

بِسْمِ اللّٰهِ الرَّحْمٰنِ الرَّحِيْمِ

سرشناسه:	رضایی تبار، وحید، ۱۳۶۳ -
عنوان و نام پدید آور:	داده کاوی با زبان برنامه نویسی R=Data Mining with R/مؤلفان:وحید رضایی تبار، شقایق خودکامه ویراست.۲
وضعیت ویراست:	مشخصات نشر:
مشخصات ظاهری:	مشخصات ظاهری:
فروش:	تهران: دانشگاه علامه طباطبائی، ۱۴۰۳
شابک:	.۵۷۴ [۲] انتشارات دانشگاه علامه طباطبائی؛
وضعیت فهرست نویسی:	۵۰۵۰ ص: مصور(بخشی رنگی)، جدول، نمودار، وزیری.
موضوع:	۹۷۸-۹۶۴-۲۱۷-۶۰۹-۰ تومان
موضوع:	داده کاوی -- راهنمای آموزشی (علی)
موضوع:	Data Mining -- Study and Teaching (Higher)
موضوع:	آر (زبان برنامه نویسی کامپیوتر)
موضوع:	R (Computer Program Language)
شناسه افزوده:	-۱۳۷۴ خودکامه، شقایق،
شناسه افزوده:	دانشگاه علامه طباطبائی
ردیفندی کنگره:	ردیفندی کنگره: ۱۴۰۳ /۰۲ /۹۷/۶۷ QA
ردیفندی دیبوئی:	ردیفندی دیبوئی: ۰۰۶۳۱۲
شماره کتابخانه ملی:	شماره کتابخانه ملی: ۹۸۱۶۴۸۸

داده‌کاوی با زبان برنامه‌نویسی R

مؤلفان:

دکتر وحید رضائی تبار

عضو هیئت علمی دانشگاه علامه طباطبائی

شقایق خودکامه

ویرایش دوم، چاپ اول



داده کاوی با زبان برنامه نویسی R

مؤلفان:

دکتر وحید رضائی تبار
شقایق خودکامه

زیر نظر معاونت پژوهشی دانشگاه

شابک: ۹۷۸-۹۶۴-۲۱۷-۶۰۹-

سرویراستار: دکتر نازبلا فرمانی انوشه صفحه آر: دکتر وحید رضائی تبار، شقایق خودکامه

طراح جلد: سمیرا حاجی گلدبی
ناظر فنی: عبدالرضا گودرزی

مرکز چاپ و انتشارات دانشگاه علامه طباطبائی؛ صندوق پستی: ۱۵۸۱۵/۳۴۸۷

فروشگاه آنلاین انتشارات book.atu.ac.ir

سایت مرکز چاپ و انتشارات press.atu.ac.ir

سایت فروشگاه الکترونیکی mybooket.com/atu

آدرس فروشگاه مرکزی: تهران انتهای بلوار دهکده المپیک، دانشگاه علامه طباطبائی
روبه روی ساختمان مرکزی جنب باجه بانک ملی تلفن: ۰۹۹۰۲۴۹۹۹۸۵ همراه: ۴۸۳۹۲۶۵۹

ویرایش دوم، چاپ اول ۱۴۰۳ شمارگان: ۱۰۰ قیمت: ۴۵۰۰۰ تومان

تقدیم به پدر و مادر عزیزم

و

تقدیم به همسر مهربانم و نور چشمانم ایلیا و آوای

عزیزم

وحید رضائی تبار

دانشیار گروه آمار

دانشگاه علامه طباطبائی

پیشگفتار مؤلفان

کتابی که در پیش رو دارید، دربردارنده‌ی توضیحات مفیدی در رابطه با دانش نوین داده کاوی است. در سال‌های اخیر و در مواجهه با حجم بسیار زیادی از داده‌ها در حوزه‌های مختلف آموزش، بانکداری، بازاریابی و غیره، استفاده از داده کاوی به عنوان یکی از الزامات مطرح شده است. همچنین زبان‌های برنامه‌نویسی زیادی برای پیاده‌سازی الگوریتم‌های داده کاوی نیز مطرح شده که در این کتاب با تأکید بر تحلیل مسائل از دید آماری، از زبان برنامه‌نویسی R استفاده شده است که یکی از بهترین زبان‌های برنامه‌نویسی برای پیاده‌سازی مسائل داده کاوی است.

در فصل اول این کتاب، تاریخچه علم داده کاوی و اهمیت مطالعه آن، نصب و راهاندازی زبان برنامه‌نویسی R ، دستورات پایه‌ای در R و توضیح انواع نمودارها ارائه شده است. از آن جا که آماده‌سازی داده‌ها از اهمیت خاصی برخوردار است و نتیجه نهایی هر آزمایش را تحت تأثیر قرار می‌دهد، در فصل دوم در خصوص آماده‌سازی داده‌ها صحبت شده است. در فصل سوم و چهارم به ترتیب به معرفی و تشریح روش‌های رده‌بندی کلاسیک و روش‌های خوشه‌بندی و انواع الگوریتم‌های آن‌ها پرداخته‌ایم. در فصل‌های پنجم و ششم شبکه بیزی، روش‌های یادگیری شبکه بیزی که شامل یادگیری ساختاری و پارامتری هستند، معرفی شده است. شبکه‌های بیزی نوع خاصی از مدل‌های گرافیکی به نام گراف‌های بدون دور سوار (DAG) هستند که برای تعیین رابطه علت و معلولی به کار می‌روند. در این کتاب، برای هر فصل تمرینات تشریحی مطرح شده است که حل آن‌ها به دانشجویان در درک کامل متن و موضوعات مربوطه کمک می‌کند.

این کتاب خصوصاً به نیازهای آموزشی دانشجویان رشته‌های مختلف آمار، علوم کامپیوتر و مهندسی صنایع نیز پاسخ می‌دهد.

امید است انتشار این کتاب مورد استفاده دانشجویان واقع شده و به عنوان یک منبع مفید درسی به آن‌ها کمک کند و در عین حال مورد قبول اساتید نیز واقع شود و گامی هر چند کوچک در اعتلای دانش داده کاوی در کشور محسوب شود. در انتها ضمن تشکر از مدیریت و دست‌اندرکاران نشر دانشگاهی علامه طباطبائی که در طی مراحل تکمیل این کتاب نهایت همکاری را با مؤلفان داشتند، از خوانندگان گرامی درخواست می‌کنیم که ما را از انتقادها و پیشنهادهای خود در جهت ارتقای کیفی کتاب بهره‌مند سازیم و همچنین سپاسگزار خواهیم بود اگر نظرها، پیشنهادها و سوال‌های احتمالی خود را با ما به آدرس ایمیل vhrezaei@atu.ac.ir در میان بگذارید.

وحید رضائی تبار

شهریور ماه ۱۳۹۹

فهرست مطالب

۱-۱-۱	مقدمه	۱
۱-۲	تاریخچه داده کاوی	۱
۱-۳-۱	تعريف‌های داده کاوی	۱
۱-۴	روش‌های یادگیری مدل در داده کاوی	۵
۱-۵	روش‌های داده کاوی	۵
۱-۶	الگوریتم‌های داده کاوی	۹
۱-۷-۱	کاربردهای داده کاوی	۱۱
۱-۸	آمار و داده کاوی	۱۳
۱-۹	مححدودیت‌های داده کاوی	۱۶
۱-۱۰-۱	ضرورت داده کاوی	۱۶
۱-۱۱-۱	چند مثال در مورد مفهوم داده کاوی	۱۶
۱-۱۲-۱	نرم افزارهای داده کاوی	۱۸
۱-۱۳-۱	زبان برنامه نویسی <i>R</i>	۲۱

۲۲	نصب و راهاندازی	۱-۱۳-۱
۲۴	<i>RStudio</i> زبان برنامه‌نویسی	۱-۱۳-۱
۲۵	نصب و روزامدسانی بسته‌های نرم‌افزاری	۱-۱۳-۲
۲۸	بسته‌های نرم‌افزاری داده‌کاوی به تفکیک فصل	۱-۱۴
۲۹	<i>Rattle</i> بسته نرم‌افزاری	۱-۱۵-۱
۳۱	<i>Rattle</i> فراخوانی مجموعه‌داده‌ها در	۱-۱۵-۱
۳۲	فراخوانی داده‌های پیش‌فرض <i>R</i>	۱-۱۶
۳۲	محاسبات پایه	۱-۱۷-۱
۳۷	دستورهای مهم زبان برنامه‌نویسی <i>R</i>	۱-۱۸-۱
۳۷	بردار	۱-۱۸-۱
۳۷	تابع‌ها در زبان برنامه‌نویسی <i>R</i>	۱-۱۸-۲
۳۹	<i>which()</i>	۱-۱۸-۳
۳۹	<i>if</i> عملگر شرطی	۱-۱۸-۴
۴۰	<i>for</i> حلقه‌ی	۱-۱۸-۵
۴۳	<i>while</i> حلقه‌ی	۱-۱۸-۶
۴۴	<i>sapply</i> دستور	۱-۱۸-۷

۴۴	گراف ۱-۱۸-۸
۴۶	آرایه ۱-۱۸-۹
۴۷	ماتریس ۱-۱۸-۱۰
۴۸	ضرب ماتریسی ۱-۱۸-۱۱
۴۸	سایر عملگرهای ماتریسی ۱-۱۸-۱۲
۵۱	متغیرهای تصادفی ۱-۱۸-۱۳
۵۲	استانداردسازی دادهها ۱-۱۸-۱۴
۵۳	رسم نمودار ۱-۱۸-۱۵
۵۵	نمودار میله‌ای ۱-۱۸-۱۶
۵۸	نمودار دایره‌ای ۱-۱۸-۱۷
۶۱	نمودار نقطه‌ای ۱-۱۸-۱۸
۶۴	نمودارهای سه‌متغیری ۱-۱۸-۱۹
۶۷	نمودار جعبه‌ای ۱-۱۸-۲۰
۷۹	نمودار شاخه و برگ ۱-۱۸-۲۱
۷۰	نمودارهای سری زمانی ۱-۱۸-۲۲
۷۱	نمودارهای موازی ۱-۱۸-۲۳

۷۲	۱-۱۸-۲۴- نمودارهای چندضلعی
۷۳	۱-۱۸-۲۵- نمودارهای ماتریسی
۷۴	۱-۱۸-۲۶- نمودار سه بعدی
۷۵	۱-۱۸-۲۷- فلوچارت
۷۸	۱-۱۸-۲۸- رسم تابع‌های ریاضی
۸۲	۱-۱۸-۲۹- نمودار خطی
۸۳	۱-۱۸-۳۰- نقشه‌های جغرافیایی
۸۵	۱-۱۹- نکته‌ها
۹۴	۱-۲۰- خلاصه‌ی فصل
۹۵	تمرین ...
۹۹	فصل دوم: آماده‌سازی داده
۱۰۰	۲-۱- شناخت انواع داده و خصیصه‌های آن
۱۰۰	۲-۲- شاخص‌های مرکزی و پراکندگی
۱۰۱	۲-۲-۱- فراوانی
۱۰۲	۲-۲-۲- میانگین
۱۰۳	۲-۲-۳- میانه

۱۰۴	۴-۲-۲-۴- مُد
۱۰۵	۲-۲-۵- واریانس
۱۰۶	۲-۲-۶- انحراف معیار
۱۰۷	۲-۳- اهمیت آماده‌سازی داده‌ها
۱۰۸	۲-۴- کارهای عمدۀ در آماده‌سازی و پیش‌پردازش داده‌ها
۱۰۸	۲-۴-۱- پاکسازی داده‌ها
۱۱۷	۲-۴-۱-۱- تحلیل داده‌های گمشده
۱۱۹	۲-۴-۱-۳- مشخص کردن داده‌های دور افتاده
۱۲۱	۲-۴-۱-۳-۱- روش دامنه‌ای
۱۲۳	۲-۴-۱-۳-۲- روش‌های آزمون آماری
۱۲۳	۱- آزمون کیو دیکسون
۱۲۴	۲- آزمون گراب
۱۴۰	۲-۴-۱-۴- رفع مشکل افروندگی داده‌ها
۱۴۱	۲-۴-۲- یکپارچه‌سازی داده‌ها
۱۴۱	۲-۴-۳- تبدیل داده‌ها

۱۴۱	۱-۳-۴-۲ - نرمالسازی داده‌ها
۱۴۶	۲-۳-۴-۲ - تعمیم داده‌ها
۱۴۶	۳-۳-۴-۲ - انبوهش
۱۴۷	۴-۴-۲ - کاهش داده‌ها
۱۴۷	۴-۴-۱-۱ - کاهش بُعد
۱۴۷	۴-۴-۱-۱ - تحلیل مؤلفه‌های اصلی (<i>PCA</i>)
۱۶۷	۴-۴-۲-۱ - فشرده‌سازی داده‌ها
۱۷۲	۴-۴-۳-۱ - گسترش‌سازی
۱۷۳	۵-۲ - خلاصه‌ی فصل
۱۷۳	تمرین ...
۱۷۵	فصل سوم: معرفی روش‌های رده‌بندی کلاسیک
۱۷۶	۱-۳ - مقدمه
۱۷۶	۲-۳ - رگرسیون لوژستیک
۱۷۹	۳-۳ - درخت تصمیم
۱۷۹	۱-۳-۳ - استقراء درخت تصمیم
۱۷۹	۲-۳-۳ - مفهوم‌های اصلی درخت تصمیم

۱۸۰	۳-۳-۳-۳- انواع درخت تصمیم
۱۸۰	۴-۳-۳- مزیت‌های درخت تصمیم
۱۸۱	۵-۳-۳- عیب‌های درخت تصمیم
۱۸۲	۶-۳-۳- نحوه نمایش درخت تصمیم
۱۸۳	۷-۳-۳- معیارهای انتخاب صفت خاص
۱۸۳	۱-۷-۳-۳- آنتروپی
۱۸۴	۲-۷-۳-۳- اطلاع سودمندی
۱۸۵	۳-۷-۳-۳- شاخص جینی
۱۸۶	۸-۳-۳- الگوریتم‌های درخت تصمیم
۱۸۷	۱-۸-۳-۳- ID^3
۱۸۷	۲-۸-۳-۳- الگوریتم ID^4
۱۸۸	۳-۸-۳-۳- الگوریتم ID^5
۱۸۸	۴-۸-۳-۳- الگوریتم ID^5_{hat}
۱۸۸	۵-۸-۳-۳- الگوریتم $C^{4.5}$
۱۸۹	۶-۸-۳-۳- الگوریتم $C \cdot 5$
۱۸۹	۷-۸-۳-۳- الگوریتم $CART$

۱۹۰ ارزیابی درخت تصمیم	-۳-۳-۹
۱۹۱ هرس کردن درخت	-۳-۳-۱۰
۱۹۱ انواع روش‌های هرس کردن	-۳-۳-۱۰-۱
۲۰۹ k -نزدیک‌ترین همسایه	-۴-۳- k
۲۱۰ الگوریتم k -نزدیک‌ترین همسایه (KNN)	-۴-۳-۱-۱
۲۱۵ مزیت و محدودیت الگوریتم k -نزدیک‌ترین همسایه در فرایندهای پیش‌بینی....	-۳-۴-۲-۲
۲۴۷ ماشین بردار پشتیبان	-۳-۳-۵-۵
۲۵۰ حالت جدایی‌پذیر در ماشین بردار پشتیبان	-۳-۳-۵-۱-۱
۲۵۶ ماشین بردار پشتیبان در حالت‌های ناخطي و جدایی‌ناپذیر	-۳-۵-۲-۲
۲۵۶ داده‌های جدایی‌ناپذیر	-۳-۵-۲-۱-۱
۲۵۸ ماشین‌های بردار پشتیبان ناخطي	-۳-۵-۲-۲-۲
۲۶۰ ماشین بردار پشتیبان و شبکه عصبی	-۳-۵-۳-۳
۲۶۱ کاربردهای ماشین بردار پشتیبان	-۳-۵-۴-۴
۲۶۲ نقاط ضعف ماشین‌های بردار پشتیبان	-۳-۵-۵-۵
۲۶۶ روش‌های رده‌بندی ترکیبی	-۳-۶-۶
۲۶۶ الگوریتم انبوهش تصادفی	-۳-۶-۱-۱

۲۷۱	۳-۶-۲- الگوریتم بگینگ
۲۷۳	۳-۶-۳- الگوریتم بوستینگ
۲۷۵	۳-۷- خلاصه‌ی فصل
۳۰۵	تمرین
۳۰۷	فصل چهارم: خوشبندی
۳۰۸	۴-۱- مقدمه
۳۱۰	۴-۲- نقاط قوت روش خوشبندی
۳۱۰	۴-۳- نقاط ضعف روش خوشبندی
۳۱۰	۴-۴- تعیین تعداد خوش
۳۱۱	۴-۵- ارزیابی اعتبار در خوشبندی
۳۱۳	۴-۶- ماتریس تشابه و فاصله
۳۱۳	۴-۶-۱- تابع تشابه
۳۱۴	۴-۶-۲- تابع فاصله
۳۱۴	۴-۶-۳- ماتریس تشابه
۳۱۵	۴-۶-۴- ماتریس فاصله
۳۱۵	۴-۷- روش‌های اصلی خوشبندی

۳۱۵	۴-۱-۷-۱- روشن‌های افزایی
۳۱۶	۴-۱-۱-۱- الگوریتم <i>k-means</i>
۳۲۳	۴-۱-۷-۲- الگوریتم <i>k-medoids</i>
۳۲۸	۴-۲-۷-۲- روشن‌های سلسله مراتبی
۳۳۰	۴-۲-۷-۱- الگوریتم <i>Birch</i>
۳۳۲	۴-۲-۷-۲- الگوریتم <i>chameleon</i>
۳۳۴	۴-۳-۲-۷-۱- الگوریتم <i>AGNES</i>
۳۳۴	۴-۳-۲-۷-۱- خوشبندی با روش <i>single-link</i>
۳۳۵	۴-۲-۳-۲-۷-۱- خوشبندی با روش <i>complete-link</i>
۳۳۶	۴-۳-۲-۷-۱- خوشبندی با روش <i>Average-link</i>
۳۴۳	۴-۲-۷-۴- مقایسه خوشبندی سلسله مراتبی و غیر سلسله مراتبی
۳۴۳	۴-۳-۷-۳- روشن‌های مبتنی بر چگال
۳۴۴	۴-۳-۷-۱-۱- الگوریتم <i>DBSCAN</i>
۳۴۵	۴-۳-۷-۱-۱-۱- مزیت‌ها و عیوب‌های الگوریتم <i>DBSCAN</i>
۳۴۷	۴-۳-۷-۲-۱- الگوریتم <i>OPTICS</i>
۳۵۷	۴-۷-۴- روشن مشبکی مبنا

۳۵۷	۱-۴-۷-۴- الگوریتم <i>STING</i>
۳۶۲	۴-۸- خلاصه فصل
۳۶۳	تمرین
۳۶۵	فصل پنجم: مقدمه‌ای بر شبکه‌های بیزی
۳۶۶	۱-۵- مقدمه
۳۶۶	۲-۵- تاریخچه شبکه‌های بیزی
۳۶۸	۳-۵- معرفی شبکه‌های بیزی
۳۷۰	۴-۵- ساختار شبکه‌های بیزی
۳۷۱	۵-۴-۱- ویژگی مارکوفی
۳۷۱	۵-۵- استدلال با شبکه‌های بیزی
۳۷۱	۵-۵-۱- انواع استنباط
۳۷۳	۵-۵-۲- استدلال احتمال‌ها توسط توزیع احتمال مشترک (توأم)
۳۸۰	۶-۵-۶- انواع یادگیری در شبکه‌های بیزی
۳۸۲	۶-۵-۱- یادگیری ساختاری در شبکه‌های بیزی
۳۸۲	۶-۱-۱-۱- روش‌های مبتنی بر قید
۳۸۳	۶-۱-۲- روش‌های جستجو- امتیازدهی

۳۸۶	۱-۲-۱-۶-۵- روش‌های جست‌وجو
۳۸۶	۳-۱-۶-۵- روش ترکیبی
۳۸۷	۴-۲-۶-۵- یادگیری پارامتری شبکه بیزی
۳۸۸	۵-۷- نرم‌افزارهای مدل‌سازی شبکه بیزی
۳۸۹	۵-۸- خلاصه‌ی فصل
۳۸۹	تمرین
۳۹۱	فصل ششم: یادگیری ساختاری و پارامتری شبکه‌های بیزی
۳۹۲	۶-۱- مقدمه
۳۹۳	۶-۲- معرفی چند الگوریتم مهم مبتنی بر قید
۳۹۳	۱- الگوریتم PC
۳۹۵	۲- ساختن کالبد DAG
۴۰۱	۳- الگوریتم GS
۴۰۷	۴- الگوریتم $IAMB$
۴۰۹	۵- الگوریتم $Inter\ IAMB$
۴۱۱	۶-۳- الگوریتم‌های امتیازدهی
۴۱۳	۶-۱-۳- معرفی و طبقه‌بندی تابع‌های امتیاز شبکه بیزی

۶-۲-۳-۶- تابع‌های امتیاز وابسته به توزیع پیشینی ۴۱۴
۶-۳-۶- مترهای مبتنی بر مفهوم‌های نظریه اطلاع ۴۱۹
۶-۴-۳-۶- الگوریتم‌های مبتنی بر جست‌وجوی حریصانه در یادگیری ساختار ۴۲۰
۶-۵-۳-۶- الگوریتم K^2 ۴۲۱
۶-۶-۳-۶- الگوریتم جست‌وجوی ممنوع ۴۲۳
۶-۴- الگوریتم‌های ترکیبی ۴۲۶
۶-۵- یادگیری شبکه‌های بیزی ۴۳۲
۶-۶-۱- ساختار معلوم و داده‌های کامل ۴۳۳
۶-۶-۱-۱- براورد پیشینه درست‌نمایی ۴۳۳
۶-۶-۱-۲- براورد بیزی ۴۳۶
۶-۶-۲- ساختار معلوم و داده‌ها ناکامل ۴۳۸
۶-۶-۳- ساختار نامعلوم و داده‌ها کامل ۴۳۸
۶-۶-۴- ساختار نامعلوم و داده‌ها ناکامل ۴۳۸
۶-۶-۵- قضیه بیز ۴۳۹
۶-۷- رده‌بند بیزی ساده ۴۴۰
۶-۷-۱- خصیصه‌های رده‌بندی بیزی ساده ۴۴۲

۶-۸- هموارسازی	۴۴۳
۶-۹- خلاصه‌ی فصل	۴۷۲
تمرین	۴۷۵
واژه‌نامه فارسی – انگلیسی	۴۷۷
نامنامه	۴۹۰
مراجع ها	۴۹۲