

Bala Deshpande | Vijay Kotu نویسنده:

ویراست دوم

# علوم داده

مفاهیم و شیوه عمل

به همراه آموزش نرم افزار RapidMiner

مترجمان:

دکتر رامین مولاناپور

(مدرس دانشگاه‌های تهران)

مهندس عادل قسمتی

## علوم داده «مفاهیم و شیوه عمل»

به همراه آموزش نرم افزار RapidMiner

متelman: دکتر رامین مولاناپور، مهندس عادل قسمتی

ناشر: انتشارات آتی نگر

ناشر همکار: انتشارات وینا

طراحی جلد و صفحه‌آرایی: همتا بیداریان

چاپ اول، ۱۴۰۲

شمارگان: ۱۰۰ نسخه

قیمت: ۳۴۰۰۰۰۰ ریال

شابک: ۹۷۸-۶۲۲-۷۵۷۱-۴۸-۶

ISBN: 978-622-7571-48-6

حق چاپ برای انتشارات آتی نگر محفوظ است.

نشانی دفتر فروش: خیابان جمالزاده جنوبی، رو به روی کوچه رشتچی، پلاک ۱۴۴، واحد ۱

تلفن: ۰۶۵۶۵۳۳۶-۸ نمبر: ۰۶۵۶۵۳۳۷



[www.ati-negar.com](http://www.ati-negar.com) \* [info@ati-negar.com](mailto:info@ati-negar.com)

سرشناسه: کوتو، ویجی، Vijay, Kotu

علوم داده «مفاهیم و شیوه عمل» به همراه آموزش نرم افزار RapidMiner؛ نویسنده: ویجی کوتو، بالا دشپند /

متelman: رامین مولاناپور، عادل قسمتی

تهران: آتی نگر، وینا

۱۴۰۲ ص.: مصور، جدول، نمودار.

ISBN: 978-622-757-48-6

فیبا.

یادداشت: عنوان اصلی کتاب: Data science : concepts and practice, Second edition, 2019

موضوع: داده کاوی - Data mining

موضوع: مصرف کنندگان - رفتار - Consumer behavior - داده پردازی - Electronic data processing

شناسه افزوده: دشپند، بالا، Deshpande, Bala

شناسه افزوده: مولاناپور، رامن، ۱۳۵۲ - مترجم

شناسه افزوده: قسمتی، عادل، ۱۳۶۹ - مترجم

شناسه افزوده: بیداریان، همتا، ۱۳۶۱

ردیبندی کنگره

ردیبندی دیویس

شماره کتابشناسی ملی

QA76/9

۰۰۶/۳۱۲

۹۱۸۴۴۴۴۳

# فهرست مطالب

۱۱	پیشگفتار
۱۵	مقدمه
۱۹	قدراتیها
۲۱	فصل اول: مقدمه
۲۲	۱- هوش مصنوعی، یادگیری ماشین و علوم داده
۲۵	۱-۲ علوم داده چیست؟
۲۵	۱-۲-۱ استخراج الگوهای معنادار
۲۵	۱-۲-۲ ایجاد مدل‌های معرف
۲۶	۱-۲-۳ ترکیب آمار، یادگیری ماشین و رایانش
۲۷	۱-۲-۴ الگوریتم‌های یادگیری
۲۷	۱-۲-۵ زمینه‌های مرتبط
۲۹	۱-۳ دلایل مستدل در رابطه علوم داده
۲۹	۱-۳-۱ حجم
۲۹	۱-۳-۲ ابعاد
۳۰	۱-۳-۳ سوالات پیچیده
۳۱	۱-۴ طبقه‌بندی علوم داده
۳۳	۱-۵ الگوریتم‌های علوم داده
۳۵	۱-۶ نقشه راه برای این کتاب
۳۵	۱-۶-۱ شروع کار با علوم داده
۳۶	۱-۶-۲ تمرین با استفاده از RapidMiner
۳۶	۱-۶-۳ الگوریتم‌های اصلی
۴۱	منابع
۴۳	فصل دوم: فرایند علوم داده
۴۶	۲-۱ دانش پیشین
۴۶	۲-۱-۱ هدف
۴۷	۲-۱-۲ حوزه موضوعی
۴۷	۲-۱-۳ داده‌ها
۴۹	۲-۱-۴ علیت در برایر همبستگی

۴۹.....	آماده سازی داده ها.....	۲-۲
۵۰.....	کاوش داده ها.....	۲-۲-۱
۵۰.....	کیفیت داده ها.....	۲-۲-۲
۵۱.....	مقادیر ناموجود.....	۲-۲-۳
۵۲.....	انواع داده ها و تبدیل .....	۲-۲-۴
۵۲.....	تبدیل .....	۲-۲-۵
۵۲.....	نقاط پرت .....	۲-۲-۶
۵۳.....	انتخاب ویژگی .....	۲-۲-۷
۵۳.....	نمونه برداری داده ها .....	۲-۲-۸
۵۴.....	مدل سازی .....	۲-۳
۵۴.....	مجموعه داده های آموزشی و آزمایشی .....	۲-۳-۱
۵۶.....	الگوریتم های یادگیری .....	۲-۳-۲
۵۷.....	ارزیابی مدل .....	۲-۳-۳
۵۹.....	مدل سازی جمعی .....	۲-۳-۴
۵۹.....	کاربرد .....	۲-۴
۵۹.....	آمادگی تولید .....	۲-۴-۱
۶۰.....	یکپارچه سازی فنی .....	۲-۴-۲
۶۰.....	زمان پاستخگویی .....	۲-۴-۳
۶۱.....	نوسازی مدل .....	۲-۴-۴
۶۱.....	همگون سازی .....	۲-۴-۵
۶۱.....	دانش .....	۲-۵
۶۲.....	منابع .....	

۶۵

### فصل سوم: کاوش داده ها

۶۶.....	اهداف کاوش داده ها.....	۳-۱
۶۶.....	مجموعه داده ها .....	۳-۲
۶۸.....	انواع داده ها .....	۳-۲-۱
۶۹.....	آمار توصیفی .....	۳-۳
۷۰.....	کاوش تک متغیره .....	۳-۳-۱
۷۳.....	کاوش چندمتغیره .....	۳-۳-۲
۷۵.....	تصور سازی داده ها .....	۳-۴
۷۶.....	تصور سازی تک متغیره .....	۳-۴-۱
۸۱.....	تصور سازی چندمتغیره .....	۳-۴-۲
۸۴.....	تصور سازی داده های با ابعاد بالا .....	۳-۴-۳

۹۰	۳-۳ نقشه راه برای کاوش داده‌ها
۹۱	منابع

۹۲	<b>فصل چهارم: طبقه‌بندی</b>
۹۴	۴-۱ درختان تصمیم
۹۴	۴-۱-۱ نحوه کار آن
۱۰۳	۴-۱-۲ نحوه پیاده‌سازی
۱۱۶	۴-۱-۳ نتیجه‌گیری
۱۱۷	۴-۲ استنتاج قوانین
۱۲۱	۴-۲-۱ نحوه کار آن
۱۲۴	۴-۲-۲ نحوه پیاده‌سازی
۱۲۸	۴-۲-۳ نتیجه‌گیری
۱۲۸	۴-۳ الگوریتم $k$ -نزدیک‌ترین همسایه
۱۳۰	۴-۳-۱ نحوه کار آن
۱۳۸	۴-۳-۲ نحوه پیاده‌سازی
۱۴۱	۴-۳-۳ نتیجه‌گیری
۱۴۲	۴-۴ بیز ساده
۱۴۴	۴-۴-۱ نحوه کار آن
۱۵۳	۴-۴-۲ نحوه پیاده‌سازی
۱۵۵	۴-۴-۳ نتیجه‌گیری
۱۵۶	۴-۵ شبکه‌های عصبی مصنوعی
۱۶۰	۴-۵-۱ نحوه کار آن
۱۶۳	۴-۵-۲ نحوه پیاده‌سازی
۱۶۶	۴-۵-۳ نتیجه‌گیری
۱۶۷	۴-۶ ماشین‌های بردار پشتیبان
۱۷۱	۴-۶-۱ نحوه کار آن
۱۷۳	۴-۶-۲ نحوه پیاده‌سازی
۱۸۰	۴-۶-۳ نتیجه‌گیری
۱۸۱	۴-۷ یادگیرنده‌های جمعی
۱۸۳	۴-۷-۱ نحوه کار آن
۱۸۵	۴-۷-۲ نحوه پیاده‌سازی
۱۹۴	۴-۷-۳ نتیجه‌گیری
۱۹۵	منابع

## فصل پنجم: روش‌های رگرسیون

۱۹۹	۵-۱ رگرسیون خطی
۲۰۰	۵-۱-۱ نحوه کار آن
۲۰۱	۵-۱-۲ نحوه پیاده‌سازی
۲۰۷	۵-۱-۳ نقاط وارسی
۲۱۴	۵-۲ رگرسیون لجستیک
۲۲۰	۵-۲-۱ نحوه کار آن
۲۲۲	۵-۲-۲ نحوه پیاده‌سازی
۲۲۸	۵-۲-۳ خلاصه نکات
۲۳۱	۵-۳ نتیجه‌گیری
۲۳۲	منابع

## فصل ششم: تحلیل وابستگی

۲۳۳	۶-۱ قوانین کاوش وابستگی
۲۳۵	۶-۱-۱ مجموعه اقلام
۲۳۷	۶-۱-۲ تولید قانون
۲۴۰	۶-۲ الگوریتم Apriori
۲۴۱	۶-۲-۱ نحوه کار آن
۲۴۲	۶-۲-۲ نحوه کار آن
۲۴۶	۶-۳ الگوریتم رشد الگوهای پر تکرار
۲۴۶	۶-۳-۱ نحوه کار آن
۲۴۷	۶-۳-۲ نحوه پیاده‌سازی
۲۵۰	۶-۴ نتیجه‌گیری
۲۵۵	منابع

## فصل هفتم: خوشه‌بندی

۲۵۷	۷-۱ خوشه‌بندی K-MEANS
۲۶۳	۷-۱-۱ نحوه کار آن
۲۶۵	۷-۱-۲ نحوه پیاده‌سازی
۲۷۱	۷-۲ خوشه‌بندی DBSCAN
۲۷۶	۷-۲-۱ نحوه کار آن
۲۷۷	۷-۲-۲ نحوه پیاده‌سازی
۲۸۱	۷-۳ نقشه‌های خودسازمان دهنده
۲۸۴	۷-۳-۱ نحوه کار آن
۲۸۶	۷-۳-۲ نحوه پیاده‌سازی
۲۸۹	

## فصل هشتم: ارزیابی مدل

۲۹۹	۸-۱ ماتریس درهم ریختگی
۳۰۱	۸-۲ AUC و ROC
۳۰۳	۸-۳ منحنی‌های برا
۳۰۶	۸-۴ نحوه پیاده‌سازی
۳۱۰	۸-۵ نتیجه‌گیری
۳۱۵	منابع

## فصل نهم: متن کاوی

۳۲۰	۹-۱ نحوه کار آن
۳۲۰	۹-۱-۱ تکرار عبارت- معکوس تکرار در سند
۳۲۲	۹-۱-۲ مجموعه اصطلاحات
۳۲۷	۹-۲ نحوه پیاده‌سازی
۳۲۷	۹-۲-۱ پیاده‌سازی ۱: خوشبندی کلمات کلیدی
۳۳۲	۹-۲-۲ پیاده‌سازی ۲: پیش‌بینی جنسیت نویسنده‌گان و بlags
۳۴۱	۹-۳ نتیجه‌گیری
۳۴۲	منابع

## فصل دهم: یادگیری ژرف

۳۴۷	۱۰-۱ زمستان هوش مصنوعی
۳۵۳	۱۰-۲ نحوه کار آن
۳۵۴	۱۰-۲-۱ مدل‌های رگرسیون به عنوان شبکه‌های عصبی
۳۵۶	۱۰-۲-۲ گرادیان نزولی
۳۶۰	۱۰-۲-۳ نیاز به پس انتشار
۳۶۱	۱۰-۲-۴ طبقه‌بندی بیش از ۲ دسته: تابع بیشینه هموار (Softmax)
۳۶۳	۱۰-۲-۵ شبکه‌های عصبی کانولوشن
۳۷۰	۱۰-۲-۶ لایه متراکم
۳۷۱	۱۰-۲-۷ لایه حذف تصادفی
۳۷۱	۱۰-۲-۸ شبکه‌های عصبی تراجی
۳۷۳	۱۰-۲-۹ خودمرمزگذارها
۳۷۴	۱۰-۲-۱۰ مدل‌های هوش مصنوعی مرتبط
۳۷۵	۱۰-۳ نحوه پیاده‌سازی

۳۸۱.....	۱۰-۴ نتیجه‌گیری
۳۸۱.....	منابع

۳۸۳	فصل یازدهم: موتورهای پیشنهاددهنده
۳۸۷.....	۱۱-۱ مفاهیم موتور پیشنهاددهنده
۳۹۲.....	۱۱-۱-۱ انواع موتورهای پیشنهاددهنده
۳۹۵.....	۱۱-۱-۲ پالایش مشارکتی
۳۹۶.....	۱۱-۲-۱ روش‌های مبتنی بر همسایگی
۴۰۹.....	۱۱-۲-۲ تجزیه عاملی ماتریس
۴۱۶.....	۱۱-۳ پالایش مبتنی بر محظوظا
۴۱۹.....	۱۱-۳-۱ محاسبه نمایه کاربر
۴۲۷.....	۱۱-۳-۲ مدل‌های یادگیری با ناظر
۴۳۳.....	۱۱-۴ پیشنهاددهندهای ترکیبی
۴۳۵.....	۱۱-۵ نتیجه‌گیری
۴۳۸.....	منابع

۴۴۱	فصل دوازدهم: پیش‌بینی سری‌های زمانی
۴۴۶.....	۱۲-۱ تجزیه سری‌های زمانی
۴۴۹.....	۱۲-۱-۱ تجزیه کلاسیک
۴۵۰.....	۱۲-۱-۲ نحوه پیاده‌سازی
۴۵۳.....	۱۲-۲ روش مبتنی بر هموارسازی
۴۵۳.....	۱۲-۲-۱ روش‌های پیش‌بینی ساده
۴۵۵.....	۱۲-۲-۲ هموارسازی نمایی
۴۵۸.....	۱۲-۲-۳ نحوه پیاده‌سازی
۴۶۰.....	۱۲-۳ روش‌های مبتنی بر رگرسیون
۴۶۱.....	۱۲-۳-۱ رگرسیون
۴۶۲.....	۱۲-۳-۲ رگرسیون با تغییرات فصلی
۴۶۵.....	۱۲-۳-۳ میانگین متحرک یکپارچه خودهمبسته
۴۷۲.....	۱۲-۳-۴ ARIMA فصلی
۴۷۵.....	۱۲-۴ روش‌های یادگیری ماشین
۴۷۶.....	۱۲-۴-۱ پنجره‌بندی
۴۷۸.....	۱۲-۴-۲ شبکه عصبی خودهمبسته
۴۸۵.....	۱۲-۵ ارزیابی عملکرد
۴۸۵.....	۱۲-۵-۱ مجموعه داده اعتبارسنجی

۴۸۸ .....	۱۲-۵-۲ اعتبارسنجی پنجره کشویی
۴۸۸ .....	۱۲-۶ نتیجه‌گیری
۴۸۹ .....	۱۲-۶-۱ پیش‌بینی به روش‌ها
۴۹۰ .....	منابع

## فصل سیزدهم: تشخیص ناهنجاری

۴۹۱ .....	۱۳-۱ مفاهیم
۴۹۲ .....	۱۳-۱-۱ علل نقاط پرت
۴۹۲ .....	۱۳-۱-۲ تکنیک‌های تشخیص ناهنجاری
۴۹۵ .....	۱۳-۲ تشخیص نقطه پرت بر اساس فاصله
۴۹۸ .....	۱۳-۲-۱ نحوه کار آن
۴۹۹ .....	۱۳-۲-۲ نحوه پیاده‌سازی
۴۹۹ .....	۱۳-۳ تشخیص نقطه پرت بر اساس چگالی
۵۰۳ .....	۱۳-۳-۱ نحوه کار آن
۵۰۳ .....	۱۳-۳-۲ نحوه پیاده‌سازی
۵۰۶ .....	۱۳-۴ عامل پرت محلی
۵۰۶ .....	۱۳-۴-۱ نحوه کار آن
۵۰۷ .....	۱۳-۴-۲ نحوه پیاده‌سازی
۵۰۹ .....	۱۳-۵ نتیجه‌گیری
۵۱۰ .....	منابع

## فصل چهاردهم: انتخاب ویژگی

۵۱۱ .....	۱۴-۱ طبقه‌بندی روش‌های انتخاب ویژگی
۵۱۳ .....	۱۴-۲ تحلیل مؤلفه اصلی
۵۱۴ .....	۱۴-۲-۱ نحوه کار آن
۵۱۵ .....	۱۴-۲-۲ نحوه پیاده‌سازی
۵۱۷ .....	۱۴-۳ پالایش مبتنی بر نظریه اطلاعات
۵۲۲ .....	۱۴-۴ پالایش مبتنی بر خود
۵۲۴ .....	۱۴-۵ انتخاب ویژگی نوع روکش
۵۲۷ .....	۱۴-۵-۱ حذف پرسو
۵۲۹ .....	۱۴-۶ نتیجه‌گیری
۵۳۲ .....	منابع
۵۳۳ .....	

## فصل پانزدهم: شروع کار با RapidMiner

۵۳۵	۱۵-۱ رابط کاربر و مجموعه اصطلاحات
۵۳۶	۱۵-۲ ابزارهای وارد کردن و خروجی گرفتن از دادهها
۵۴۲	۱۵-۳ ابزارهای مصورسازی دادهها
۵۴۵	۱۵-۴ ابزارهای تبدیل داده
۵۴۸	۱۵-۵ ابزارهای نمونه برداری و مقدار ناموجود
۵۵۳	۱۵-۶ ابزارهای بهینه سازی
۵۵۷	۱۵-۷ یکپارچگی با R
۵۶۴	۱۵-۸ نتیجه گیری
۵۶۵	منابع

## ۵۶۷ مقایسه الگوریتم های علوم دادهها

۵۶۷	رگرسیون: پیش بینی یک متغیر هدف کمی رسته ای
۵۶۹	رگرسیون: پیش بینی یک متغیر هدف کمی عددی
۵۶۹	تحلیل وابستگی: فرایند بدون ناظر برای یافتن روابط بین اقلام
۵۷۰	خوشه بندی: فرایندی بدون ناظر برای یافتن گروه های با معنی در دادهها
۵۷۱	تشخیص ناهنجاری: تکنیک های با ناظر و بدون ناظر برای یافتن نقاط پرت در دادهها
۵۷۲	یادگیری ژرف: آموزش با استفاده از لایه های متعدد بازنمایی دادهها
۵۷۳	پیشنهاددهندها: یافتن ارجحیت کاربر برای یک قلم
۵۷۴	پیش بینی سری های زمانی: پیش بینی مقدار آتی یک متغیر
۵۷۵	انتخاب ویژگی: انتخاب مهم ترین صفات خاصه