

# دانش‌یابی داده‌ها

مقدمه‌ای بر داده‌کاوی

نویسنده:

دانیل ت. لاروس

مدیر گروه داده‌کاوی دانشگاه کانک تیکات آمریکا

برگردان:

علی زینل همدانی

استاد دانشکده مهندسی صنایع و سیستمها

دانشگاه صنعتی اصفهان

مهندس فرهاد ابراهیمیان

مهندس حدیث یعقوبزاده



انتهایی اندیشه صنعتی اسلامی

گروه فنی و مهندسی ۵۱

شماره کتاب ۱۳۳

### دانش‌یابی داده‌ها مقدمه‌ای بر داده‌کاوی

دانلیل ت. لاروس	.....	پدیدآور.....
علی زینل همدانی، فرهاد ابراهیمیان، حدیث یعقوبزاده	.....	برگردان.....
آتوسا سعادتی	.....	ویراستار ادیب.....
زحل شیروانی	.....	صحنه آی و طراح جلد.....
چاپخانه دانشگاه صنعتی اصفهان	.....	لیتوگرافی، چاپ و صحافی.....
انتشارات دانشگاه صنعتی اصفهان	.....	ناشر.....
تابستان ۱۳۹۷	.....	چاپ دوم.....
جلد ۳۰۰	.....	شمارگان.....
۹۷۸-۹۶۴-۸۴۷۶-۹۰-۳	.....	شابک.....
۱۵۰۰۰ ریال	.....	قیمت.....

سرشناسه : لارز، دلیل نی، / - م. ۱۹۵ - م. Daniel T. Larose :

عنوان و نام پدیدآور : دانش‌یابی داده‌ها: مقدمه‌ای بر داده‌کاوی  
مشخصات نشر : اصفهان: دانشگاه صنعتی ام‌همان، انتشارات، ۱۳۹۳.  
مشخصات ظاهری : هشت، ۲۹۵ ص. + یک ۱ جلد.  
شابک : ۹۷۸-۹۶۴-۸۴۷۶-۹۰-۳

وضعیت فهرست نویسی : فیبای مختصر  
یادداشت : این مدرک در آدرس <http://opac.nlai.ir> قابل دسترسی است.  
عنوان اصلی : یادداشت

Discovering Knowledge in Data and introduction to data mining

شناسه افزوده	: همدانی، علی زینل، ۱۳۳۱ - مترجم
شناسه افزوده	: ابراهیمیان، فرهاد، ۱۳۵۵ - مترجم
شناسه افزوده	: یعقوبزاده، حدیث، ۱۳۶۶ - ، مترجم
شماره کتابشناسی ملی	: ۳۷۷۴۳۳۵

حق چاپ برای انتشارات دانشگاه صنعتی اصفهان محفوظ است.

اصفهان: دانشگاه صنعتی اصفهان - انتشارات - کدپستی ۸۴۱۵۶-۸۳۱۱۱ تلفن: (۰۳۱) ۳۳۹۱۲۵۰۹  
دورنگار: ۳۳۹۱۲۵۲۸ برای خرید اینترنتی کلیه کتاب‌های منتشره انتشارات می‌توانید به وبگاه <http://publication.iut.ac.ir> مراجعه و یا مستقیماً از کتابفروشی انتشارات واقع در کتابخانه مرکزی دانشگاه صنعتی اصفهان (تلفن ۳۳۹۱۳۹۵۲) خریداری فرمائید.

## پیشنهاد مترجمین

امروزه بست واقعیت‌نابار صورت برنامه‌ریزی نمده و یا روابط عادت در بهداشت‌های اقتصادی، اجتماعی، علمی، هنری و فرهنگی امری خادی محظوظ می‌گردید و تجمع این بهداوهای دست گردیده است که مکملان و برنامه‌ریزان را ترغیب نماید تا از ابوجه عموماً مطمئن آنها انشی اتحاج شود که بولند بسیاری از این اثناخت و واقعیت‌های پنهان شده در پیشه‌های مختلف را در قالب اطلاعات اتحاج نمایند.

سال ناتناروش تحملی داده‌ها و اتحاج اطلاعات استفاده از روش‌های آماری بوده است که لزوماً با جم زیادی از داده‌ها رو بوده‌اند و حتی بصنایعی دسترسی سریع بنتایج وبالا بودن قوان تجزیه داده‌ها از برآوردهای بودن جم زیاد آنها بوده‌اند انتخاب نمونه خواری می‌گردیده است.

روش‌های داده‌کاوی که با گنریتی نوبه مسد و اسخراج داشت نهفته در دل داده‌های پردازده در حال حاضر به عنوان  
هم‌ترین فناوری جست‌بهره‌برداری موثر از داده‌های چشم محبوب گردیده و ایمت آن به عنوان علم تجزیه و تحلیل داده‌ها  
الگوف داشت نهفته در آنها را به فزونی است

امروزه نظریه روشن‌بودی داده‌کاوی و گفت و انش دکنار روشن‌بودی آماری و بوسی مصنوعی برای بررسی‌های  
دانشگاهی ضروری داشته‌اند اما گنریتی ای شدت‌تر ترجمه‌ی از مناسب‌ترین کتاب منتشر شده در این زمینه را انجام دیم  
به این امید که بتواند مورد اسحاق همت‌آین و انسخوان رشته‌های مختلف تحصیلی قرار گیرد.  
از مردمیت نشر ثورای نشر و عید هنگامه‌ی انتشارات دانشگاه صنعتی اصفهان که قبول زحمت فرموده و د  
جهت چاپ این کتاب را حیات کرده است. هنگامه‌ی انتشار این کتاب می‌شود لازم است به طور خاص از داوران محترم و  
ویراستار ادبی این کتاب سرکار خانم سعادتی و پچنین از سرکاران زمینه را ایجاد کرده است صفحه‌آرایی و تصحیح و طراحی جلد  
کتاب را کشیده‌اند، پاسکننداری شود از خوانندگان محترم تعاضداریم که نسخه نظریه روشن‌بودی که احتملا وجود دارد را در  
بسود چاپ‌های بعدی کتاب را بسیاری فریانند.

به امید توفیق الهی

مترجمین

# فهرست مطالب

یک

## فصل ۱ : مقدمه‌ای بر داده‌کاوی

۳	داده‌کاوی چیست؟
۵	چرا داده‌کاوی؟
۶	نیاز داده‌کاوی به هدایت توسط انسان
۷	فرایند استاندارد میان صنعتی: CRISP-DM
۸	مراحل شش گانه CRISP-DM
۹	مطالعه موردنی اول: تحلیل دعایی ضمانت خودرو
۱۰	برداشت‌های غلط از داده‌کاوی
۱۱	داده‌کاوی چه کارهایی را می‌تواند انجام دهد؟
۱۲	تشریح
۱۳	برآوردهایی
۱۴	پیش‌بینی (پیشگویی)
۱۵	طبقه‌بندی
۱۶	خوشبندی
۱۷	وابستگی

۲۲	مطالعه موردی دوم: پیش‌بینی سود غیرعادی بازار بورس
۲۵	مطالعه موردی سوم: استخراج قوانین انجمنی
۲۷	مطالعه موردی چهارم: پیش‌بینی ورشکستگی شرکت‌ها
۳۰	مطالعه موردی پنجم: توصیف بازار گردشگری
۳۱	منابع
۳۳	تمرین‌ها

## فصل ۲: پیش‌پردازش داده‌ها

۳۶	چرا نیاز مدل پردازش داده‌ها هستیم؟
۳۶	مرتب‌سازی داده‌ها
۳۹	مدیریت داده‌های نمشد
۴۳	شناسایی طبقه‌بندی داده‌ها
۴۴	روش‌های گرافیکی شناسایی داده‌های پرت
۴۶	تبدیل داده‌ها
۴۶	روش نرمال‌سازی Min-Max
۴۸	استانداردسازی به روش Z-Score
۴۹	روش‌های عددی شناسایی داده‌های پرت
۵۱	منابع
۵۱	تمرین‌ها
۵۲	تحلیل عملی

## فصل ۳: تحلیل اکتشافی داده‌ها

۵۴	شناخت مجموعه داده‌ها
۵۷	پرداختن به متغیرهای همبسته
۵۹	کاوش در متغیرهای رسته‌ای
۶۵	استفاده از EDA جهت آشکارسازی فیلدهای غیرعادی
۶۶	کاوش در متغیرهای عددی
۷۶	کاوش در ارتباطات چند متغیره
۷۸	انتخاب زیرمجموعه‌های موردنظر غلبه بر داده‌ها، جهت بررسی بیشتر
۷۹	دسته‌بندی
۸۰	خلاصه
۸۱	منابع
۸۱	تمرین‌ها

## ۸۲ ..... تحلیل عملی

### فصل ۴ : روش‌های آماری برای برآوردهای و پیشگویی

۸۶ ..... وظایف داده‌کاوی برای کشف دانش موجود در داده‌ها
۸۷ ..... روش‌های آماری جهت برآوردهای و پیشگویی
۸۷ ..... روش‌های یک متغیره : معیارهای مرکزی و پراکندگی
۹۰ ..... استیباط آماری
۹۲ ..... تا چه میزان، برآوردهای خود می‌توان اطمینان کرد؟
۹۳ ..... برآوردهای ساصل اطمینان
۹۵ ..... روش‌های دو متغیره : رگرسیون خطی ساده
۱۰۰ ..... خطوط برآوردهایی
۱۰۲ ..... فواصل اطمینان بردار ممکن است $\pm$ به شرط داده شدن $X$
۱۰۲ ..... فواصل پیش‌بینی برای یک مدل اراده صادقی $\pm$ به شرط داده شدن $X$
۱۰۵ ..... رگرسیون چندگانه
۱۰۸ ..... صحه‌گذاری فرض‌های مدل
۱۱۲ ..... منابع
۱۱۲ ..... تمرین‌ها
۱۱۲ ..... تحلیل عملی

### فصل ۵ : الگوریتم $k$ نزدیک‌ترین همسایه

۱۱۶ ..... روش‌های با ناظر و بدون ناظر
۱۱۷ ..... متداول‌تری مدل‌سازی با ناظر
۱۱۹ ..... موازنۀ واریانس - اریبی
۱۲۱ ..... عمل طبقه‌بندی
۱۲۳ ..... الگوریتم $k$ نزدیک‌ترین همسایه
۱۲۶ ..... تابع فاصله
۱۲۹ ..... تابع ترکیب
۱۲۹ ..... رأی گیری غیر وزن‌دار ساده
۱۳۰ ..... رأی گیری وزن‌دار
۱۳۳ ..... ملاحظات پایگاه داده
۱۳۵ ..... انتخاب $k$
۱۳۶ ..... منابع
۱۳۶ ..... تمرین‌ها

## فصل ۶ : درخت‌های تصمیم

۱۴۱	..... درخت‌های طبقه‌بندی و رگرسیون (CART)
۱۴۸	..... الگوریتم C4.5
۱۵۶	..... قوانین تصمیم
۱۵۷	..... مقایسه کاربرد الگوریتم‌های C4.5 و (CART) با استفاده از داده‌های واقعی
۱۶۳	..... منابع
۱۶۳	..... تمرین‌ها
۱۶۴	..... تحلیل عملی

## فصل ۷ : شبکه‌های عصبی

۱۷۱	..... کاربرد شبکه‌های عصبی در برآوردیابی و پیشگویی
۱۷۲	..... مثال ساده‌ای از یک شبکه عصبی
۱۷۴	..... تابع فعال‌سازی سیگموید
۱۷۵	..... پس انتشار خطأ
۱۷۶	..... روش کاهش گرادیان
۱۷۷	..... قوانین پس انتشار خطأ
۱۷۸	..... مثالی از پس انتشار خطأ
۱۸۰	..... شرایط توقف
۱۸۲	..... نرخ یادگیری
۱۸۳	..... اصطلاح اندازه حرکت
۱۸۵	..... تحلیل حساسیت
۱۸۶	..... کاربرد مدل‌سازی شبکه عصبی
۱۸۸	..... منابع
۱۸۸	..... تمرین‌ها
۱۸۹	..... تحلیل عملی

## فصل ۸ : خوشبندی سلسله مراتبی و k-MEANS

۱۹۱	..... خوشبندی
۱۹۵	..... روش‌های خوشبندی سلسله مراتبی
۱۹۶	..... خوشبندی ارتباط منفرد
۱۹۷	..... خوشبندی ارتباط کامل
۱۹۹	..... خوشبندی k-MEANS
۱۹۹	..... مثالی از کاربرد خوشبندی k-MEANS

۲۰۶	.....	انجام خوشبندی k-MEANS با استفاده از نرم افزار
۲۰۹	.....	پیش بینی رویگردانی مشتری با استفاده از عضویت در خوش
۲۱۱	.....	منابع
۲۱۱	.....	تمرین ها
۲۱۲	.....	تحلیل عملی

## فصل ۹ : شبکه های کوهن

۲۱۲	.....	شبکه های خو ساز مانده (SOM)
۲۱۶	.....	شبکه های ک من
۲۱۷	.....	مثالی از خو ش بندی به وسیله شبکه های کوهن
۲۲۲	.....	اعتبار خوش
۲۲۳	.....	کاربرد خوشبندی به وسیله شبکه های کوهن
۲۲۵	.....	تفسیر خوش ها
۲۲۹	.....	مشخصات خوش
۲۳۱	.....	استفاده از خروجی خوشبندی، عنوان وردی سایر مدل های داده کاوی
۲۳۲	.....	منابع
۲۳۳	.....	تمرین ها
۲۳۴	.....	تحلیل عملی

## فصل ۱۰ : قوانین انجمنی

۲۳۵	.....	تحلیل همسایگی و تحلیل سبد بازار
۲۳۸	.....	نمایش داده ها برای تحلیل سبد بازار
۲۳۹	.....	پشتیان، اطمینان، مجموعه اقلام مکرر و الگوریتم استقرایی
۲۴۳	.....	الگوریتم استقرایی چگونه کار می کند؟
۲۴۳	.....	الگوریتم استقرایی چگونه کار می کند؟
۲۴۷	.....	تعمیم داده های صفر یا یک به داده های رسته ای عمومی
۲۴۹	.....	رویکرد نظریه اطلاعات: روش تعیم یافته استنتاج قانون
۲۵۱	.....	کاربرد روش تعیم یافته استنتاج قانون
۲۵۳	.....	چه موقع ناید از قوانین انجمنی استفاده کرد؟
۲۵۷	.....	آیا قوانین انجمنی، یادگیری باناظر یا بدون ناظر را ارائه می کند؟
۲۵۷	.....	الگوهای محلی در مقابل مدل های عمومی
۲۵۹	.....	منابع
۲۵۹	.....	تمرین ها

## تحلیل عملی

### فصل ۱۱ : روش‌های ارزیابی مدل

۲۵۹	تکنیک‌های ارزیابی مدل برای وظیفه تشریح
۲۶۴	تکنیک‌های ارزیابی مدل برای وظایف برآوردهایی و پیشگویی
۲۶۵	تکنیک‌های ارزیابی مدل برای وظیفه طبقه‌بندی
۲۶۷	نرخ خطأ، مشتّت‌های کاذب و منفی‌های کاذب
۲۶۷	تقطیم دینه طبقه‌بندی غلط، برای انعکاس مشکلات دنیای واقعی
۲۷۰	تحلیل زینه /سود تصمیم
۲۷۲	نمودارهای ترفیع و نمودارهای بهره
۲۷۳	ترکیب اوزن‌بایو / ساخت مدل
۲۷۷	تلاقی نتایج: کارگری نک مجموعه متناسب از مدل‌ها
۲۷۸	منابع
۲۷۹	تحلیل گروهی
۲۸۳	واژه‌نامه فارسی به انگلیسی
۲۸۹	واژه‌یاب