکزی: میانگین، میانه و مُد ۵۰
۲-۲-۲ محاسبه‌ی شاخص‌های پراکندگی داده‌ها ۵۳
- ۳-۲-۲ نمایش گرافیکی آمار توصیفی ۵۷
۳-۲ مصورسازی داده‌ها ۶۱
۱-۳-۲ روش‌های مصورسازی پیکسل‌گرا ۶۲
۲-۳-۲ تکنیک‌های تصویر کردن هندسی ۶۳
۳-۳-۲ تکنیک‌های مصورسازی مبتنی بر شمایل ۶۶
۴-۳-۲ تکنیک‌های سلسله‌مراتبی مصورسازی ۶۷
۵-۳-۲ مصورسازی داده‌ها و روابط پیچیده ۶۸
۴-۲ محاسبه‌ی تشابه و عدم تشابه ۷۱
۱-۴-۲ ماتریس داده‌ها در مقابل ماتریس عدم تشابه ۷۱
۲-۴-۲ محاسبه‌ی همسایگی برای صفات خاصه‌ی اسمی ۷۳
۳-۴-۲ محاسبه‌ی همسایگی برای صفات خاصه‌ی دودویی ۷۴

۷-۴-۲	تشابه کسینوسی	۷۶-۴-۲	عدم تشابه داده‌های عددی: فاصله‌ی مینکوفسکی
۸۳-۴-۲	خلاصه	۷۸-۴-۲	محاسبه‌ی همسایگی برای صفات خاصه‌ی ترتیبی
۸۴-۴-۲	تمرین‌ها	۶-۴-۲	محاسبه‌ی عدم تشابه برای مخلوطی
		۷۹-۴-۲	از انواع صفات خاصه

فصل ۳

پیش‌پردازش داده‌ها

۸۷-۴-۳	رگرسیون و مدل‌های لگاریتمی-خطی:	۸۸-۴-۳	مروری بر پیش‌پردازش داده‌ها
۱۰۸-۴-۳	کاهش پارامتری داده‌ها	۸۸-۴-۳	کیفیت داده‌ها: چرا پیش‌پردازش داده‌ها؟
۱۰۸-۴-۳	هیستوگرام‌ها	۸۹-۴-۳	وظایف اصلی در پیش‌پردازش داده‌ها
۱۱۰-۴-۳	خوشه‌بندی	۹۱-۴-۳	پالایش داده‌ها
۱۱۰-۴-۳	نمونه‌گیری	۹۱-۴-۳	مقادیر ناموجود
۱۱۲-۴-۳	تجمیع در مکعب داده‌ها	۹۳-۴-۳	داده‌های نویزدار
۱۱۳-۴-۳	تبدیل داده‌ها و گسسته‌سازی داده‌ها	۹۴-۴-۳	پالایش داده‌ها به عنوان یک فرایند
۱۱۳-۴-۳	مروری اجمالی بر روی راهبردهای تبدیل داده‌ها	۹۷-۴-۳	یکپارچه‌سازی داده‌ها
۱۱۵-۴-۳	تبدیل داده‌ها با کمک نرمال‌سازی	۹۷-۴-۳	مشکل شناسایی موجودیت
۱۱۷-۴-۳	گسسته‌سازی با کمک بسته‌بندی	۹۸-۴-۳	افزونگی و تحلیل همبستگی
۱۱۷-۴-۳	گسسته‌سازی با تحلیل هیستوگرام	۱۰۱-۴-۳	تکرار تاپل
۱۱۷-۴-۳	گسسته‌سازی با کمک خوشه‌بندی، درخت	۱۰۲-۴-۳	تشخیص و حل تصادم میان مقادیر داده‌ها
۱۱۸-۴-۳	تصمیم و تحلیل همبستگی	۱۰۲-۴-۳	کاهش داده‌ها
۱۱۹-۴-۳	ایجاد سلسله‌مراتب مفهومی برای داده‌های اسمی	۱۰۲-۴-۳	مروری بر راهبردهای کاهش داده‌ها
۱۲۲-۴-۳	خلاصه	۱۰۳-۴-۳	تبدیل‌های موجک
۱۲۳-۴-۳	تمرین‌ها	۱۰۵-۴-۳	تحلیل مولفه‌های اصلی
		۱۰۶-۴-۳	انتخاب زیرمجموعه‌ای از صفات خاصه

فصل ۴

انبارش داده‌ها و پردازش تحلیلی برخط

۱۲۷-۴-۴	طراحی و کاربرد انبار داده‌ها	۱۲۷-۴-۴	انبار داده‌ها: مفاهیم پایه‌ای
۱۴۹-۴-۴	یک چارچوب تحلیل کسب‌وکار برای طراحی انبار داده‌ها	۱۲۸-۴-۴	انبار داده‌ها چیست؟
۱۴۹-۴-۴	فرایند طراحی انبار داده‌ها	۱۲۸-۴-۴	تفاوت‌های میان سیستم‌های پایگاه داده عملیاتی و انبارهای داده‌ها
۱۵۰-۴-۴	کاربرد انبار داده‌ها در پردازش اطلاعات	۱۳۰-۴-۴	چرا یک انبار داده‌ی مجزا؟
۱۵۱-۴-۴	از پردازش تحلیلی برخط تا داده‌کاوی چندبُعدی	۱۳۱-۴-۴	انبارش داده‌ها: یک معماری چند ردیفی
۱۵۳-۴-۴	پیاده‌سازی انبار داده‌ها	۱۳۲-۴-۴	مدل‌های انبار داده‌ها
۱۵۴-۴-۴	مروری بر محاسبه‌ی مؤثر مکعب داده‌ها	۱۳۲-۴-۴	مدل‌های انبار داده‌ها: استخراج، تبدیل و بارگذاری
۱۵۵-۴-۴	شاخص‌بندی داده‌های OLAP: شاخص نگاشت‌بیتی و شاخص پیوند	۱۳۵-۴-۴	مخزن مِتاداده‌ها
۱۶۱-۴-۴	پردازش مؤثر پرسش‌های OLAP	۱۳۵-۴-۴	مدل‌سازی انبار داده‌ها: مکعب داده‌ها و OLAP
۱۶۲-۴-۴	معماری‌های سرویس‌دهنده‌ی OLAP	۱۳۶-۴-۴	مکعب داده‌ها: یک مدل داده‌ای چندبُعدی
۱۶۳-۴-۴	تعمیم داده‌ها با کمک استنتاج صفت‌گرا	۱۳۶-۴-۴	ستاره‌ای، برفکونه و صورفلکی:
۱۶۵-۴-۴	استنتاج صفت‌گرا برای توصیف داده‌ها	۱۳۹-۴-۴	شماهایی برای مدل‌های داده‌ای چندبُعدی
۱۷۰-۴-۴	پیاده‌سازی مؤثر فرایند استنتاج صفت‌گرا	۱۴۲-۴-۴	ابعاد: نقش سلسله‌مراتب‌های مفهومی
۱۷۲-۴-۴	استنتاج صفت‌گرا برای مقایسه‌های کلاسی	۱۴۴-۴-۴	سنججه‌ها: طبقه‌بندی و محاسبه‌ی آنها
۱۷۶-۴-۴	خلاصه	۱۴۵-۴-۴	عملیات معمول OLAP
۱۷۸-۴-۴	تمرین‌ها	۱۴۵-۴-۴	مدل شبکه ستاره‌ای برای پرس‌وجو از پایگاه داده‌های چندبُعدی

فصل ۵ تکنولوژی مکعب داده‌ها ۱۸۳

- ۱-۵ محاسبه‌ی مکعب داده‌ها: مفاهیم مقدماتی ۱۸۴
- ۱-۱-۵ ساخت مکعب: مکعب کامل، مکعب کوه یخی، مکعب بسته و پوسته‌ی مکعب ۱۸۴
- ۲-۱-۵ راهبردهای کلی برای محاسبه‌ی مکعب داده‌ها ۱۸۸
- ۲-۵ روش‌های محاسبه‌ی مکعب داده‌ها ۱۹۰
- ۱-۲-۵ جمع‌چندراهه آرایه برای محاسبه‌ی مکعب کامل ۱۹۰
- ۲-۲-۵ BUC: محاسبه‌ی مکعب‌های کوه‌یخی از شبه‌مکعب نوک به سمت پایین ۱۹۵
- ۳-۲-۵ Star-Cubing: محاسبه‌ی مکعب‌های کوه‌یخی با کمک ساختار ستاره-درختی ۱۹۹
- ۴-۲-۵ پیش‌محاسبه‌ی قطعات پوسته برای تسریع OLAP در کار با ابعاد بالا ۲۰۴
- ۳-۵ پردازش گونه‌های پیشرفته‌ای از پرسش‌ها ۲۱۳
- ۱-۳-۵ مکعب‌های نمونه‌گیری: کاوش مبتنی بر OLAP بر روی داده‌های نمونه ۲۱۳
- ۲-۳-۵ مکعب‌های رتبه‌بندی: محاسبه‌ی موثر پرسش‌های رتبه‌ای ۲۱۹
- ۴-۵ تحلیل داده‌های چندبُعدی در فضای مکعب ۲۲۱
- ۱-۴-۵ مکعب‌های پیش‌گویانه: کاوش پیش‌گویانه در فضای مکعب ۲۲۱
- ۲-۴-۵ مکعب‌های چندجانبه: جمع‌بندی پیچیده در چندین سطح از داده‌بندی ۲۲۳
- ۳-۴-۵ کاوش اکتشافی و مبتنی بر استثناء فضای مکعب ۲۲۵
- ۵-۵ خلاصه ۲۲۹
- ۶-۵ تمرین‌ها ۲۳۰

فصل ۶ کاوش الگوهای مکرر، مشارکت‌ها و همبستگی‌ها: مفاهیم پایه و روش‌ها ۲۳۵

- ۱-۶ مفاهیم پایه ۲۳۶
- ۱-۱-۶ تحلیل سبد خرید: یک مثال انگیزشی ۲۳۶
- ۲-۱-۶ مجموعه‌اقدام مکرر، مجموعه اقدام بسته و قوانین انجمنی ۲۳۷
- ۲-۶ روش‌های کاوش مجموعه‌اقدام مکرر ۲۴۰
- ۱-۲-۶ الگوریتم Apriori: یافتن مجموعه‌اقدام مکرر با کمک تولید کاندیداهای محدود ۲۴۰
- ۲-۲-۶ ایجاد قوانین انجمنی از مجموعه‌اقدام مکرر ۲۴۵
- ۳-۲-۶ بهبود Apriori ۲۴۵
- ۴-۲-۶ رویکرد رشد الگو برای کاوش مجموعه‌اقدام مکرر ۲۴۷
- ۵-۲-۶ کاوش مجموعه‌اقدام مکرر با کمک قالب عمودی داده‌ها ۲۵۱
- ۶-۲-۶ کاوش الگوهای بسته و ماکسیمال ۲۵۳
- ۳-۶ کدامیک از الگوها جالب هستند؟ روش‌های ارزشیابی الگو ۲۵۵
- ۱-۳-۶ قوانین قوی الزاماً جالب نیستند ۲۵۵
- ۲-۳-۶ از تحلیل مشارکت تا تحلیل همبستگی ۲۵۶
- ۳-۳-۶ مقایسه‌ای میان سنج‌های ارزشیابی الگو ۲۵۸
- ۴-۶ خلاصه ۲۶۲
- ۵-۶ تمرین‌ها ۲۶۳

فصل ۷ مباحث پیشرفته در کاوش الگو ۲۶۷

- ۱-۷ کاوش الگو: یک نقشه‌ی راه ۲۶۷
- ۲-۷ کاوش الگو در فضای چندسطحی و چندبُعدی ۲۷۰
- ۱-۲-۷ کاوش مشارکت‌های چندسطحی ۲۷۱
- ۲-۲-۷ کاوش مشارکت‌های چندبُعدی ۲۷۴
- ۳-۲-۷ کاوش قوانین انجمنی کمی ۲۷۶
- ۴-۲-۷ کاوش الگوهای نادر و الگوهای منفی ۲۷۸
- ۳-۷ کاوش الگوی مکرر مبتنی بر محدودیت ۲۸۱
- ۱-۳-۷ کاوش قوانین انجمنی با کمک متقاعدده ۲۸۲
- ۲-۳-۷ تولید الگو بر اساس محدودیت: هرس نمودن فضای الگو و فضای داده‌ها ۲۸۳
- ۴-۷ کاوش داده‌هایی با ابعاد بالا و الگوهای بسیار طولانی ۲۸۸
- ۱-۴-۷ کاوش الگوهای بسیار طولانی با کمک ائتلاف الگوها ۲۸۹
- ۵-۷ کاوش الگوهای فشرده یا تقریبی ۲۹۳
- ۱-۵-۷ کاوش الگوهای فشرده با کمک خوشه‌بندی الگو ۲۹۴
- ۲-۵-۷ استخراج k الگوی برتر با آگاهی از افزونگی ۲۹۷
- ۶-۷ کشف و کاربرد الگو ۲۹۹
- ۱-۶-۷ حاشیه‌نویسی معنایی الگوهای مکرر ۲۹۹
- ۲-۶-۷ کاربردهایی از کاوش الگو ۳۰۴
- ۷-۷ خلاصه ۳۰۷
- ۸-۷ تمرین‌ها ۳۰۸

فصل ۸ دسته‌بندی: مفاهیم پایه

- ۳۱۱-----
- ۱-۵-۸ متریک‌هایی جهت ارزشیابی کارایی دسته‌بند- ۳۴۷
- ۲-۵-۸ روش holdout و زیرنمونه‌گیری تصادفی- ۳۵۳
- ۳-۵-۸ اعتبارسنجی متقابل- ۳۵۳
- ۴-۵-۸ روش bootstrap- ۳۵۴
- ۵-۵-۸ انتخاب مدل با کمک آزمون‌های آماری- ۳۵۴
- ۶-۵-۸ مقایسه‌ی مدل‌ها بر اساس هزینه-سود و منحنی‌های ROC- ۳۵۶
- ۶-۸ تکنیک‌هایی جهت بهبود صحت دسته‌بندی- ۳۶۰
- ۱-۶-۸ معرفی روش‌های تلفیقی- ۳۶۰
- ۲-۶-۸ روش bagging- ۳۶۱
- ۳-۶-۸ روش boosting و الگوریتم AdaBoost- ۳۶۲
- ۴-۶-۸ جنگل‌های تصادفی- ۳۶۵
- ۵-۶-۸ بهبود صحت دسته‌بندی در داده‌هایی با عدم‌تعادل کلاس- ۳۶۶
- ۷-۸ خلاصه- ۳۶۹
- ۸-۸ تمرین‌ها- ۳۷۰
- ۱-۸ مفاهیم پایه- ۳۱۱
- ۱-۱-۸ دسته‌بندی چیست؟- ۳۱۲
- ۲-۱-۸ رویکرد عمومی برای دسته‌بندی- ۳۱۲
- ۲-۸ استقراء درخت تصمیم- ۳۱۴
- ۱-۲-۸ استقراء درخت تصمیم- ۳۱۵
- ۲-۲-۸ سنجش‌های انتخاب صفت‌خاصه- ۳۱۹
- ۳-۲-۸ هرس نمودن درخت- ۳۲۷
- ۴-۲-۸ مقیاس‌پذیری و استقراء درخت تصمیم- ۳۲۹
- ۵-۲-۸ کاوش بصری جهت استقراء درخت تصمیم- ۳۳۱
- ۳-۸ روش‌های دسته‌بندی بیزی- ۳۳۳
- ۱-۳-۸ نظریه‌ی بیز- ۳۳۳
- ۲-۳-۸ دسته‌بندی بیزی ساده- ۳۳۴
- ۴-۸ دسته‌بندی مبتنی بر قاعده- ۳۳۸
- ۱-۴-۸ دسته‌بندی با کمک قواعدی به شکل IF-THEN- ۳۳۸
- ۲-۴-۸ استخراج قواعد از یک درخت تصمیم- ۳۴۰
- ۳-۴-۸ استقراء قاعده با کمک الگوریتم پوشای متوالی- ۳۴۲
- ۵-۸ ارزشیابی و انتخاب مدل- ۳۴۷

فصل ۹ دسته‌بندی: روش‌های پیشرفته

- ۳۷۳-----
- ۱-۵-۹ دسته‌بندی‌های k همسایه‌ی نزدیک- ۴۰۲
- ۲-۵-۹ استدلال مبتنی بر مورد- ۴۰۴
- ۶-۹ روش‌های دسته‌بندی دیگر- ۴۰۴
- ۱-۶-۹ الگوریتم‌های ژنتیک- ۴۰۵
- ۲-۶-۹ رویکرد مجموعه‌های راف- ۴۰۵
- ۳-۶-۹ مجموعه‌های فازی- ۴۰۶
- ۷-۹ چند موضوع دیگر درباره‌ی دسته‌بندی- ۴۰۸
- ۱-۷-۹ دسته‌بندی چندکلاسه- ۴۰۹
- ۲-۷-۹ دسته‌بندی نیمه‌ناظر- ۴۱۰
- ۳-۷-۹ یادگیری فعال- ۴۱۲
- ۴-۷-۹ یادگیری انتقالی- ۴۱۳
- ۸-۹ خلاصه- ۴۱۶
- ۹-۹ تمرین‌ها- ۴۱۸
- ۱-۹ شبکه‌های بیز- ۳۷۳
- ۱-۱-۹ مفاهیم و مکانیزم‌ها- ۳۷۴
- ۲-۱-۹ آموزش شبکه‌های بیز- ۳۷۵
- ۲-۹ دسته‌بندی با کمک الگوریتم پس‌انتشار- ۳۷۷
- ۱-۲-۹ شبکه‌ی عصبی پیش‌خور چندلایه- ۳۷۸
- ۲-۲-۹ معرفی یک توپولوژی شبکه- ۳۷۹
- ۳-۲-۹ پس‌انتشار- ۳۸۰
- ۴-۲-۹ درون جعبه سیاه: پس‌انتشار و قابلیت تفسیرپذیری- ۳۸۵
- ۳-۹ ماشین‌های بردار پشتیبان- ۳۸۷
- ۱-۳-۹ تفکیک‌پذیری خطی داده‌ها- ۳۸۷
- ۲-۳-۹ تفکیک‌پذیری غیرخطی- ۳۹۱
- ۴-۹ دسته‌بندی با کمک الگوهای مکرر- ۳۹۴
- ۱-۴-۹ دسته‌بندی مشارکتی- ۳۹۵
- ۲-۴-۹ دسته‌بندی مبتنی بر الگوهای مکرر جداساز- ۳۹۸
- ۵-۹ یادگیرنده‌های کُند (یادگیری از طریق همسایه‌ها)- ۴۰۱

فصل ۱۰ تحلیل خوشه: مفاهیم پایه و روش‌ها

- ۴۲۱-----
- ۱-۲-۱۰ روش k-means: تکنیکی مبتنی بر گرانیگاه- ۴۲۹
- ۲-۲-۱۰ روش k-medoids: تکنیکی مبتنی بر شی نماینده- ۴۳۲
- ۳-۱۰ روش‌های سلسله‌مراتبی- ۴۳۵
- ۱-۳-۱۰ خوشه‌بندی سلسله‌مراتبی تجمیعی در مقابل تقسیمی- ۴۳۷
- ۱-۱۰ تحلیل خوشه- ۴۲۲
- ۱-۱-۱۰ تحلیل خوشه چیست؟- ۴۲۲
- ۲-۱-۱۰ نیازهای تحلیل خوشه- ۴۲۳
- ۳-۱-۱۰ مروری بر روش‌های پایه‌ی خوشه‌بندی- ۴۲۶
- ۲-۱-۱۰ روش‌های افراز- ۴۲۸

- ۲-۳-۱۰ سنجه‌های فاصله در روش‌های الگوریتمیک-----۴۳۹
- ۳-۳-۱۰ الگوریتم BIRCH: خوشه‌بندی
سلسله‌مراتبی چندمرحله‌ای با کمک
درختان ویژگی خوشه‌بندی-----۴۴۱
- ۴-۳-۱۰ الگوریتم Chameleon: خوشه‌بندی
سلسله‌مراتبی چندمرحله‌ای با کمک
مدل‌سازی پویا-----۴۴۳
- ۵-۳-۱۰ خوشه‌بندی سلسله‌مراتبی احتمالاتی-----۴۴۵
- ۴-۱۰ روش‌های مبتنی بر چگالی-----۴۴۸
- ۱-۴-۱۰ DBSCAN: خوشه‌بندی مبتنی بر چگالی
بر اساس اتصال مناطقی با چگالی بالا-----۴۴۸
- ۲-۴-۱۰ OPTICS: نظم‌دهی نقاط برای شناسایی
ساختار خوشه‌بندی-----۴۵۱
- ۳-۴-۱۰ DENCLUE: خوشه‌بندی بر اساس
توابع توزیع چگالی-----۴۵۳
- ۵-۱۰ روش‌های مبتنی بر گرید-----۴۵۶
- ۱-۵-۱۰ تکنیک STING-----۴۵۶
- ۲-۵-۱۰ CLIQUE: یک تکنیک خوشه‌بندی
زیرفضا مشابه Apriori-----۴۵۸
- ۶-۱۰ ارزشیابی خوشه‌بندی-----۴۶۱
- ۱-۶-۱۰ ارزیابی روند خوشه‌بندی-----۴۶۱
- ۲-۶-۱۰ تعیین تعداد خوشه‌ها-----۴۶۳
- ۳-۶-۱۰ اندازه‌گیری کیفیت خوشه‌بندی-----۴۶۴
- ۷-۱۰ خلاصه-----۴۶۸
- ۸-۱۰ تمرین‌ها-----۴۶۹

فصل ۱۱ روش‌های پیشرفته‌ی تحلیل خوشه-----۴۷۳

- ۱-۳-۱۱ کاربردها و چالش‌ها-----۴۹۷
- ۲-۳-۱۱ سنجه‌های تشابه-----۵۰۰
- ۳-۳-۱۱ روش‌های خوشه‌بندی گراف-----۵۰۳
- ۴-۱۱ خوشه‌بندی همراه با محدودیت‌ها-----۵۰۶
- ۱-۴-۱۱ گروه‌بندی محدودیت‌ها-----۵۰۶
- ۲-۴-۱۱ روش‌هایی برای خوشه‌بندی همراه با
محدودیت‌ها-----۵۱۰
- ۵-۱۱ خلاصه-----۵۱۳
- ۶-۱۱ تمرین‌ها-----۵۱۵
- ۱-۱۱ خوشه‌بندی بر پایه‌ی مدل احتمالاتی-----۴۷۳
- ۱-۱-۱۱ خوشه‌های فازی-----۴۷۵
- ۲-۱-۱۱ خوشه‌های مبتنی بر مدل احتمالاتی-----۴۷۷
- ۳-۱-۱۱ الگوریتم EM-----۴۸۰
- ۲-۱۱ خوشه‌بندی داده‌هایی با ابعاد بالا-----۴۸۳
- ۱-۲-۱۱ خوشه‌بندی داده‌هایی با ابعاد بالا: مشکلات،
چالش‌ها و متدولوژی‌های اصلی-----۴۸۴
- ۲-۲-۱۱ روش‌های خوشه‌بندی زیرفضا-----۴۸۵
- ۳-۲-۱۱ خوشه‌بندی دوگانه-----۴۸۷
- ۴-۲-۱۱ روش‌های کاهش ابعاد و خوشه‌بندی طیفی-----۴۹۵
- ۳-۱۱ خوشه‌بندی گراف و شبکه-----۴۹۷

فصل ۱۲ شناسایی داده‌های پرت-----۵۱۷

- ۱-۱۲ داده‌های پرت و تحلیل آنها-----۵۱۸
- ۱-۱-۱۲ داده‌ی پرت چیست؟-----۵۱۸
- ۲-۱-۱۲ انواع داده‌های پرت-----۵۱۹
- ۳-۱-۱۲ چالش‌های تشخیص داده‌های پرت-----۵۲۲
- ۲-۱۲ روش‌های تشخیص داده‌های پرت-----۵۲۳
- ۱-۲-۱۲ روش‌های باناظر، نیمه‌ناظر و بی‌ناظر-----۵۲۳
- ۲-۲-۱۲ روش‌های آماری و روش‌های مبتنی بر
مجاورت و روش‌های مبتنی بر خوشه‌بندی-----۵۲۵
- ۳-۱۲ رویکردهای آماری-----۵۲۷
- ۱-۳-۱۲ روش‌های پارامتری-----۵۲۷
- ۲-۳-۱۲ روش‌های بدون پارامتر-----۵۳۲
- ۴-۱۲ رویکردهای مبتنی بر مجاورت-----۵۳۴
- ۱-۴-۱۲ تشخیص داده‌های پرت بر مبنای فاصله و
یک روش با حلقه‌ی تودرتو-----۵۳۴
- ۲-۴-۱۲ یک روش مبتنی بر گرید-----۵۳۶
- ۳-۴-۱۲ تشخیص داده‌های پرت بر مبنای چگالی-----۵۳۸
- ۴-۴-۱۲ رویکردهای مبتنی بر دسته‌بندی-----۵۴۴
- ۷-۱۲ کاوش داده‌های پرت بافتاری و گروهی-----۵۴۶
- ۱-۷-۱۲ تبدیل روش‌های تشخیص داده‌های پرت
بافتاری به روش‌های متداول برای این کار-----۵۴۷
- ۲-۷-۱۲ مدلسازی رفتار نرمال با توجه به بافتها-----۵۴۸
- ۳-۷-۱۲ کاوش داده‌های پرت گروهی-----۵۴۸
- ۸-۱۲ تشخیص داده‌های پرت در داده‌هایی با ابعاد بالا-----۵۵۰
- ۱-۸-۱۲ تعمیم روش‌های متداول تشخیص داده‌های پرت-----۵۵۱
- ۲-۸-۱۲ یافتن داده‌های پرت در زیرفضاها-----۵۵۳
- ۳-۸-۱۲ مدل‌سازی داده‌های پرت با ابعاد بالا-----۵۵۳
- ۹-۱۲ خلاصه-----۵۵۵
- ۱۰-۱۲ تمرین‌ها-----۵۵۷

فصل ۱۳ روندها و مرزهای تحقیق در حوزه‌ی داده‌کاوی

- ۵۵۹-----
- ۵۸۲----- ۲-۳-۱۳ داده‌کاوی برای صنایع خرده‌فروشی و مخابرات
- ۵۸۴----- ۳-۳-۱۳ داده‌کاوی در علوم و مهندسی
- ۵۸۶----- ۴-۳-۱۳ داده‌کاوی برای تشخیص نفوذ و پیشگیری از آن
- ۵۸۸----- ۵-۳-۱۳ داده‌کاوی و سیستم‌های پیشنهاد دهنده
- ۵۹۱----- ۴-۱۳ داده‌کاوی و جامعه
- ۱-۴-۱۳ داده‌کاوی غیرمحسوسی که در همه جا حضور دارد----- ۵۹۱
- ۲-۴-۱۳ محرمانگی، امنیت و اثرات اجتماعی داده‌کاوی----- ۵۹۳
- ۵-۱۳ روندهای داده‌کاوی----- ۵۹۶
- ۶-۱۳ خلاصه----- ۵۹۸
- ۷-۱۳ تمرین‌ها----- ۵۹۹
- ۱-۱۳ کاوش انواع داده‌های پیچیده----- ۵۵۹
- ۱-۱-۱۳ کاوش توالی‌ها: سری‌های زمانی، توالی‌های نمادین و توالی‌های زیستی----- ۵۶۰
- ۲-۱-۱۳ کاوش گراف‌ها و شبکه‌ها----- ۵۶۵
- ۳-۱-۱۳ کاوش گونه‌های دیگر داده‌ها----- ۵۶۸
- ۲-۱۳ متدولوژی‌های دیگر داده‌کاوی----- ۵۷۲
- ۱-۲-۱۳ داده‌کاوی آماری----- ۵۷۲
- ۲-۲-۱۳ دیدگاه‌هایی در مورد پایه‌های نظری داده‌کاوی----- ۵۷۴
- ۳-۲-۱۳ داده‌کاوی تصویری و صوتی----- ۵۷۵
- ۳-۱۳ کاربردهای داده‌کاوی----- ۵۸۱
- ۱-۳-۱۳ داده‌کاوی برای تحلیل داده‌های مالی----- ۵۸۱

فصل ۱

مقدمه

این کتاب مقدمه‌ای بر موضوع داده‌کاوی^۱ است که گاه با عنوان کشف دانش^۲ از داده‌ها نیز شناخته می‌شود. داده‌کاوی حوزه‌ی جدیدی تلقی می‌شود و به سرعت در حال رشد است. کتاب بر روی مفاهیم و تکنیک‌های اصلی داده‌کاوی متمرکز می‌شود که با کمک آنها می‌توان الگوهای جالبی را از داده‌ها و در برنامه‌های کاربردی متفاوت کشف نمود. تاکید ما بر روی تکنیک‌هایی است که جهت توسعه‌ی ابزارهای داده‌کاوی موثر، کارا و مقیاس‌پذیر از آنها استفاده می‌شود.

فصل حاضر بدین ترتیب سازمان‌دهی شده است. در بخش ۱-۱ خواهیم آموخت که چرا تقاضا برای داده‌کاوی افزایش یافته است و چرا بخشی از سیر تکامل تکنولوژی اطلاعات تلقی می‌شود. بخش ۱-۲ به معرفی داده‌کاوی در فرایند کشف دانش می‌پردازد. در ادامه جنبه‌های زیادی در مورد داده‌کاوی را بحث خواهیم کرد، نوع داده‌هایی که می‌توانند کاوش شوند (بخش ۱-۳)، گونه‌هایی از دانش که نتیجه‌ی کاوش است (بخش ۱-۴)، انواع تکنولوژی‌هایی که استفاده می‌شوند (بخش ۱-۵) و برنامه‌های کاربردی که در آنها از داده‌کاوی استفاده شده است (بخش ۱-۶). بدین ترتیب از منظرهای متفاوتی به داده‌کاوی خواهیم پرداخت. در بخش ۱-۷ که بخش پایانی فصل را تشکیل می‌دهد چالش‌های تحقیقاتی این حوزه بررسی می‌شوند.

۱-۱ چرا داده‌کاوی؟

ما در دنیایی زندگی می‌کنیم که روزانه حجم وسیعی از داده‌ها تولید می‌شود. یکی از نیازهای مهم تحلیل این داده‌های حجیم است. بخش ۱-۱-۱ نگاهی به این موضوع دارد که چگونه داده‌کاوی با مهیا ساختن ابزارهایی جهت کشف دانش از داده‌ها این نیاز را مرتفع می‌سازد. در بخش ۱-۱-۲ نیز مشاهده خواهیم کرد که چگونه داده‌کاوی را می‌توان نتیجه سیر تکاملی تکنولوژی اطلاعات دانست.

۱-۱-۱ حرکت به سوی عصر اطلاعات

جمله‌ی "ما در حال زندگی در عصر اطلاعات هستیم" شاید یک جمله‌ی رایج تلقی شود، اما بهر حال و در واقع ما در حال زندگی در عصر داده‌ها هستیم. قرار گرفتن حجم وسیع داده‌ها در مقیاس ترابایت و یا حتی پتابایت بر روی شبکه‌های کامپیوتری، شبکه جهانی وب و رسانه‌های ذخیره‌سازی متفاوت که داده‌ها را از علوم متفاوتی جمع‌آوری و نگهداری می‌کنند، این موضوع را به خوبی اثبات می‌کند. رشد انفجاری حجم داده‌ها نتیجه‌ی مکانیزه شدن جامعه‌ی ماست و همچنین توسعه‌ی ابزارهای قدرتمند برای جمع‌آوری و نگهداری داده‌ها است. کسب و کارهای

^۱ Data Mining

^۲ Knowledge Discovery