



به نام خداوند ن و الْقَلَم



علم داده کاربردی

تحلیل داده‌های واقعی با اسکریپت نویسی **پایتون**

تألیف:

ابوالفضل خورشیدی

نیاز دانش

سرشناسه	: خورشیدی، ابوالفضل، ۱۳۶۸ -
عنوان و نام پدیدآور	: علم داده کاربردی تحلیل داده‌های واقعی با اسکریپت‌نویسی بش / تالیف ابوالفضل خورشیدی.
مشخصات نشر	: تهران: نیاز دانش، ۱۳۹۷.
مشخصات ظاهری	: ۳۳۶ص:، جدول، نمودار.
شابک	: 978-600-8906-37-7
وضعیت فهرست‌نویسی	: فیپا
یادداشت	: کتابنامه.
موضوع	: بش (زبان برنامه‌نویسی کامپیوتر)
موضوع	: Bash (Computer program language):
موضوع	: زبان‌های نوشتاری (کامپیوتر) -- آزمون‌ها و تمرین‌ها
موضوع	: Scripting languages (Computer science) -- Examinations, questions, etc.:
رده‌بندی کنگره	: ۱۳۹۷ خ ۹۵ب / QA۷۶/۷۳
رده‌بندی دیویی	: ۰۵/۱۳۳
شماره کتابشناسی ملی	: ۵۴۸۱۵۵۵



نام کتاب	: علم داده کاربردی-تحلیل داده‌های واقعی با اسکریپت نویسی بش
نویسنده	: ابوالفضل خورشیدی
طراح جلد، حروف چینی و صفحه آرایی	: فاطمه سادات کسائیان زاده مهابادی
مدیر اجرایی - ناظر بر چاپ	: حمیدرضا محمد شیرازی - محمد شمس
ناشر	: نیاز دانش
نوبت چاپ	: اول - ۱۳۹۷
شمارگان	: ۱۰۰
قیمت	: ۳۶۰۰۰۰ ریال

ISBN:978-600-8906-37-7

شابک: ۹۷۸-۶۰۰-۸۹۰۶-۳۷-۷

هرگونه چاپ و تکثیر (اعم از زیراکس، بازنویسی، ضبط کامپیوتری و تهیه‌ی CD) از محتویات این اثر بدون اجازه کتبی ناشر ممنوع است. متخلفان به موجب بند ۵ از ماده ۲ قانون حمایت از مؤلفان، مصنفان و هنرمندان تحت پیگرد قانونی قرار می‌گیرند.

آدرس انتشارات: تهران - میدان انقلاب - خیابان ۱۲ فروردین - تقاطع وحید نظری - پلاک ۲۵۵ - طبقه ۱ - واحد ۲

کلیه حقوق این اثر برای مولف محفوظ است.

تماس با انتشارات: ۰۹۱۲۷۰۷۳۹۳۵-۰۹۱۲۷۰۷۳۹۳۵-۰۶۶۴۷۸۱۰۸-۰۶۶۴۷۸۱۰۶-۰۶۶۴۷۸۱۰۶

www.Niye-Danesh.com

مشاوره جهت نشر: ۰۹۱۲ - ۲۱۰۶۷۰۹

فهرست مطالب

پیش‌گفتار.....	۱۵
بخش اول.....	۳۳
فصل اول: آشنایی با یونیکس و ابزارهای پوسته.....	۳۵
۱,۱ مقدمه‌ای بر یونیکس.....	۳۵
۱,۱,۱ تاریخچه یونیکس.....	۳۷
۲,۱,۱ معماری یونیکس.....	۳۸
۳,۱,۱ ساختار فایل یونیکس.....	۳۹
۴,۱,۱ نام مسیر.....	۴۰
۲,۱ پوسته‌ی بَش.....	۴۱
۱,۲,۱ انواع پوسته‌ها.....	۴۱
۲,۲,۱ نصب و اجرای بَش.....	۴۴
۳,۱ شروع کار با خط فرمان.....	۴۵
۱,۳,۱ کار با پوشه‌ها.....	۴۶
۴,۱ حل یک مسئله واقعی با کمک خط فرمان.....	۴۹
۱,۴,۱ روش اول: استفاده از ساختارهای برنامه‌نویسی زبان بَش.....	۴۹
۲,۴,۱ روش دوم: استفاده از یک دستور ساده.....	۵۰
۵,۱ خود آزمایی.....	۵۲
۶,۱ جواب خود آزمایی.....	۵۳
فصل دوم: مبانی اسکریپت‌نویسی پوسته.....	۵۵
۱,۲ دستور چیست؟.....	۵۵
۱,۱,۲ شکل‌های متخلف دستورات.....	۵۶

۵۷.....	۲,۱,۲ سلسه مراتب دستورات و نوع دستور
۵۹.....	۲,۲ آماده‌سازی اولیه پوسته
۶۰.....	۱,۲,۲ راه‌اندازی پوسته بُش
۶۰.....	۲,۲,۲ فایل‌های راه‌انداز بُش
۶۵.....	۳,۲ اسکریپت ساده «سلام، دنیا!»
۶۵.....	۱,۳,۲ آماده‌سازی محیط و ویرایشگر
۶۸.....	۲,۳,۲ اسکریپت اول: چاپ یک پیغام ثابت
۷۰.....	۳,۳,۲ بهبود اسکریپت برای چاپ پیغام‌های پویا
۷۲.....	۴,۳,۲ اهمیت استفاده بجای از نقل‌قول‌ها
۷۳.....	۵,۳,۲ چاپ پیغام به همراه نمایش نام اسکریپت
۷۵.....	۴,۲ خودآزمایی
۷۶.....	۵,۲ جواب خودآزمایی
۷۷.....	فصل سوم: مدیریت سیستم یونیکس
۷۷.....	۱,۳ فایل
۷۸.....	۱,۱,۳ نوع فایل
۷۸.....	۲,۱,۳ فایل مخفی
۷۹.....	کار با فایل‌ها
۸۰.....	۲,۳ پروسه
۸۰.....	۱,۲,۳ پروسه‌های پس‌زمینه و پیش‌زمینه
۸۲.....	۲,۲,۳ فهرست پروسه‌ها
۸۳.....	۳,۲,۳ پروسه‌های والد و فرزند
۸۴.....	۴,۲,۳ سیگنال
۸۵.....	۳,۳ مالکیت و سطوح دسترسی
۸۵.....	۱,۳,۳ مجوز فایل‌ها
۸۶.....	۲,۳,۳ مجوزهای پوشه

۸۶.....	تغییر مجوزهای فایل و پوشه ۳,۳,۳
۸۸.....	خودآزمایی ۴,۳
۹۰.....	جواب خودآزمایی ۵,۳
۹۳.....	فصل چهارم: ورودی و خروجی
۹۴.....	۱,۴ ورودی
۹۴.....	۱,۱,۴ مرور محتویات یک فایل
۹۷.....	۲,۱,۴ خواندن ورودی کاربر
۹۹.....	۳,۱,۴ سند اینجایی
۱۰۳.....	۲,۴ خروجی
۱۰۳.....	۱,۲,۴ چاپ روی ترمینال
۱۰۶.....	۳,۴ جریان‌های ورودی و خروجی استاندارد و بازهدایت
۱۰۷.....	۱,۳,۴ واصف فایل
۱۰۸.....	۲,۳,۴ عملگرهای بازهدایت <, > و >>
۱۱۰.....	۳,۳,۴ دستور exec
۱۱۱.....	۴,۴ خطوط لوله
۱۱۳.....	۵,۴ خودآزمایی
۱۱۴.....	۶,۴ جواب خودآزمایی
۱۱۷.....	فصل پنجم: متغیرها، توابع و جانشینی
۱۱۸.....	۱,۵ متغیرها
۱۱۹.....	۱,۱,۵ ایجاد و حذف متغیر
۱۱۹.....	۲,۱,۵ انواع متغیرها
۱۲۲.....	۳,۱,۵ متغیر عدد صحیح
۱۲۲.....	۴,۱,۵ الصاق مقدار جدید به یک متغیر اسکالر
۱۲۳.....	۵,۱,۵ آرایه‌ها
۱۲۶.....	۲,۵ توابع

۱۲۶	۱,۲,۵	تعریف تابع
۱۲۶	۲,۲,۵	پارامترهای مکانی
۱۲۷	۳,۲,۵	حوزه تعریف
۱۲۹	۴,۲,۵	دستور source
۱۲۹	۳,۵	جانشینی
۱۲۹	۱,۳,۵	نقل قول
۱۳۱	۲,۳,۵	ارزیابی عبارتهای ریاضی
۱۳۴	۳,۳,۵	جانشینی دستور
۱۳۴	۴,۳,۵	جانشینی نام فایل
۱۳۶	۵,۳,۵	جداسازی کلمات
۱۳۷	۶,۳,۵	بسط پارامتر یا متغیر
۱۴۱	۷,۳,۵	متغیرهای ویژه
۱۴۳	۸,۳,۵	زیرپوسته
۱۴۴	۴,۵	خودآزمایی
۱۴۷	۵,۵	جواب خودآزمایی
۱۵۳		فصل ششم: شرطی ها و حلقه ها
۱۵۴	۱,۶	وضعیت خروج و کدهای بازگشت
۱۵۴	۱,۱,۶	آزمایش یک عبارت
۱۵۵	۲,۱,۶	آزمایش عدد صحیح
۱۵۶	۳,۱,۶	آزمایش رشته ها
۱۵۸	۲,۶	شرطی ها
۱۵۸	۱,۲,۶	ساختار if
۱۶۲	۲,۲,۶	ساختار case
۱۶۵	۳,۶	حلقه ها
۱۶۵	۱,۳,۶	حلقه for

۱۶۷ while حلقه ۲,۳,۶
۱۷۰ خودآزمایی ۴,۶
۱۷۲ جواب خودآزمایی ۵,۶
۱۷۳	فصل هفتم: ابزارهای پردازش متن و عبارت با قاعده
۱۷۴ فیلترهای متنی مقدماتی ۱,۷
۱۷۴ head و tail دستورات ۱,۱,۷
۱۷۶ tr دستور ۲,۱,۷
۱۷۷ grep دستور ۳,۱,۷
۱۸۱ عبارت با قاعده ۲,۷
۱۸۲ انواع عبارت‌های با قاعده ۱,۲,۷
۱۸۲ عبارت‌های با قاعده مقدماتی ۲,۲,۷
۱۸۶ تطبیق چندین نویسه با یک عبارت ۳,۲,۷
۱۸۷ عبارت با قاعده تعمیم‌یافته ۴,۲,۷
۱۸۸ ارجاعات برگشتی ۵,۲,۷
۱۹۰ خودآزمایی ۳,۷
۱۹۱ جواب خودآزمایی ۴,۷
۱۹۳	فصل هشتم: ابزارهای پردازش متن پیشرفته sed و awk
۱۹۴ ابزار sed ۱,۸
۱۹۵ آرگومان‌های انتخابی ۱,۱,۸
۱۹۵ دستور جانشینی ۲,۱,۸
۱۹۷ دستور حذف ۳,۱,۸
۱۹۸ دستور چاپ ۴,۱,۸
۱۹۹ ترکیب دستورات ۵,۱,۸
۲۰۰ عبارت‌های باقاعده تعمیم‌یافته ۶,۱,۸
۲۰۱ دستور ترجمه ۷,۱,۸

- ۲,۸ ابزار awk ۲۰۲
- ۱,۲,۸ نحو awk ۲۰۲
- ۲,۲,۸ ویرایش فیلدها ۲۰۳
- ۳,۲,۸ چاپ خروجی قالب‌بندی شده ۲۰۷
- ۴,۲,۸ توابع ۲۰۸
- ۵,۲,۸ اسکریپت‌نویسی awk ۲۰۹
- ۶,۲,۸ متغیرهای درونی ۲۰۹
- ۷,۲,۸ شرطی‌ها و حلقه‌ها ۲۱۰
- ۳,۸ عبارت باقاعده در عمل ۲۱۲
- ۴,۸ خودآزمایی ۲۱۸
- ۵,۸ جواب خودآزمایی ۲۲۰

بخش دوم ۲۲۳

فصل نهم: عادات‌های غذایی مشتریان یک رستوران ۲۲۵

- سؤال (۱) چه تعداد از یک غذای خاص سفارش داده شده است؟ ۲۲۸
- گام اول: پیدا کردن تمام استیک بوریتوهای موجود در فایل ۲۲۹
- گام دوم: شمارش تعداد خطوط ۲۳۰
- سؤال (۲) همراه با سیب‌زمینی سرخ‌شده، بیشتر کدام چاشنی استفاده می‌شود؟ ۲۳۴
- گام اول: استخراج خطوط شامل سفارش‌های سیب‌زمینی سرخ‌شده با سس ۲۳۵
- گام دوم: شمارش تعداد هر نوع سس ۲۳۶

فصل دهم: نظام رتبه‌بندی دانشگاه‌ها ۲۴۱

- سؤال (۱) کالجها چند درصد از این فهرست را تشکیل میدهند؟ ۲۴۴
- گام اول: پیدا کردن تمام کالجها ۲۴۴
- گام دوم: شمارش تعداد کالجها ۲۴۵
- سؤال (۲) کدام ایالت دارای بیشترین تعداد دانشگاه می‌باشد؟ ۲۴۶
- گام اول: جستجو و شمارش دانشگاه‌های ایالت کالیفرنیا ۲۴۶

- گام دوم: استخراج ستون نام دانشگاه‌ها و ایالت‌ها..... ۲۴۶
- گام سوم: مرتب‌سازی ایالت‌ها بر حسب تعداد دانشگاه‌ها ۲۴۷
- سوال (۳) آیا بین رتبه دانشگاه و شهریه تحصیلی آن همبستگی وجود دارد؟..... ۲۴۹

فصل یازدهم: داده‌کاوی شبکه‌های اجتماعی ۲۵۱

- سؤال (۱) از هر نوع پست چه تعداد وجود دارد؟..... ۲۵۳
- گام اول: استخراج ستون نوع پست ۲۵۳
- گام دوم: شمارش تعداد مقادیر یکتا از هر نوع پست ۲۵۴
- سؤال (۲) معروف‌ترین پست این فهرست کدام است؟ ۲۵۵
- گام اول: استخراج ستون‌های عددی ۲۵۵
- گام دوم: محاسبه‌ی مجموع واکنش‌های هر پست ۲۵۶
- گام سوم: مرتب‌سازی ستون مربوط به مجموع واکنش‌ها ۲۵۷
- گام چهارم: جستجوی نوع، عنوان و پیغام معروف‌ترین پست ۲۵۷
- سوال (۳) آیا می‌توان با اجرای یک دستور ساده معروف‌ترین پست را پیدا کرد؟ ۲۵۸
- گام اول: ترکیب دستورات با کمک پایپ ۲۵۸
- گام دوم: نوشتن یک تابع ۲۵۸

فصل دوازدهم: بهترین شهرهای یک کشور بر اساس امنیت ۲۶۱

- سؤال (۱) بدون استفاده از ابزار csvkit، چگونه می‌توان تعداد سطرها و ستون‌ها را شمارش کرد؟..... ۲۶۴
- روش سخت‌تر ۲۶۴
- روش ساده و آسان ۲۶۵
- سؤال (۲) کدام جرم در سراسر کشور دارای بیشترین تعداد ارتکاب بوده است؟ ۲۶۶
- گام اول: مرتب‌سازی ستون مجموع جرائم شهرها ۲۶۶
- گام دوم: چاپ فهرست مرتب‌شده‌ی نام جرائم و تعداد ارتکاب آن‌ها ۲۶۷
- گام سوم: چاپ بیشترین جرم مرتکب‌شده در کشور ۲۶۷
- سؤال (۳) در هر شهر چه نوع جرمی دارای بیشترین تعداد ارتکاب است؟ ۲۶۸
- گام اول: مرتب‌سازی ستون جرائم یک شهر ۲۶۸

- گام دوم: نوشتن یک تابع برای پیدا کردن بیشترین جرم انجام گرفته در یک شهر ۲۶۹
- سؤال (۴) متوسط تعداد جرائم هر شهر چقدر است؟ ۲۷۰
- سؤال (۵) بهترین شهر استرالیا کدام است؟ ۲۷۱
- فصل سیزدهم: کاوشی در آثار ادبی دوره‌ی رنسانس** ۲۷۵
- سؤال (۱) چه تعداد آثار منظوم (اشعار) و یا نمایش نامه وجود دارد؟ ۲۷۷
- گام اول: استخراج خط اول ۲۷۸
- گام دوم: تبدیل ستون‌های جدا شده با ویرگول به سطرهای چندگانه ۲۷۸
- گام سوم: استخراج تمامی آثار از یک نوع خاص ۲۷۹
- گام چهارم: شمارش سطرها ۲۸۰
- سؤال (۲) آثار ادبی هر نویسنده چه تعداد است؟ ۲۸۰
- سه گام در یک خط: ۲۸۱
- سؤال (۳) برای یک کتاب خاص، کدام واژه بیشتر از همه استفاده شده است؟ ۲۸۲
- نوشتن یک تابع: ۲۸۲
- سؤال (۴) برای یک نویسنده‌ی خاص، کدام واژه بیشتر از همه استفاده شده است؟ ۲۸۳
- گام اول: حذف عناوین آثار برای یک نویسنده‌ی خاص ۲۸۴
- گام دوم: استخراج تمامی آثار شکسپیر ۲۸۴
- گام سوم: جمع افقی تمامی ستون‌ها با عنوان یکسان ۲۸۶
- فصل چهاردهم: گزارش وضعیت حساب‌های جاری** ۲۸۹
- سؤال (۱) موجودی حساب جاری بعد از هر تراکنش چقدر است؟ ۲۹۲
- سؤال (۲) موجودی نهایی حساب بعد از پایان تمامی تراکنش‌ها چقدر است؟ ۲۹۶
- سؤال (۳) اگر موجودی حساب کمتر از یک مقدار مشخص شود، چگونه به کاربر هشدار دهیم؟ ۲۹۷
- فصل پانزدهم: قرآن کاوی رایانه‌ای** ۲۹۹
- سؤال (۱) تعداد آیات هر سوره را شمارش کنید ۳۰۳
- گام اول: استخراج ستون‌های مربوط به شماره سوره و شماره آیه ۳۰۳
- گام دوم: شمارش تعداد تکرارهای شماره مربوط به هر سوره ۳۰۳

سوال ۲) نمودار فراوانی کلمات هر سوره را ترسیم کنید ۳۰۴

گام اول: تولید یک فایل داده جامع ۳۰۴

گام دوم: شمارش تعداد کلمات هر سوره ۳۰۸

گام آخر: رسم نمودار تعداد کلمات هر سوره ۳۱۱

سؤال ۳) نمودار فراوانی آیات هر سوره را ترسیم کنید ۳۱۴

گام اول: رسم نمودار تعداد آیات بر اساس ترتیب سوره ۳۱۴

گام دوم: نمودار تعداد آیات بر اساس ترتیب نزول ۳۱۵

سوال ۴) کدام واژگان بیشترین تعداد استفاده را در متن عربی قرآن داشته‌اند؟ ۳۱۶

گام اول: یافتن کلمات با بیشترین تعداد تکرار ۳۱۶

گام دوم: رسم نمودار بیشترین تعداد کلمات در متن عربی قرآن ۳۱۸

فصل شانزدهم: پروژه‌هایی برای تمرین بیشتر ۳۲۱

چالش ۱) حافظ‌پژوهی با رویکرد داده‌کاوی ۳۲۲

چالش ۲) پژوهشی در آثار شکسپیر ۳۲۳

چالش ۳) عرضه و فروش نرم‌افزارهای کاربردی ایرانی برای اندروید ۳۲۴

چالش ۴) گزارش هواشناسی ۳۲۴

چالش ۵) پروازهای لغو شده‌ی یک شرکت هواپیمایی ۳۲۵

چالش ۶) عرضه و تقاضا برای تاکسی درون شهری ۳۲۶

چالش ۷) بررسی حقوق اساتید دانشگاه ۳۲۸

چالش ۸) مطالعه ژنوم‌ها با کمک علم داده ۳۲۸

سخن پایانی ۳۳۱

مراجع و متون پیشنهادی ۳۳۳

پیش‌گفتار

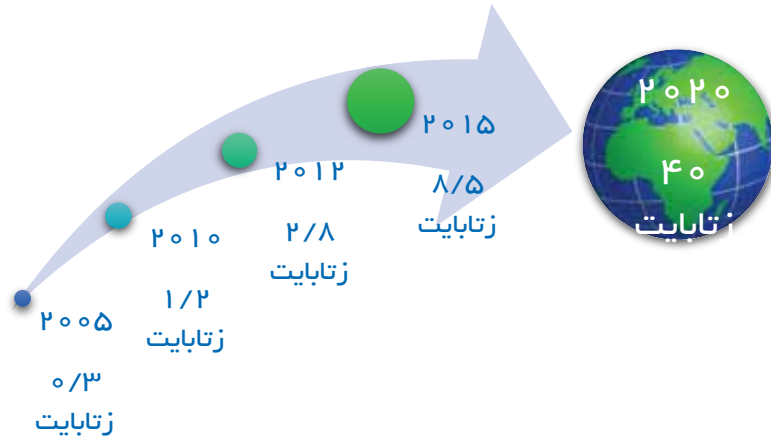
قبل از هر چیز لازم است تا از شما خواننده‌ی گرامی تشکر کنم و جای بسی خرسندی است که تصمیم گرفته‌اید تا بخشی از وقت ارزشمندتان را برای مطالعه‌ی این کتاب سرمایه‌گذاری نمایید. در این کتاب گام به گام همراه شما خواهیم بود تا با رویکردی عملگرایانه به ماجراجویی در دنیای علم داده بپردازیم و بتوانیم بینش‌های ناشناخته را کشف و آینده را پیش‌بینی کنیم. برای آن که علم داده برای تمامی افراد با هر سطح دانش و تجربه‌ی کاری فراگیر شود، لازم است تا در وهله‌ی نخست، داده‌ها را ساده، در دسترس و ملموس کنیم. برای رسیدن به این مهم، در این کتاب اتفاقات و مسائلی را که در کار و زندگی روزمره با آن‌ها سر و کار داریم از دریچه‌ای متفاوت مورد بررسی قرار داده‌ایم و تلاش کرده‌ایم با کمک ابزارهای ساده و کوچک، اما قدرتمند، نگاهی نوین، علمی و مبتنی بر کدنویسی به چالش‌های زندگی واقعی داشته باشیم.

واقعیت این است که بینش‌های ناشناخته به صورت داده‌های خام و بعضاً حجیم و کلان ذخیره می‌شوند. پیش‌بینی آینده‌ی یک کسب و کار نیازمند به کارگیری این داده‌ها می‌باشد. به بیان دیگر، داده‌ها نقش بسیار مهمی در تصمیم‌سازی، مدیریت بهتر منابع و سودآوری برای یک شرکت یا سازمان داشته و هم‌چنین موجب ارائه‌ی خدمات بهتر به مشتریان شده و رضایت و به تبع آن وفاداری مشتری را به همراه خواهند داشت. مطالعات مؤسسه IDC^۱ نشان می‌دهد که تا سال ۲۰۲۰ حدود ۴۰ زتابایت (۴۰ هزار میلیارد گیگابایت)^۲ داده در سراسر دنیا تولید می‌شود که ۳۰۰ برابر بیشتر از حجم داده‌ی دنیا در سال ۲۰۰۵ است [۱]. در این مطالعه اشاره شده است که بیش‌تر این داده‌ها نه توسط انسان بلکه توسط ماشین‌ها تولید خواهد شد. گوشی‌های تلفن همراه، شبکه‌های

^۱ International Data Corporation

^۲ 40 ZB = 40 × 10²¹ bytes

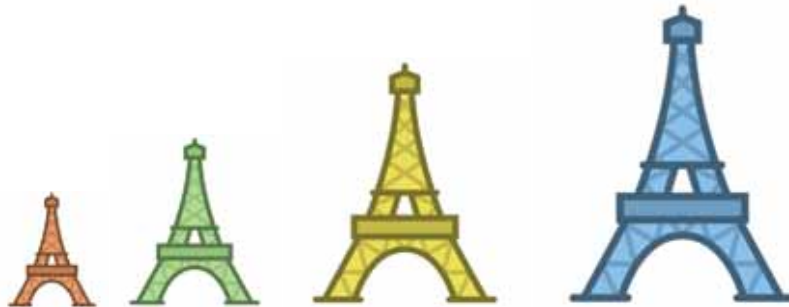
حسگری، دوربین‌ها، دستگاه‌های کنترل تردد، و سائل حمل و نقل، وب سایت‌ها و فرو شگاه‌های اینترنتی، لاگ‌های نرم‌افزاری و ... ماشین‌هایی هستند که در افزایش حجم داده نقش بسزایی دارند.



تعداد سلول‌های بدن انسان حدود ۴۰ هزار میلیارد تخمین زده می‌شود.

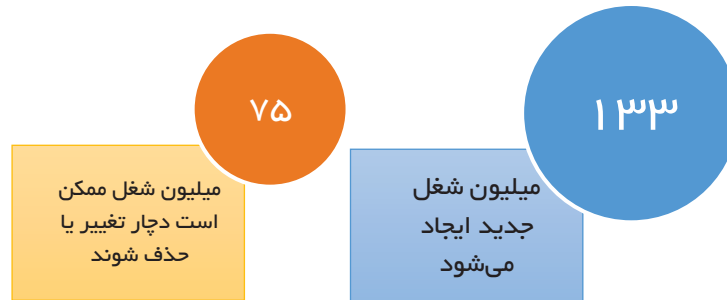
اگر هر بیت را معادل یک سلول بگیریم، ۴۰ زتابایت معادل مجموع سلول‌های کل جمعیت بشر بر روی کره زمین می‌باشد؛ یعنی کل سلول‌های بدن ۸ میلیارد انسان.

شرکت DOMO با تمرکز بر روی داده‌های تولید شده در سراسر دنیا گزارشی منتشر کرده است که نشان می‌دهد این تولیدکنندگان جهانی داده در حال تولید ۲/۵ اگزابایت (۲/۵ میلیارد گیگابایت) داده جدید در هر روز هستند. پیش‌بینی می‌شود این حجم تولید داده با رشد اینترنت اشیا (IoT) شتاب فزاینده‌ای به خود بگیرد [۲].



۲۵ اگزابایت (میلیارد گیگابایت) معادل داده‌ی پر شده بر روی ۱۰ میلیون دیسک بلوری است. اگر این تعداد دیسک را روی یکدیگر قرار دهیم ارتفاع آن ۴ برابر ارتفاع برج ایفل می‌شود.

بر اساس پیش‌بینی‌های مجمع جهانی اقتصاد (WEF)^۱ تا سال ۲۰۲۲، حدود ۱۳۳ میلیون جایگاه شغلی وجود خواهد داشت که از این میان، صنایع مختلف به تحلیل‌گران و دانشمندان داده بیشتر از هر متخصص دیگری نیاز دارند.



دانشمندان و تحلیل‌گران داده و متخصصان هوش مصنوعی و یادگیری ماشین در رتبه‌های اول و دوم فهرست شغل‌های پر تقاضا تا ۲۰۲۲ قرار دارند. همچنین متخصصان کلان‌داده، متخصصان فناوری‌های نوین و هوشمند و متخصصان مخابرات و انتقال دیجیتال در بین ۱۰ شغل برتر تا سال ۲۰۲۲ حضور دارند.

علم داده به عنوان جذاب‌ترین شغل قرن ۲۱ معرفی شده است [۳]. این حوزه با یک نرخ نجومی در حال رشد بوده و شرکت‌ها و سازمان‌های متعددی به دنبال استخدام نیروهای متخصص هستند تا بتوانند معدن طلای داده‌ها را کاوش کرده و بدین ترتیب با کمک داده‌ها موجب هدایت بهتر و مؤثرتر کسب و کار شوند. IBM پیش‌بینی کرده است که تا سال ۲۰۲۰ تعداد جایگاه‌های شغلی مرتبط با علم داده در ایالت متحده آمریکا به میزان ۳۶۴۰۰۰ شغل افزایش یافته و به عدد ۲۰۷۲۰۰۰۰ شغل خواهد رسید.

از سوی دیگر، در سال‌های اخیر با پدیده‌ی رو به رشد دیگری مواجه بوده‌ایم، و همچنان هستیم، و آن ظهور رشته‌ی دانشگاهی علوم داده در دانشگاه‌های معتبر و بزرگ آمریکا چون هاروارد [۴]، استنفورد [۵] و کالیفرنیا برکلی [۶] و ... می‌باشد. در کشورمان نیز حرکت‌های امیدبخشی را شاهد بوده‌ایم و آن تصویب گرایش داده‌کاوی زیرمجموعه رشته‌ی علوم کامپیوتر [۷] و گرایش علوم داده زیرمجموعه رشته‌ی ریاضی کاربردی [۸] است.

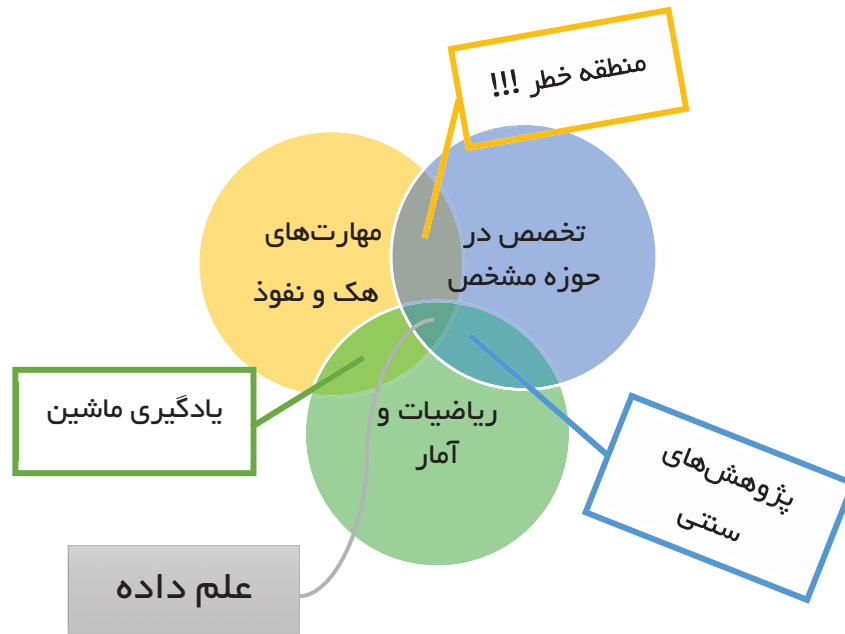
^۱ World Economic Forum



تا سال ۲۰۲۰ تعداد ۳۶۴۰۰۰ شغل جدید در حوزه علم داده در ایالت متحده آمریکا ایجاد می‌شود. حداقل حقوق سالانه ۸۰۲۶۵ دلار بوده و ۸۱ درصد این مشاغل به سابقه کار ۳ تا ۵ سال یا بیشتر نیاز دارند.

بعد از ارائه‌ی این آمار و ارقام در مورد این حوزه علمی و فضای کاری جدید آن، به سراغ تعریف علم داده می‌رویم. واقعیتی که در صنعت وجود دارد این است که علم داده آنقدر جدید است که شرکت‌ها هنوز یک مسیر مشخص و ثابتی برای تربیت نیروی متخصص این حوزه ندارند. این حوزه هنوز در دوران طفولیت خود بسر می‌برد و لذا تعاریف متنوعی از مرزهای آن وجود دارد. شاید مهم‌ترین سوالی که ذهن شما را درگیر نموده است این باشد که چگونه در حوزه علم داده مشغول به کار شویم؟ اجازه دهید این سوال را این‌گونه مطرح کنیم که مباحث نظری و مهارت‌هایی که باید بر آن‌ها مسلط شویم، چیست؟ خوب، پاسخ این است که کمی از همه چیز، یعنی ریاضیات، آمار، یادگیری ماشین، علوم کامپیوتر، هوش تجاری، مدل‌سازی و برنامه‌نویسی و ... و البته یک مهارت اصلی و ضروری و آن تخصص در یک حوزه کاری خاص، برای مثال: اقتصاد، مالی، علوم اجتماعی، بیوانفورماتیک، مهندسی رباتیک، مدیریت سرور، فیزیک، پزشکی و ... نمودار و تعریف شده توسط درو کانوی^۱ به خوبی نشان می‌دهد که علم داده چیست و چرا دانشمندان و تحلیلگران داده از ارزش و اهمیت بالایی در دنیای امروز برخوردار هستند.

^۱ Drew Conway's Data Science Venn Diagram



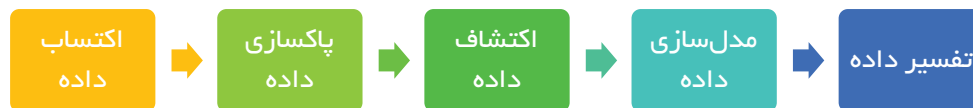
منظور از «منطقه خطر» در این نمودار شاید به هکرها و کدنویسانی اشاره دارد که با نفوذ به سیستم‌های کامپیوتری، امنیت آن‌ها را به مخاطره می‌اندازند.

- با توجه به این نمودار مشخص است که علم داده از تقاطع سه حوزه اصلی تشکیل شده است:
- **ریاضیات و آمار:** استفاده از فرمول‌ها و روابط برای مدل‌سازی مسئله و تحلیل آن.
 - **برنامه‌نویسی:** توانایی کدنویسی و توسعه برنامه‌ها برای ایجاد خروجی مطلوب بر روی کامپیوتر.
 - **دانش وابسته به حوزه کاری:** یعنی همان شناخت و درک درست حوزه‌ی مسئله (رباتیک، پزشکی، مالی، علوم اجتماعی و ...).

خوب است در این جا مقایسه‌ای بین آمار و علم داده داشته باشیم. در واقع به این فکر کنید که آیا علم داده همان آمار است؟ یا آیا علم داده بدون آمار نیز ممکن است؟ در مقاله‌ی «۵۰ سال علم داده» [۹] به اختلاف سلیقه‌ی بین آماردان‌ها و دانشمندان علوم کامپیوتر در مورد واژه‌ی علم داده به تفصیل پرداخته شده است. به نظر برخی آماردان‌ها دانشمند داده صرفاً یک واژه‌ی جذاب برای

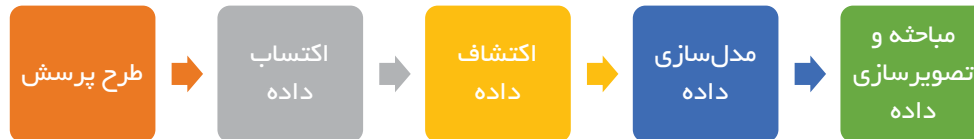
آماردان بوده و علم داده نامی است که در سال‌های اخیر در بازار کار و صنعت برای علم آمار بکار برده می‌شود. با این وجود، همانطور که در این کتاب خواهیم دید، علم داده بدون آمار، نه تنها امکان‌پذیر بوده بلکه لذت‌بخش‌تر خواهد بود. به بیان روشن‌تر علم آمار فقط یک قسمت اندک، اما بسیار مهم، از علوم داده را تشکیل می‌دهد.

رویکردی که در این کتاب در پیش گرفته شده است یک رویکرد کاملاً عملگراییانه می‌باشد. لذا برای این کتاب یک تعریف خیلی کاربردی از علم داده را که توسط میسون و ویگینز در سال ۲۰۱۰ ارائه شده است [۱۰-۱۲]، برگزیده‌ایم. آن‌ها علم داده را بر اساس پنج گام زیر توصیف کرده‌اند: ۱- اکتساب داده، ۲- پاکسازی داده، ۳- اکتشاف داده، ۴- مدل‌سازی داده، ۵- تفسیر داده.



تعریف شماره یک از علم داده در عمل. در این تعریف مرحله «پاکسازی داده» جزو مراحل اصلی فرایند علم داده می‌باشد.

آنچه که بسیار اهمیت دارد این است که حوزه‌ای که مسئله در آن تعریف شده است را به خوبی بشناسیم تا درک درستی از مسئله برای ما حاصل شود. این مرحله شامل مطرح کردن سؤالات مهم و قابل توجه، هم‌فکری و مشورت با افراد متخصص آن حوزه، تهیه مجموعه داده‌های مورد نیاز، تعریف متغیرهای کلیدی و مهم‌تر از همه تعیین هدف پروژه، یعنی آنچه که باید پیش‌بینی یا کشف شود، می‌باشد. لذا می‌توانیم قبل از این مراحل پنج‌گانه یک پیش-مرحله جدید با عنوان طرح سؤالات قابل توجه اضافه کنیم. در [۱۳، ۱۴] این مراحل پنج‌گانه اندکی متفاوت تعریف شده‌اند: ۱- طرح پرسش، ۲- اکتساب داده، ۳- اکتشاف داده، ۴- مدل‌سازی داده، ۵- مباحثه و تصویرسازی داده.



تعریف شماره دو از علم داده در عمل. تفاوت این تعریف با تعریف قبلی اضافه شدن یک پیش-مرحله با عنوان «طرح پرسش» می‌باشد.

این تعریف پاکسازی داده را جزو مراحل اصلی ذکر نکرده است اما همانطور که در [۱۵] بیان شده است ۸۰ درصد کار در یک پروژه علم داده مربوط به مرحله‌ی پاکسازی داده‌ها می‌باشد. از آن جا که خط فرمان یونیکس ابزارهای متعدد و متنوعی را برای این منظور در اختیار ما قرار می‌دهد، در این کتاب این مرحله به تفسیر و در قالب پروژه‌های عملی مختلف توضیح داده می‌شود. بنابراین اگر بگوییم این کتاب لازمه‌ی پیش‌برد بیش از ۸۰ درصد از هر پروژه‌ی علم داده است، سخن به گزافه نگفته‌ایم. خوب است در اینجا به قسمتی از سخنرانی دکتر سید جواد کاظمی تبار در همایش آشنایی با علوم داده که آذرماه ۹۶ در دانشگاه شریف [۱۶] ایراد گردیده است، اشاره کنیم. ایشان در شروع صحبت درباره اهمیت اسکریپت‌نویسی خط فرمان و استفاده از عبارتهای باقاعده چنین گفتند:

«... ذهنیتی که در دانشگاه در مورد داده‌کاوی وجود دارد این است که داده‌کاوی به طریق خاصی انجام می‌شود و حتماً باید از روش‌های آماری و یادگیری ماشین در آن استفاده کنیم. اما زمانی که وارد یک پروژه واقعی می‌شوید، آنجاست که متوجه می‌شوید لزوماً این‌گونه نیست. بلکه داده‌کاوی یعنی در درجه‌ی اول یادگیرییم با داده‌ها راحت باشیم، با داده‌ها بازی کنیم و از حجم زیاد آن‌ها نترسیم... در بین دوستان وطنی که در این حوزه مشغول‌اند، کم‌تر دیده‌ام که به متن‌کاوی توجه کنند. این در حالی بود که در شرکت‌های بزرگ دنیا اولین چیزی که از شما خواسته می‌شود تا یاد بگیرید، متن‌کاوی و کار با عبارتهای باقاعده (Regex) است. برای یادگیری عبارتهای باقاعده مجبورید تا یک زبان اسکریپت‌نویسی بیاموزید. چرا که اولاً این داده‌ها فایل‌های متنی (مثلاً CSV) هستند و ثانیاً شما اغلب اجازه‌ی ذخیره آن‌ها را روی سیستم خود نداشته و باید به یک سرور لینوکسی متصل و از طریق ترمینال آن روی پروژه کار کنید. لذا قواعد متن‌کاوی و نیز دستورات یونیکس را باید به خوبی فراگیرید... شما برای آن که متخصص داده شوید، لازم نیست مهندس کامپیوتر بشوید. اما باید یک سری ابزارها را به طور هوشمندانه انتخاب کرده و یاد بگیرید و یکی از مهم‌ترین آن‌ها یادگیری یک زبان اسکریپت‌نویسی است.»

بدون هیچ داده‌ای، علم داده‌ی چندانی وجود نخواهد داشت. بنابراین گام نخست بدست آوردن مجموعه داده‌ها می‌باشد. داده‌ها از منابع مختلف و به شیوه‌های متنوعی بدست می‌آیند. برای مثال لازم است تا در آغاز هر پروژه یک یا چند گام زیر را انجام دهیم:

- دانلود داده روی یک سرور یا صفحه وب (معمولاً فایل‌های فشرده هستند)
- بازیابی داده از روی یک دیتابیس (برای مثال MySQL) و یا API (برای مثال تلگرام، توئیتر)
- استخراج داده از یک فایل صفحه گسترده (برای مثال CSV)
- خودمان داده‌ها را تولید کنیم، برای مثال خواندن اطلاعات سنسورها و یا استفاده از شیوه‌های سنتی مثل جمع‌آوری پرسش‌نامه.

گام بعدی پاکسازی داده‌ها می‌باشد. معمولاً داده‌ی اکتساب شده دارای مقادیر از دست رفته، خطاها، نویسه‌های عجیب و غریب، سطرهای خالی و یا حتی ستون‌های زاید می‌باشد. منظور از ستون‌های زاید، ستون‌هایی است که در تحلیل مورد نظر علاقه‌ای به آن‌ها نداریم. کارهایی که معمولاً در این قسمت انجام می‌گیرد شامل موارد زیر است:

- تبدیل فایل‌ها از یک فرمت (پسونند) به فرمت دیگر
- پالایش سطور و استخراج ستون‌های خاص
- جایگذاری مقادیر جدید و تصحیح مقادیر از دست‌رفته

در گام سوم سراغ اکتشاف بر روی داده‌ها می‌رویم. اینجا جایی است که برای اولین بار با داده‌ها و چیزی که بیان‌گر آن هستند سروکار پیدا می‌کنیم. شناخت داده، بدست آوردن آماره‌های آن و تصویرسازی داده‌های خام در این قسمت انجام می‌گیرد. در این مرحله با کمک روش‌های آماری و نمودارها به دنبال یافتن الگوها و مشخصه‌های قابل توجه در مجموعه داده هستیم. اغلب یک نمودار ساده از داده‌های خام بینش‌های مهم و بسیار مفیدی در اختیار قرار می‌دهد که ادامه مسیر پروژه تحلیل داده را به ما دیکته می‌کند. در این کتاب از روش‌های آماری ساده و از طریق ترسیم نمودارها و با کمک فیلترهای متنی یونیکس مرحله‌ی اکتشاف داده را انجام می‌دهیم. علاوه

بر این‌ها در این مرحله با استفاده از کدنویسی کیفیت توصیف مجموعه داده را بهبود داده و آن را برای مرحله‌ی مدل‌سازی آماده می‌کنیم.

در گام چهارم سراغ توصیف ریاضی داده‌ها و یا اینکه چه چیزی را پیش‌بینی می‌کنند، کمک مدل‌سازی آماری می‌رویم. در این قسمت روش‌های یادگیری ماشین چون خوشه‌بندی، دسته‌بندی، رگرسیون و کاهش ابعاد بروی داده‌ها اعمال می‌شود. علاوه بر انتخاب و برازش مدل، در این مرحله از معیارهای ریاضی برای سنجش اعتبار مدل انتخابی بهره می‌بریم. برای این منظور داده‌ها به گروه‌های داده آموزش، داده آزمایش و داده اعتبارسنجی متقابل^۱ تقسیم می‌شوند. هنگامی که مهارت کافی در یادگیری ماشین پیدا کنید، قادر خواهید بود تا الگوریتم‌های مناسب برای یک کاربرد خاص را براحتی شناسایی نمایید. در این کتاب به جزء گام چهارم، تمامی گام‌های دیگر در پروژه‌های مورد بررسی وجود دارند. دلیل این امر آن است که خط فرمان یونیکس برای پیاده‌سازی مدل‌های آماری و یادگیری ماشین از ابتدا مناسب نیست. لذا مرسوم است تا برای پیاده‌سازی الگوریتم‌های یادگیری ماشین در پروژه‌های علم داده از زبان‌هایی چون پایتون و R استفاده می‌شود.

گام آخر تفسیر داده‌ها و تحلیل نتایج است. این مرحله شامل نتیجه‌گیری، بررسی و ارزیابی نتایج و بحث و تبادل نظر با صاحب‌نظران و متخصصان شرکت یا سازمان می‌باشد تا مدیران بتوانند تصمیمات مناسبی اتخاذ کرده و موجب رونق و شکوفایی کسب و کارشان شوند. شما باید قادر باشید تا نتایج را به بهترین شکل که بیانگر واقعیت‌های موجود باشد، تصویرسازی نمایید. علاوه بر این شما باید مهارت داستان‌سرایی داده^۲ را نیز کسب نمایید. اهمیت این قسمت تا آنجا است که کتاب‌هایی با همین عنوان چاپ شده‌اند. بر اساس این واقعیت که اکثر مدیران دارای دانش آمار و علم داده نیستند، لازم است تا هر آنچه در دل داده‌ها نهفته است در قالب یک داستان جذاب و واقعی بیان شده تا توسط یک فرد تجاری قابل درک باشند.

^۱ cross-validation set

^۲ data storytelling

تا اینجا به تفصیل در مورد علم داده و اهمیت آن صحبت کردیم. در بخش پایانی این پیش‌گفتار به فصل‌بندی این کتاب می‌پردازیم. آنچه که ما در این کتاب روی آن تمرکز کرده‌ایم، پردازش و کار روی داده‌های متنی است. اما دلیل اهمیت زبان بَش در حوزه علم داده و یا به بیان دیگر دلیل کاربرد بَش در تحلیل داده‌های متنی را باید در تاریخچه‌ی سیستم یونیکس ریشه‌یابی کرد. در سال‌های اولیه‌ی توسعه‌ی سیستم یونیکس، علاقه‌ی شدیدی روی پردازش متن با کمک این سیستم وجود داشت؛ به طوری که اولین محصول تولید شده بر اساس سیستم یونیکس، یک ابزار پردازش متن برای استفاده در دیپارتمان ثبت اختراعات آزمایشگاه بل بود.

این کتاب در دو بخش تدوین شده است: ۱- **اسکرپت نویسی بَش (Bash)** و ۲- **تحلیل چالش‌های دنیای واقعی**. بخش اول کتاب به مبانی اسکرپت نویسی بَش، عبارات‌های با قاعده، و پردازش متن با کمک فیلترهای یونیکس می‌پردازد. این بخش دانش کافی برای مطالعه و درک این کتاب و پروژه‌هایی که در بخش دوم معرفی می‌شوند، در اختیار شما قرار می‌دهد. با این وجود خواننده‌ی علاقه‌مند برای مطالعه‌ی مباحث پیشرفته‌تر خط فرمان می‌تواند به مراجعی که در انتهای کتاب معرفی شده‌اند، مراجعه نماید. این بخش شامل ۸ فصل ابتدایی کتاب بوده که محتویات هر فصل در زیر خلاصه شده است. در انتهای هر فصل یک قسمت خودآزمایی همراه با جواب قرار داده شده است.

در **فصل اول** ابتدا تاریخچه‌ی یونیکس و بَش ارائه می‌شود و سپس با معرفی چند دستور ساده، شما را با محیط خط فرمان آشنا می‌کنیم. در انتهای این فصل از شما می‌خواهیم تا در مورد استفاده از دستورات کوچک و ساده خط فرمان برای ساختن برنامه‌های کاربردی فکر کنید. به طور خاص‌تر از شما خواسته می‌شود تا مسئله‌ی ساده و پرکاربرد تغییر نام گروهی هزاران فایل را با کمک ابزارهای موجود، حل نمایید.

فصل دوم در مورد اجراهای مختلف دستورات، نوع دستورات، فایل‌های راه‌انداز پوسته و نحوه‌ی صدور متغیر به آن‌ها، می‌باشد. فایل `profile` و `bashrc`. که به ترتیب محل نگهداری

متغیرهای محیطی و محل نگهداری توابع و جانشین‌ها هستند، به طور کامل بررسی می‌شوند. در ادامه فصل اسکرپت سنتی «سلام دنیا» را تمرین می‌کنیم و گام‌به‌گام با معرفی دستورات جدید بَش، ویژگی‌های جدیدی به آن اضافه می‌کنیم. بعلاوه در این فصل به متغیرهای محیطی در محیط ویندوز و نحوه ویرایش آن‌ها اشاره می‌شود.

در **فصل سوم** ابزارهای ضروری برای مدیریت یک سیستم یونیکس معرفی می‌شوند. در ابتدا مباحث فایل، پروسه و سیگنال بررسی می‌شوند. پروسه‌های پیش‌زمینه و پس‌زمینه، پروسه‌های والد و فرزند، و زیر-پوسته (زیر-شیل) از جمله مواردی هستند که در این قسمت توضیح داده شده‌اند. در قسمت انتهایی فصل به مبحث بسیار مهم مجوزهای دسترسی و حق مالکیت فایل‌ها و پوشه‌ها پرداخته می‌شود. تغییر مجوزها به دو روش نمادین و هشت‌هشتی آموزش داده شده است.

ورودی‌ها و خروجی‌ها در زبان بَش موضوع **فصل چهارم** کتاب است. خواندن ورودی صفحه کلید، چاپ بر روی ترمینال، بازهدایت به یک فایل / از یک فایل، و نیز بازهدایت خروجی یک دستور به ورودی دستور دیگر در این فصل به طور مفصل بحث شده‌اند. در ابتدا ابزارهای یونیکس برای مرور و مشاهده محتویات فایل‌های متنی معرفی می‌شوند. سپس برای خواندن ورودی دو روش مرسوم آموزش داده شده است. دستور read برای خواندن ورودی کاربر و ذخیره آن در یک متغیر استفاده می‌شود. عملگر «سند اینجایی» یکی دیگر از شیوه‌های خواندن اطلاعات در زبان بَش است. تمامی خطوط مشخص شده بین کلمه‌ی جداکننده یک سند اینجایی یکی پس از دیگری خوانده می‌شوند. برای چاپ روی خروجی استاندارد که به صورت پیش‌فرض ترمینال است، از دو دستور echo و printf بهره می‌بریم. جریان‌های ورودی، خروجی و خطای استاندارد، که به ترتیب با stdin، stdout و stderr مشخص می‌شوند، با کمک عملگرهای بازهدایت و واصف فایل‌ها ارجاع داده می‌شوند. واصف فایل و نحوه تغییر آن‌ها توسط دستور exec در ادامه بحث می‌شوند. در انتهای فصل نیز به لوله یا «پایپ» برای خواندن ورودی یک دستور به خروجی دستور، دیگر به همراه چندین مثال توضیح داده می‌شود.

فصل پنجم یکی از مهم‌ترین فصل‌های این کتاب است که به انواع «جانشینی» در بَش می‌پردازد. متغیرها و توابع نیز در این فصل بحث شده‌اند. در ابتدا سراغ متغیرها و نحوه تنظیم و حذف آن‌ها می‌رویم. در بَش برای تعریف یک متغیر نیازی به اعلان نوع آن نیست. هم‌چنین مقادیر متغیرها به صورت یک رشته ذخیره می‌شوند. در این قسمت انواع متغیرهای اسکالر، آرایه‌ها و اعداد صحیح بحث شده‌اند. بعلاوه به روش‌های مختلف بسط یک متغیر آرایه و نیز آرایه‌های شرکت‌پذیر می‌پردازیم. متغیرهای زبان بَش در سه دسته مختلف یعنی، متغیرهای محیطی، متغیرهای محلی و متغیرهای پوسته مقایسه می‌شوند. محیط به مجموعه‌ای از متغیرها و توابع پوسته اطلاق می‌شود. تابع و محدوده‌ی تعریف متغیرها موضوع قسمت بعدی است. در این قسمت پارامترهای مکانی استفاده شده در داخل بدنه توابع نیز بحث می‌شوند. قسمت پایانی این فصل به انواع مختلف جانشینی یا بسط در بَش اختصاص یافته است. اهمیت استفاده بجا از نقل قول‌های یگانه، دوگانه و نقل قول برگشتی در ابتدا توضیح داده می‌شود. در ادامه سراغ عملگرهای ریاضی و شیوه‌های مختلف ارزیابی عبارت‌های ریاضی می‌رویم. سپس در مورد جانشینی دستور، جانشینی نام‌فایل و جداسازی کلمات بحث می‌کنیم. مبحث مهم بسط پارامتر یا متغیر در انتهای فصل آورده شده است. این قسمت مکمل قسمت بسط آرایه‌ها می‌باشد. در این قسمت شکل‌های مختلف بسط یک پارامتر به تفصیل توضیح داده شده است. در پایان اشاره‌ای به متغیرهای ویژه زبان بَش می‌کنیم.

فصل ششم در مورد ساختارهای کنترل جریان می‌باشد. در واقع این ویژگی است که یک زبان برنامه‌نویسی را می‌سازد. در ابتدا وضعیت خروج و آزمایش بولی توضیح داده می‌شود. وضعیت خروج یا نتیجه‌ی یک آزمایش بولی در ساختارهای if، while و until و نیز عملگرهای NOT، AND و OR مورد استفاده قرار می‌گیرد. دستور test و یا معادل آن دستور [...] برای ارزیابی فایل‌ها، رشته و مقادیر عددی مورد استفاده قرار می‌گیرد. گروه‌های دوگانه [...] و پرازنترهای دوگانه (..) به ترتیب برای ارزیابی عبارت‌های باقاعده تعمیم‌یافته و عبارت‌های ریاضی کاربرد دارند. در قسمت دوم این فصل عبارت‌های شرطی بحث می‌شوند. نتایج بدست آمده از

آزمایش یک عبارت، برای تصمیم‌گیری در مورد ادامه‌ی جریان یک اسکریپت استفاده می‌شود. شکل‌های مختلف ساختار `if` و نیز عبارت `case` در این قسمت توضیح داده شده است. قسمت پایانی این فصل به بحث در مورد حلقه‌های `for`، `while` و `until` اختصاص دارد.

در **فصل هفتم** به سراغ فیلترهای متنی مقدماتی در `bash` و عبارت‌های باقاعده می‌رویم. عبارت‌های باقاعده یک حالت قوی‌تری از تطابق الگو نسبت به بسط نام فایل هستند و می‌توانند محدوده‌ی خیلی وسیع‌تری از الگوها را با دقت بیش‌تری بیان کنند. ابزارهای زیادی را در یونیکس می‌توان یافت که قدرت‌شان را از عبارت‌های باقاعده می‌گیرند، از جمله `grep`، `sed` و `awk`. دستور `grep` به همراه سایر ابزارهای پردازش متن مقدماتی در ابتدای این فصل بحث می‌شوند. قسمت دوم این فصل به عبارت‌های باقاعده اختصاص دارد. انواع عبارت‌های باقاعده از جمله عبارت باقاعده مقدماتی و تعمیم‌یافته با یکدیگر مقایسه و کاربردهای هر یک با ذکر مثال توضیح داده می‌شود. کلاس نویسه، لنگرها و ارجاعات برگشتی سایر مباحثی هستند که در این قسمت بررسی می‌شوند.

در **فصل هشتم** به دو ابزار پیشرفته یونیکس برای پردازش متن، یعنی `awk` و `sed` اختصاص دارد. ویرایشگر `sed` برخلاف ویرایشگرهای مرسوم با صفحه‌کلید می‌توانید با آن در تعامل باشید، یک ویرایش‌گر جریان است. به این معنی که قبل از پردازش یک فایل متنی، قوانین و دستورات لازم برای ویرایش آن را باید وارد کنیم. این ابزار بوسیله‌ی محدوده آدرس تعریف شده برای آن تعیین می‌کند روی کدام خطوط باید تغییرات صورت گیرد. دستوراتی که برای ویرایش یک جریان داده می‌توان استفاده کرد شامل جانشینی، حذف، چاپ و انتقال است. `awk` یک زبان برنامه‌نویسی کامل مبتنی بر فیلد بوده و از ساختارهای برنامه‌نویسی کامل‌تری نسبت دستورات ساده‌ی ویرایشگر `sed` برخوردار است. در این زبان برنامه‌نویسی شما می‌توانید با تعریف متغیرها، داده‌ها را ذخیره نمایید. از ساختارهای شرطی و حلقه‌ها کمک بگیرید. بر روی رشته کار کنید. عملیات ریاضی انجام دهید. و در نهایت گزارش‌های قالب‌بندی شده اختصاصی تهیه نمایید.

بخش اول کتاب همچون پلی است که فاصله‌ی بین علم و عمل را می‌پوشاند. بعد از یادگیری ابزارها و دستورات معرفی شده در این بخش، شما آماده‌ید تا آموخته‌هایتان را روی پروژه‌های واقعی بخش دوم کتاب اعمال نمایید. در فصول نهم تا پانزدهم، گام‌های علم داده را در چالش‌های واقعی بکار می‌بریم. فصل شانزدهم نیز شامل ۸ چالش دیگر خواهد بود که از حوزه‌های مختلف در صنعت، خدمات و تجارت انتخاب شده و از خواننده خواسته شده تا با تفکر و کار تحلیلی بر روی این پروژه‌ها دانش خود در این حوزه تثبیت نماید.

اگر می‌خواهید بر روی داده‌های تولید شده توسط برنامه‌های سفارش آنلاین غذا مثل اسنپ‌فود، چیلیوری، ریچون و یا ... کار کنید، **فصل نهم** برای شما نوشته شده است. در این فصل چون دسترسی به داده رستوران‌های ایرانی مقدور نبود، از یک مجموعه داده آزاد موجود در اینترنت برای انجام تحلیل بهره برده‌ایم. پروژه تعریف شده در این فصل به تحلیل عادت‌های غذایی مشتریان یک رستوران زنجیره‌ای در ایالت متحده آمریکا می‌پردازد. در حال حاضر نه تنها در تهران بلکه در اکثر شهرهای کشور، امکان سفارش غذا به صورت آنلاین از طریق برنامه‌های کاربردی تلفن‌های همراه و یا وبسایت‌های بی‌شماری که برای این منظور توسعه پیدا کرده‌اند، میسر شده است. صاحبان چنین کسب و کارهایی، اعم از توسعه‌دهندگان برنامه‌های کاربردی و نیز صاحبان رستوران‌ها، می‌توانند با کمک علم داده خدمات مطلوب‌تری (غذای با کیفیت، قیمت مناسب، کاهش زمان تحویل و بسته‌بندی سالم) را به مشتریان خود ارائه دهند. آن‌ها می‌توانند با در نظر گرفتن ذائقه و سلیقه‌ی مشتریان، ضمن بهبود رضایت و به تبع آن کسب وفادرای مشتریان، به موفقیت کسب و کار خود کمک کنند.

اگر قصد ادامه تحصیل در خارج از کشور را دارید و می‌خواهید بر اساس واقعیت‌های نهفته و مجهول در دل داده‌ها، بهترین دانشگاه‌ها را مشخص نمایید، **فصل دهم** برای شماست. نظام رتبه‌بندی دانشگاه‌ها در عصر حاضر بسیار متداول بوده و توسط موسسه‌های مختلفی در سراسر دنیا انجام می‌پذیرد. معمولاً شاخص‌های ارزیابی شامل مواردی چون جوایز معتبر دریافت شده توسط دانشجویان و اساتید، تعداد مقالات، تعداد ارجاع‌ها (استناد به مقالات)، کیفیت هیئت علمی، و

نسبت هیئت علمی به دانشجویان هستند. در این فصل بروی یک مجموعه داده آزاد که شامل رتبه‌بندی دانشگاه‌های کشور ایالات متحده آمریکا در سال ۲۰۱۷ میلادی است، کار می‌کنیم.

اگر به تحقیق بروی رفتار افراد در شبکه‌های اجتماعی و تاثیرگذاری سایرین و موضوعات مختلف بر نظر و عقاید آن‌ها علاقه دارید، مطمئناً **فصل یازدهم** برای شما جذاب خواهد بود. تلفیق علم داده، علوم اجتماعی و شبکه‌های پیچیده حوزه‌ی جدیدی با نام علوم اجتماعی محاسباتی یا واقعیت کاوی را بوجود آورده است، که البته موضوع این فصل نیست. حوزه دیگر مرتبط با شبکه‌های اجتماعی، بازاریابی و تبلیغات شبکه‌های اجتماعی است. این شبکه‌ها نقش مهمی در شکل‌دهی و آینده‌ی کسب و کارها بازی می‌کنند. لذا لازم است با تجزیه و تحلیل دقیق داده‌های صفحه‌ی شرکت یا سازمان در شبکه‌های اجتماعی، آینده کسب و کار را پیش‌بینی نماییم. در این فصل ابتدا ابزاری برای استخراج داده‌های موردنیاز از صفحات فیس‌بوک معرفی می‌شود و سپس بر اساس داده‌های تولیدشده، به بررسی پست‌های یک صفحه خاص می‌پردازیم.

اگر به کاربرد علم داده در حوزه‌های سیاسی و اجتماعی و به طور خاص مبحث مهم امنیت اجتماعی در یک کشور علاقه‌مند هستید، مطالعه‌ی **فصل دوازدهم** بسیار برای شما مفید و جالب خواهد بود. تا کنون مؤسسات مختلفی در سراسر دنیا، به مقایسه و امتیازبندی شهرهای جهان بر اساس معیارهایی چون امنیت، آزادی کسب و کار، آموزش، حمل و نقل، تامین اجتماعی، درآمد سرانه، سلامت و درمان، آب و هوا و اوقات فراغت پرداخته‌اند. یکی از معیارهای مهم در این ارزشیابی‌ها، معیار محیط سیاسی و اجتماعی می‌باشد که به ثبات سیاسی، ارتکاب جرم و جنایت و اجرای قانون اشاره دارد. در این فصل چون به داده‌های ملی در حوزه‌ی امنیت اجتماعی دسترسی نداشتیم، از یک مجموعه داده که مربوط به کشور استرالیا می‌باشد، استفاده کرده‌ایم. با کمک این مجموعه داده و با بررسی جرائم و جنایات انجام شده، نرخ ناامنی در شهرهای مختلف را پیدا می‌کنیم.

فصل سیزدهم برای اهالی فرهنگ و ادب نگاشته شده است و سعی دارد تا با رویکردی عملگرایانه و مبتنی بر کدنویسی، آثار منظوم و نمایشنامه‌های دوره‌ی رنسانس در اروپا را از زاویه‌ی متفاوت مورد نقد و واکاوی قرار دهد. به طور خاص، در این فصل به بررسی و تحلیل آثار ادبی دوره‌ی رنسانس در انگلستان می‌پردازیم و می‌خواهیم دریابیم که شاعران معروفی هم‌چون شکسپیر بیشتر از چه واژگانی در اشعار و نمایش‌نامه‌های خود بهره برده‌اند. مجموعه داده‌ای که برای انجام این پروژه مورد استفاده قرار گرفته است، نتیجه‌ی یک کار پژوهشی در دانشگاه نیوکاسل استرالیا می‌باشد.

اگر به بانکداری، بورس و بازارهای مالی علاقه دارید، فصل **چهاردهم** با یک مثال ساده قابلیت‌های بش در تولید و مدیریت صورت‌های مالی را به شما نشان می‌دهد. صورت‌های مالی از جمله معروف‌ترین داده‌هایی هستند که بعد از پایان هر عملیات تجاری و غیرتجاری تولید یا روزرسانی می‌شوند. در واقع داده‌های موجود در یک صورت حساب مالی مهم‌ترین و اساسی‌ترین اطلاعات را برای مدیریت و هدایت هر کسب و کاری فراهم می‌آورند. در حوزه بانکداری یکی از عملیاتی که هر روز و یا در هر ساعت انجام می‌گیرد، بالانس حساب جاری می‌باشد. در این فصل بر اساس یک صورت مالی که به طور دستی تولید شده است، میزان واریز، برداشت و چک‌های صادره از حساب جاری را مدیریت می‌کنیم. پس از پایان این فصل قادر خواهید بود تا برای صورت‌های مالی بزرگ و حجیمی که توسط سامانه‌های اینترنتی بانکی یا بورس تولید می‌شوند، برنامه‌هایی مشابه به زبان بش بنویسید.

فصل پانزدهم یک فصل کاملاً منحصر به فرد و ویژه بوده و علاقمندان خاص خود را در حوزه‌های هوش مصنوعی، علوم کامپیوتر و علوم و فنون قرآن کریم دارد. عنوان این فصل را قرآن کاوی یا قرآن پژوهی رایانه‌ای انتخاب کرده‌ایم تا به خوبی منعکس‌کننده‌ی آنچه در این اینجا به دنبال آن هستیم، باشد. به بیان دیگر، رویکردی نوین و مبتنی بر کدنویسی در حوزه علوم و فنون قرآن کریم معرفی می‌شود. در این فصل بجای استفاده از متن انگلیسی قرآن، روی متن عربی آن

متمرکز شده‌ایم تا بدین وسیله بینش و دانش دقیق و کاملی از تحلیل کلام الله مجید در اختیار ما قرار بگیرد.


آخرین فصل کتاب، یعنی **فصل شانزدهم**، به پروژه‌هایی برای تمرین بیشتر اختصاص داده شده است. در انتخاب پروژه‌ها سعی شده است تا چالش‌های متنوعی از حوزه‌های گوناگون انتخاب شود. تمرین‌های متعدد مطرح شده در این فصل، با این هدف انتخاب شده‌اند که خواننده را به شوق آورده و به چالش برانگیزد تا ضمن مطرح نمودن سوالات جالب و اساسی، به شرح و توصیف بیشتر مسائل پردازد، و در نهایت بتواند با تحلیل‌های آگاهانه و منطقی، ضمن افزایش بینش خود نسبت به موضوع مورد بررسی، آینده را پیش‌بینی نماید. این تمرین‌ها مثال‌های واقعی از کاربرد علم داده در صنعت، تجارت و تجارب زندگی روزانه می‌باشد.

کتاب «**علم داده کاربردی با کمک اسکریپت‌نویسی پش**» اولین کتاب فارسی می‌باشد که در مورد زبان پش نوشته شده است. با توجه به اهمیت سیستم‌های یونیکس و کاربرد روزافزون آن‌ها چه در صنعت و چه در دانشگاه، نگارش کتابی که به طور اختصاصی به زبان برنامه‌نویسی خط فرمان پردازد، لازم و ضروری بود. بنابراین مطالعه‌ی بخش اول این کتاب علاوه بر متخصصان علم داده برای سایر متخصصان و پژوهشگران در علوم مختلف از جمله هوش مصنوعی، رباتیک، فیزیک و ریاضی بسیار مفید خواهد بود. بعد از مطالعه‌ی مطالب این کتاب، خواننده باید پایه‌ی خوبی برای درک علم داده و کاربردهای آن یافته و قادر به انجام پروژه‌های واقعی و چالش‌های پیچیده‌تر با کمک اسکریپت‌نویسی پش و قواعد متن کاوی آن باشد. بعلاوه و مهم‌تر از همه، خواننده باید بتواند گام‌های اجرایی این کتاب را بدرستی پیاده‌سازی کرده تا بدین ترتیب، کشف بینش‌های نهفته در دل داده‌ها او را به واقعیت کاوی و یا پیش‌بینی آینده سوق دهد.

از شما خواننده محترم تقاضا دارم پس از مطالعه و استفاده از این کتاب نظرات و پیشنهادات خود در مورد این اثر را به نشانی ایمیل akhorshidi@live.com ارسال کرده و یا از طریق صفحه‌ی اختصاصی کتاب در سایت ناشر، کاربردهای این اثر در پروژه‌های کاری و شخصی خود

را با سایر خوانندگان به اشتراک بگذارید. گرچه که ما تمام سعی خود را به کار گرفته‌ایم تا اثری عاری از نقص را تقدیم حضورتان کنیم، اما اشتباهات هیچ وقت قابل اجتناب نبوده و نیستند. لذا خواهشمندیم اشتباهات ظاهری و باطنی اثر را برای ما مکاتبه نمایید. پیشاپیش از حسن نظر و دقت شما سپاسگزاریم. کدها، فایل‌ها، مجموعه داده‌ها، مثال‌ها، پاسخ تمرینات و حل پروژه‌های کتاب در مخزن اختصاصی زیر قرار داده می‌شود.

<https://github.com/akhorshidi/Bash4DSBook>

برای دریافت فایل‌های کتاب کافی است روی دکمه  کلیک نموده و گزینه Download Zip را انتخاب نمایید. همچنین می‌توانید از طریق وبسایت انتشارات و صفحه اختصاصی این کتاب به این محتویات دسترسی پیدا کنید.

در پایان از همسر عزیزم خانم فاطمه سادات کسائیان‌زاده مهابادی به خاطر سعی و بردباری‌های زیادی که در مدت زمان نگارش این کتاب از خود نشان دادند، و نیز زحمت حروف‌چینی، صفحه‌آرایی و ویرایش ادبی این اثر را متحمل شدند، صمیمانه سپاسگزاری می‌کنم. از کلیه برنامه‌نویسان و اساتید ارجمندی که با ارائه پیشنهادات و نقطه نظرات خود، سطح کیفی این اثر را ارتقاء داده‌اند، خالصانه تقدیر و تشکر می‌کنم. از همکاری کلیه کارکنان نشر ... بویژه مدیریت محترم آن جناب آقای ... کمال تشکر را دارم. سرانجام، این اثر را تقدیم به همه‌ی مردان و زنانی می‌کنم که برای اعتلای نام این سرزمین خالصانه تلاش می‌کنند. و این قول حتمی خداوند هست که می‌فرماید:

« وَأَنْ لَّيْسَ لِلْإِنْسَانِ إِلَّا مَا سَعَى » نجم/۳۹.