

کلدان‌ها

و اینترنت اشیاء

اصول و آبزارها

سعدهون عزیزی
یاسین امینی



کلان داده‌ها و اینترنت اشیاء

اصول و ابزارها

کلان داده‌ها و اینترنت اشیاء

اصول و ابزارها

مؤلفان:

سعدون عزیزی

عضو هیات علمی گروه مهندسی کامپیوتر دانشگاه کردستان

یاسین امینی

فارغ‌التحصیل مهندسی کامپیوتر دانشگاه کردستان



انتشارات دانشگاه کردستان

۱۳۹۹

عنوان و نام پدیدآور	عزیزی، سعدون	سرشناسه
مشخصات نشر	کلان داده‌ها و اینترنت اشیاء: اصول و ابزارها/ سعدون عزیزی، یاسین امینی.	
مشخصات ظاهری	سنندج، دانشگاه کردستان، انتشارات، ۱۳۹۹	
شابک	۳۱۷ص..، مصور، جدول، نمودار.	
وضعیت فهرست‌نویسی	۹۷۸-۶۲۲-۶۷۰۲-۲۴-۹	
یادداشت	فیپا.	
یادداشت	واژه نامه.	
یادداشت	کتاب نامه.	
موضوع	کلان داده‌ها، اینترنت اشیاء	
شناسه افزوده	امینی، یاسین.	
شناسه افزوده	دانشگاه کردستان، انتشارات.	
رده‌بندی کنگره	QA۷۶/۹:	
رده‌بندی دیوی	۰۰۵/۷:	
شماره کتابشناسی ملی	۶۱۸۴۰۰۷:	

نام کتاب: کلان داده‌ها و اینترنت اشیاء: اصول و ابزارها
 مولغین: سعدون عزیزی، یاسین امینی
 طراحی جلد: کوروش عنبری
 نوبت چاپ: اول
 تیراژ: ۵۰۰ جلد
 ناشر: انتشارات دانشگاه کردستان - سنندج، بلوار پاسداران، دانشگاه کردستان
 قیمت: ۴۵۰۰۰ تومان
 تلفن جهت سفارس کتاب: ۰۸۷۳۳۶۲۴۰۰۸ (خانم کاردوسیان)
 شابک: ۹۷۸-۶۲۲-۶۷۰۲-۲۴-۹
 حق چاپ برای انتشارات دانشگاه کردستان محفوظ است.

در حال حاضر داده به عنوان بالارزش ترین منبع جهان مطرح است و با ظهور اینترنت اشیاء اهمیت آن چند برابر شده است. مجموعه داده‌های نوین به اندازه‌ای پیچیده هستند که نمی‌توان آنها را توسط سخت‌افزارها و نرم‌افزارهای سنتی مدیریت کرد. حجم، سرعت و تنوع سه ویژگی اصلی داده‌های تولید شده‌ی امروزی هستند که باعث ایجاد مفهومی به نام کلان داده‌ها شده است. چنین ویژگی‌هایی روال‌های دریافت، ذخیره‌سازی، پردازش، تحلیل و بصری‌سازی کلان داده‌ها را به یک امر چالش‌برانگیز تبدیل کرده است. در دنیای رقابتی امروز، تحلیل کلان داده‌ها از اهمیت چشمگیری برخوردار است. اهمیت کلان داده‌ها به معنای مقدار داده در اختیار یک شرکت یا سازمان نیست بلکه به چگونگی استفاده از آن داده‌ها بستگی دارد. پردازش و تحلیل داده‌های جمع‌آوری شده به شرکت‌ها و سازمان‌ها کمک می‌کند تا بینش‌های لازم را کسب کرده و از آن در راستای تصمیم‌گیری‌های استراتژیک بهره بگیرند. طی چند سال گذشته، چارچوب‌ها و ابزارهای نوینی برای ذخیره‌سازی، پردازش و تحلیل کلان داده‌ها ارائه شده است که آشنایی و کار با آنها می‌تواند زمینه‌های پژوهشی و فرصت‌های شغلی متنوعی را پیش‌روی متخصصان این حوزه فراز دهد. با جستجو و مطالعه‌ی منابع مختلف، نویسنده‌گان بی‌بردنده که منبع جامعی برای آشنایی با اصول و ابزارهای کلان داده‌ها و اینترنت اشیاء وجود ندارد. به همین منظور، تأليف چنین کتابی را ضروری دانستند. دانشجویان و علاقمندان به رشته‌های مهندسی کامپیوتر، فناوری اطلاعات، علوم کامپیوتر و آمار در مقاطع مختلف تحصیلی می‌توانند از مطالب این کتاب استفاده کنند. علاوه بر این، شرکت‌ها و سازمان‌های فعال در حوزه کلان داده‌ها و اینترنت اشیاء می‌توانند مخاطب کتاب پیش‌رو باشند. در مورد نحوه‌ی نگارش این کتاب بایستی به این نکته اشاره شود که نویسنده‌گان با تحقیق و مطالعه‌ی منابع معتبر، که در انتهای کتاب ذکر شده‌اند، و همچنین دانش حاصل از تدریس درس کلان داده‌ها و اینترنت اشیاء این اثر را تولید کرده‌اند. امید است تأليف این کتاب بتواند به نوبه‌ی خود سهم برجسته‌ای در افزایش توان علمی جامعه‌ی دانشگاهی و حرفه‌ای کشور داشته باشد. مسلماً این اثر بدون اشکال نیست. لطفاً جهت کمک به بهبود و ارتقای بیشتر کیفیت مطالب، پیشنهادات و انتقادات خود را به آدرس iot@uok.ac.ir ارسال نمائید.

سعدون عزیزی

عضو هیات علمی گروه مهندسی کامپیوتر دانشگاه کردستان

فهرست مطالب

۱	۱. مقدمه‌ای بر کلان داده‌ها و اینترنت اشیاء
۳	۱-۱ مقدمه
۴	۲-۱ کلان داده‌ها
۵	۱-۲-۱ ویژگی‌های کلان داده‌ها
۱۰	۲-۲-۱ اهمیت کلان داده‌ها
۱۳	۳-۲-۱ حوزه‌های مرتبط با کلان داده‌ها
۱۸	۳-۱ منابع کلان داده‌ها
۲۰	۴-۱ چالش‌های کلان داده‌ها
۲۲	۵-۱ اینترنت اشیاء
۲۳	۱-۵-۱ تعریف اینترنت اشیاء
۲۵	۲-۵-۱ اهمیت اینترنت اشیاء
۳۰	۳-۵-۱ چالش‌های اینترنت اشیاء
۳۲	۴-۵-۱ معماری اینترنت اشیاء
۳۲	۱-۵-۴-۱ معماری سه لایه
۳۲	۱-۵-۴-۲ معماری چهار لایه
۳۳	۱-۵-۴-۳ معماری پنج لایه
۳۵	۶-۱ ترکیب اینترنت اشیاء با کلان داده‌ها
۳۶	۷-۱ جمع‌بندی و نتیجه‌گیری
۳۹	۲. پایگاه داده‌های غیررابطه‌ای
۴۱	۱-۲ مقدمه
۴۱	۲-۲ تاریخچه
۴۲	۳-۲ مدل داده BASE و ACID
۴۳	۴-۲ دسته‌بندی داده‌ها

کلان داده‌ها و اینترنت اشیاء: اصول و ابزارها

۴۵	CAP اصل ۵-۲
۴۶	۶-۲ اهمیت پایگاه داده‌های غیرابطه‌ای
۴۹	۷-۲ مقیاس‌پذیری عمودی و افقی
۵۰	۸-۲ ویژگی‌های پایگاه داده‌های غیرابطه‌ای
۵۳	۹-۲ چالش‌های پایگاه داده‌های غیرابطه‌ای
۵۵	۱۰-۲ دسته بندی سیستم‌های مدیریت داده غیرابطه‌ای
۵۶	۱۰-۲ پایگاه داده‌های کلید/مقدار
۵۷	۱۰-۲ پایگاه داده‌های سندگرا
۶۲	۱۰-۲ پایگاه داده‌های ستونی
۶۴	۱۰-۲ پایگاه داده‌های مبتنی بر گراف
۶۶	۱۱-۲ تفاوت پایگاه داده‌های رابطه‌ای و غیرابطه‌ای
۷۰	۱۲-۲ پایگاه داده‌های NewSQL
۷۲	۱۳-۲ نتیجه‌گیری
۷۵	۳. داده کاوی در کلان داده‌ها
۷۷	۱-۳ مقدمه
۷۸	۲-۳ داده کاوی
۷۹	۳-۳ مراحل داده کاوی
۷۹	۱-۳-۳ شناسایی منابع اطلاعاتی و تعریف مسئله
۸۰	۲-۳-۳ جمع‌آوری و پیش‌پردازش داده
۸۱	۳-۳-۳ استخراج الگو و مدل از داده‌ها
۸۲	۴-۳-۳ شناسایی موارد کلیدی استخراج شده
۸۲	۵-۳-۳ تفسیر و گزارش نتایج
۸۳	۴-۳ انواع روش‌های داده کاوی
۸۳	۵-۳ یادگیری ماشین
۸۴	۱-۵-۳ یادگیری با نظارت

فهرست

ج

۸۵	۲-۵-۳ یادگیری بدون نظارت
۸۶	۳-۵-۳ یادگیری نیمه نظارتی
۸۷	۴-۵-۳ یادگیری تقویتی
۸۸	۶-۳ دسته‌بندی
۸۹	۱-۶-۲ رگرسیون
۹۱	۲-۶-۳ درخت تصمیم
۹۵	۳-۶-۳ شبکه‌های عصبی
۹۹	۴-۶-۳ K-نزدیک‌ترین همسایه
۱۰۲	۷-۳ خوشبندی
۱۰۵	۱-۷-۲ کاربردهای خوشبندی
۱۰۶	۲-۷-۳ الگوریتم‌های خوشبندی
۱۰۶	۱-۲-۷-۳ الگوریتم K-میانگین
۱۰۹	۲-۲-۷-۳ الگوریتم CLARA
۱۱۰	۸-۳ استخراج قواعد انجمنی
۱۱۱	۱-۸-۳ پشتیبان و اطمینان
۱۱۳	۲-۸-۳ الگوریتم‌های استخراج قواعد انجمنی
۱۱۴	۳-۸-۲-۱ AIS الگوریتم
۱۱۶	۳-۸-۳ الگوریتم پیشینار یا Apriori
۱۱۷	۴-۸-۳ الگوریتم AprioriTid
۱۱۷	۵-۸-۲ الگوریتم Apriorihybrid
۱۱۷	۹-۳ ابزارهای داده کاوی
۱۱۸	۱-۹-۲ زبان برنامه‌نویسی پایتون
۱۱۹	۲-۹-۲ زبان برنامه‌نویسی R
۱۲۰	۳-۹-۲ محیط برنامه‌نویسی متلب
۱۲۱	۴-۹-۲ نرم‌افزار ریپیدماینر
۱۲۲	۵-۹-۲ نرم‌افزار کلمنتاین
۱۲۳	۶-۹-۲ نرم‌افزار ارنج

۱۲۴	۷-۹-۳ نرم‌افزار و کا
۱۲۵	۱۰-۳ جمع‌بندی و نتیجه‌گیری
۱۲۷	۴. اکوسیستم هدوپ.
۱۲۹	۱-۴ مقدمه
۱۳۰	۲-۴ داستانی از هدوپ
۱۳۲	۳-۴ تاریخچه هدوپ
۱۳۳	۴-۴ هدوپ چیست؟
۱۳۵	۵-۴ هدوپ چگونه کار می‌کند؟
۱۳۶	۶-۴ معماری هدوپ
۱۳۷	۷-۴ ویژگی‌ها و مزایای هدوپ
۱۴۱	۸-۴ اکوسیستم هدوپ
۱۴۲	۱-۸-۴ سیستم فایل توزیع شده هدوپ (HDFS)
۱۴۷	۲-۸-۴ زمانبند هدوپ (Yarn)
۱۴۹	۳-۸-۴ مدل پردازشی نگاشت/کاهش
۱۵۲	۴-۸-۴ HBASE
۱۵۳	۵-۸-۴ HIVE
۱۵۴	۶-۸-۴ آپاچی پیگ (PIG)
۱۵۶	۷-۸-۴ تر
۱۵۷	۸-۸-۴ اسکوپ (Sqoop)
۱۵۸	۹-۸-۴ آپاچی فلوم (Flume)
۱۵۹	۱۰-۸-۴ دریل (Drill)
۱۶۰	۱۱-۸-۴ آپاچی آمباری (Ambari)
۱۶۰	۱۲-۸-۴ زوکپر (ZooKeeper)
۱۶۱	۱۳-۸-۴ اووزی (Oozie)
۱۶۱	۱۴-۸-۴ سامزا (Samza)
۱۶۳	۱۵-۸-۴ ماہوت (Mahout)

۱۶۳	۹-۴ توسعه دهنده‌گان هدوپ
۱۶۴	۱۰-۴ معاييـ هدوپ
۱۶۵	۱۱-۴ نصب هدوپ
۱۶۸	۱۲-۴ شروع کار با هدوپ (شمارش کلمات)
۱۷۰	۱۳-۴ نتیجه‌گيري و جمع‌بندی
۱۷۱	۵. پردازش درون حافظه‌اي با استفاده از گريـدـگـين و اـسـپـارـكـ
۱۷۳	۱-۵ مقدمه
۱۷۳	۲-۵ گـريـدـگـين
۱۷۴	۱-۲-۵ موـتوـرـ پـرـداـزـشـي
۱۷۵	۲-۲-۵ فـاـيـلـ سـيـسـتـم
۱۷۷	۳-۲-۵ تـرـكـيـبـ هـدوـپـ وـ گـريـدـگـين
۱۷۸	۳-۵ آـپـاـچـيـ اـسـپـارـكـ
۱۷۹	۱-۳-۵ مـزـيـتـهـايـ اـسـپـارـكـ
۱۸۵	۲-۳-۵ اـسـپـارـكـ درـ دـنـيـايـ وـاقـعـيـ
۱۸۶	۳-۳-۵ بـرـرسـيـ اـجـزـايـ اـسـپـارـكـ
۱۹۰	۴-۳-۵ چـالـشـهـايـ اـسـپـارـكـ
۱۹۱	۵-۳-۵ نـصـبـ اـسـپـارـكـ
۱۹۳	۴-۵ شـروعـ کـارـ باـ اـسـپـارـكـ وـ زـيـانـ بـرـنـامـهـ نـوـيـسـيـ اـسـكـالـاـ
۲۰۲	۵-۵ جـمـعـبـندـيـ وـ نـتـيـجـهـ گـيرـي
۲۰۵	۶. پـرـداـزـشـ جـرـيـانـيـ وـ اـبـزارـهـايـ آـنـ
۲۰۷	۱-۶ مـقـدـمه~
۲۰۷	۲-۶ پـرـداـزـشـ جـرـيـانـي~
۲۱۲	۳-۶ آـپـاـچـيـ استورـم
۲۱۳	۱-۳-۶ مـعـارـىـ آـپـاـچـيـ استورـم
۲۱۶	۲-۳-۶ مـزـيـتـهـايـ آـپـاـچـيـ استورـم

کلان داده‌ها و اینترنت اشیاء: اصول و ابزارها

۲۱۸	۳-۳-۶ نصب آپاچی استورم
۲۲۳	۴-۶ آپاچی فلینک
۲۲۵	۱-۴-۶ معماری آپاچی فلینک
۲۲۷	۲-۴-۶ اجزای آپاچی فلینک
۲۲۸	۳-۴-۶ مراحل اجرای برنامه‌ها در فلینک
۲۲۹	۴-۴-۶ مزایا و ویژگی‌ها
۲۳۱	۵-۴-۶ نصب و پیکربندی
۲۳۳	۵-۶ آپاچی نایفای
۲۳۵	۱-۵-۶ نحوه کار آپاچی نایفای
۲۳۵	۲-۵-۶ ویژگی‌های آپاچی نایفای
۲۳۷	۳-۵-۶ کاربردهای آپاچی نایفای
۲۳۷	۴-۵-۶ برنامه‌نویسی جریانی
۲۳۸	۵-۵-۶ نصب آپاچی نایفای
۲۳۹	۶-۶ نتیجه‌گیری و جمع‌بندی
۲۴۱	۷. سیستم پیامرسانی کافکا
۲۴۳	۱-۷ مقدمه
۲۴۳	۲-۷ سیستم‌های پیامرسانی
۲۴۸	۳-۷ هدف کافکا
۲۵۰	۴-۷ اجزای آپاچی کافکا
۲۵۴	۵-۷ مزایای آپاچی کافکا
۲۵۶	۶-۷ چالش‌های آپاچی کافکا
۲۵۷	۷-۷ انواع سیستم‌های مبتنی بر آپاچی کافکا
۲۵۹	۸-۷ ابزارهای کافکا
۲۶۰	۹-۷ نحوه کار آپاچی کافکا
۲۶۱	۱۰-۷ نصب و پیکربندی آپاچی کافکا

۱۱-۷ شروع کار با کافکا.....	۲۶۳
۱۲-۷ نتیجه‌گیری و جمع بندی.....	۲۶۹
۸. ابزارهای پردازشی نوین.....	۲۷۱
۱-۸ مقدمه.....	۲۷۳
۲-۸ آپاچی بیم.....	۲۷۳
۱-۲-۸ مزایا و چالش‌ها.....	۲۷۵
۲-۲-۸ نصب آپاچی بیم.....	۲۷۶
۳-۸ آپاچی آپکس.....	۲۷۸
۱-۳-۸ معماری آپاچی آپکس.....	۲۷۹
۲-۳-۸ مدل برنامه‌نویسی.....	۲۸۰
۴-۸ آپاچی ایگنایت.....	۲۸۱
۱-۴-۸ نحوه کار کرد آپاچی ایگنایت.....	۲۸۲
۲-۴-۸ معماری آپاچی ایگنایت.....	۲۸۴
۳-۴-۸ مزیت‌های استفاده از آپاچی ایگنایت.....	۲۸۵
۵-۸ نتیجه‌گیری و جمع بندی.....	۲۸۹
واژه‌نامه.....	۲۹۱
اختصارات.....	۲۹۶
منابع و مراجع.....	۲۹۷

فهرست شکل‌ها

..... ۶	شکل ۱-۱: ویژگی‌های اصلی کلان داده‌ها
..... ۱۱ شکل ۱-۲: در هر دقیقه چه اتفاقی در اینترنت می‌افتد؟
..... ۱۷ شکل ۱-۳: ارتباط یادگیری ماشین و کلان داده
..... ۲۶ شکل ۱-۴: حوزه‌های کاربردی اینترنت اشیاء
..... ۲۸ شکل ۱-۵: سوسمک هوشمند برای نجات زندگانی زیر آوار
..... ۳۰ شکل ۱-۶: درصد درآمد هر یک از حوزه‌های اینترنت اشیاء تا سال ۲۰۲۵
..... ۳۴ شکل ۱-۷: لایه‌های مختلف معماری پنج لایه
..... ۴۶ شکل ۲-۱: نمونه‌هایی از سیستم مدیریت داده که در آنها دو مورد از ویژگی‌های اصل CAP رعایت شده است
..... ۴۸ شکل ۲-۲: رشد داده‌های غیرساختاریافته
..... ۵۰ شکل ۲-۳: مقیاس‌پذیری عمودی (up-scale) در مقابل مقیاس‌پذیری افقی (out-scale)
..... ۶۰ شکل ۲-۴: رابطه بین جدول‌ها در پایگاه داده‌های رابطه‌ای
..... ۶۵ شکل ۲-۵: پایگاه داده مبتنی بر گراف با ارتباط دوستی یا علاقه‌مندی‌ها
..... ۶۶ شکل ۲-۶: روند روبه رشد استفاده از پایگاه داده‌های غیررابطه‌ای به ویژه Neo4j
..... ۸۳ شکل ۳-۱: انواع روش‌های داده کاوی
..... ۸۵ شکل ۳-۲: یادگیری با نظارت
..... ۸۶ شکل ۳-۳: یادگیری بدون نظارت
..... ۸۷ شکل ۳-۴: یادگیری تقویتی (یادگیری مبتنی بر پاداش و تنبیه)
..... ۹۱ شکل ۳-۵: دسته‌بندی داده‌ها با استفاده از الگوریتم رگرسیون
..... ۹۳ شکل ۳-۶: نمونه‌ای از یک درخت تصمیم برای تشخیص بیماری مراجعه کننده
..... ۹۵ شکل ۳-۷: شبکه‌های عصبی ساده که چند ورودی داشته و فقط یک لایه میانی دارد
..... ۹۷ شکل ۳-۸: رسم خط به کمک پرسپترون برای جداسازی سگ از گربه
..... ۹۸ شکل ۳-۹: پرسپترون چند لایه‌ای که دارای چندین لایه مخفی یا میانی است
..... ۱۰۳ شکل ۳-۱۰: خوشبندی و تفکیک مشتریان برای سودهای بیشتر

..... ۱۰۸	شکل ۳-۱۱: نمودار داده‌های جدول ۳-۵
..... ۱۱۳	شکل ۳-۱۲: سبد خریدهای مختلف یک فروشگاه
..... ۱۱۸	شکل ۳-۱۳: محیط ژوپیتر برای برنامه‌نویسی پایتون
..... ۱۱۹	شکل ۳-۱۴: محیط RStudio برای زبان برنامه‌نویسی R
..... ۱۲۰	شکل ۳-۱۵: محیط برنامه‌نویسی متلب
..... ۱۲۲	شکل ۳-۱۶: محیط نرم‌افزار ریبدماینر
..... ۱۲۳	شکل ۳-۱۷: محیط نرم‌افزار spss clementine
..... ۱۲۴	شکل ۳-۱۸: محیط نرم‌افزار ارنج
..... ۱۲۵	شکل ۳-۱۹: محیط نرم‌افزار وکا
..... ۱۳۴	شکل ۴-۱: داگ کاتینگ، خالق هدوپ
..... ۱۳۷	شکل ۴-۲: معماری هدوپ
..... ۱۴۲	شکل ۴-۳: اجزای تشکیل دهنده اکوسیستم هدوپ (باغ و حش هدوپ!)
..... ۱۴۵	شکل ۴-۴: تقسیم‌بندی داده در فایل سیستم هدوپ
..... ۱۴۶	شکل ۴-۵: نحوه انجام عمل تکرار در فایل سیستم هدوپ
..... ۱۴۶	شکل ۴-۶: نحوه تکرار داده‌ها در مرکز داده هدوپ
..... ۱۴۹	شکل ۴-۷: نحوه کار کرد آپاچی یارن (زمان بند هدوپ)
..... ۱۵۰	شکل ۴-۸: نحوه کار کرد موتور پردازشی نگاشت/کاهش، برای شمارش کلمات
..... ۱۵۳	شکل ۴-۹: دسترسي ترتیبی در مقابل دسترسي تصادفی
..... ۱۵۵	شکل ۴-۱۰: معماری آپاچی پیگ
..... ۱۵۶	شکل ۴-۱۱: معماری آپاچی تز
..... ۱۵۸	شکل ۴-۱۲: معماری آپاچی اسکوپ
..... ۱۵۹	شکل ۴-۱۳: معماری آپاچی فلوم
..... ۱۶۲	شکل ۴-۱۴: معماری و نحوه کار کرد آپاچی سامزا
..... ۱۶۴	شکل ۴-۱۵: ذخیره نتایج میانی در فایل سیستم هدوپ و کاهش سرعت پردازشی سیستم
..... ۱۷۵	شکل ۵-۱: نحوه کار کرد گریدگین
..... ۱۷۷	شکل ۵-۲: نحوه کار کرد سیستم ذخیره‌سازی گریدگین

کلان داده‌ها و اینترنت اشیاء: اصول و ابزارها

۱۷۷	شکل ۵-۳: ترکیب اکوسیستم هدوب و گرید گین
۱۸۱	شکل ۵-۴: عملیات‌های انجام شده توسط چارچوب اسپارک
۱۸۲	شکل ۵-۵: تفاوت بین ذخیره نتایج میانی در هدوب و اسپارک
۱۸۳	شکل ۵-۶: معماری آپاچی اسپارک
۱۸۶	شکل ۵-۷: معماری رایانش مه
۱۸۷	شکل ۵-۸: اجزای تشکیل‌دهنده اسپارک
۱۸۸	شکل ۵-۹: پردازش داده‌های جریانی در اسپارک با روش میکرودسته‌ای
۱۹۲	شکل ۱۰-۵: درست نصب شدن اسپارک
۲۰۹	شکل ۱۶-۱: پنجره ثابت در پردازش‌های جریانی
۲۱۰	شکل ۱۶-۲: پنجره اسلامیدی در پردازش داده‌های جریانی
۲۱۱	شکل ۱۶-۳: تفاوت بین پنجره‌های داده
۲۱۴	شکل ۱۶-۴: معماری آپاچی استورم
۲۱۵	شکل ۱۶-۵: توپولوژی آپاچی استورم
۲۱۶	شکل ۱۶-۶: نحوه کارکرد و ارتباط بخش‌های مختلف آپاچی استورم
۲۲۵	شکل ۱۶-۷: لوگوی آپاچی فلینک
۲۲۶	شکل ۱۶-۸: معماری آپاچی فلینک
۲۲۹	شکل ۱۶-۹: مراحل اجرای برنامه‌ها در آپاچی فلینک
۲۳۲	شکل ۱۰-۶: داشبورد تحت وب آپاچی فلینک
۲۳۳	شکل ۱۱-۶: نمایش اجرای برنامه‌ها در داشبورد مدیریتی آپاچی فلینک
۲۳۴	شکل ۱۲-۶: محیط توسعه آپاچی نایفای
۲۳۶	شکل ۱۳-۶: محیط یکپارچه پردازشی برای پردازش کلان داده‌ها
۲۳۹	شکل ۱۴-۶: نصب موفق آپاچی نایفای
۲۴۴	شکل ۱۷-۱: سیستم پیامرسانی نقطه به نقطه
۲۴۵	شکل ۱۷-۲: سیستم پیامرسانی انتشار/اشتراك
۲۴۸	شکل ۱۷-۳: نحوه عملکرد آپاچی کافکا در دنیای واقعی

..... ۲۴۹	شکل ۷-۴: شبکه استفاده آپاچی کافکا
..... ۲۴۹	شکل ۷-۵: ارتباط مستقیم شبکه پردازشی
..... ۲۵۰	شکل ۷-۶: ارتباط غیرمستقیم شبکه پردازشی با استفاده از آپاچی کافکا
..... ۲۵۰	شکل ۷-۷: اجزای تشکیل دهنده آپاچی کافکا
..... ۲۵۸	شکل ۷-۸: کافکا در نقش سیستم انتشار/اشتراك
..... ۲۵۸	شکل ۷-۹: کافکا در نقش سیستم صفحه
..... ۲۶۰	شکل ۷-۱۰: نحوه کار کرد کافکا استریمز
..... ۲۶۱	شکل ۷-۱۱: نحوه عملکرد آپاچی کافکا
..... ۲۷۴	شکل ۸-۱: آپاچی بیم و نحوه کار کرد آن
..... ۲۷۵	شکل ۸-۲: آپاچی بیم به عنوان یک پارچه ساز ابزارهای کلان داده
..... ۲۸۰	شکل ۸-۳: معماری آپاچی آپکس
..... ۲۸۱	شکل ۸-۴: نحوه کار کرد و اجرای برنامه در آپاچی آپکس
..... ۲۸۴	شکل ۸-۵: یک پارچه سازی ابزارهای ذخیره سازی و زبان های برنامه نویسی در آپاچی ایگنایت
..... ۲۸۵	شکل ۸-۶: نحوه ذخیره داده ها در آپاچی ایگنایت
..... ۲۸۶	شکل ۸-۷: معماری پردازش جریانی در آپاچی ایگنایت
..... ۲۸۸	شکل ۸-۸: تغییرات قیمت حافظه اصلی و حافظه جانبی با گذر زمان
..... ۱۹۳	شکل ۸-۹: دایرکتوری های لازم برای راه اندازی SBT
..... ۱۹۴	شکل ۸-۱۰: نحوه ایجاد دایرکتوری های پروژه اسکالا و اسپارک
..... ۱۹۶	شکل ۸-۱۱: موفقیت آمیز بودن ایجاد پروژه
..... ۱۹۶	شکل ۸-۱۲: آدرس فرار گیری فایل اجرایی
..... ۱۹۷	شکل ۸-۱۳: خروجی اجرای پروژه

فهرست جدول‌ها

جدول ۱-۱: تفاوت بین داده‌های سنتی و کلان داده‌ها	۱۳
جدول ۲-۱: ذخیره اطلاعات در جدول پایگاه داده رابطه‌ای	۵۸
جدول ۲-۲: اضافه کردن اطلاعات جدید به پایگاه داده‌های رابطه‌ای	۵۹
جدول ۲-۳: ذخیره اطلاعات در جدول پایگاه داده رابطه‌ای	۶۲
جدول ۲-۴: ذخیره اطلاعات در جدول پایگاه داده غیررابطه‌ای ستون‌گرا	۶۲
جدول ۲-۵: تفاوت بین پایگاه داده‌های رابطه‌ای و غیررابطه‌ای	۷۰
جدول ۳-۱: مجموعه داده‌های آموزشی	۸۹
جدول ۳-۲: مجموعه داده‌های آموزش	۱۰۱
جدول ۳-۳: فاصله ورودی جدید از تک نمونه‌ها	۱۰۱
جدول ۳-۴: تعیین برچسب برای نمونه جدید بر حسب رأی اکثریت	۱۰۲
جدول ۳-۵: مجموعه داده جهت خوشه‌بندی با روش k-میانگین	۱۰۸
جدول ۳-۶: فاصله داده‌ها از مرکز خوشه	۱۰۹
جدول ۳-۷: سبد خرید برای استخراج قواعد انجمنی	۱۱۲
جدول ۳-۸-داده‌های یک سبد خرید	۱۱۳
جدول ۳-۹: پایگاه داده اصلی	۱۱۴
جدول ۳-۱۰: الگوریتم AIS و اولین پیمایش پایگاه داده برای استخراج محصولات و تعداد آن‌ها	۱۱۵
جدول ۳-۱۱: استخراج قواعد دوتایی با شناسه‌ی تراکنش	۱۱۵
جدول ۳-۱۲: استخراج قواعد سه‌تایی با شناسه‌ی تراکنش	۱۱۶
جدول ۷-۱: مقایسه سیستم‌های پیامرسان نقطه به نقطه و انتشار/اشتراک	۲۴۵

فصل اول

مقدمه ای بر کلان داده ها و اینترنت اشیاء

مطلوبی که در این فصل پوشش داده می‌شوند:

- ✓ کلان داده چیست و چه ویژگی‌هایی دارد؟
- ✓ چرا مطالعه کلان داده‌ها اهمیت دارد؟
- ✓ اینترنت اشیاء چیست؟
- ✓ چرا مطالعه اینترنت اشیاء اهمیت دارد؟
- ✓ کلان داده‌ها و اینترنت اشیاء چه حوزه‌ها و فناوری‌هایی را تحت تأثیر قرار می‌دهند؟
- ✓ رابطه کلان داده‌ها و اینترنت اشیاء چیست؟

۱-۱ مقدمه

در دهه‌ی اخیر داده به عنوان بالارزش‌ترین منبع جهان مطرح شده است و این ارزش و اهمیت داده روز به روز افزایش می‌یابد تا جایی که در سال ۲۰۱۵ برای اولین بار ارزش داده از نفت پیشی گرفت!

طبق آمارهای منتشر شده، ۹۰ درصد داده‌های ذخیره شده در دنیا در دو سال اخیر تولید شده‌اند! روزانه ۲.۵ کوین‌تليون^۱ داده در دنیا تولید می‌شود. در سال ۲۰۱۱ یک رکورد جدید ثبت گردید و در هر دو روز به اندازه کل داده‌های تولید شده تا سال ۲۰۰۳، داده تولید شد! این رکورد در سال ۲۰۱۵ به دو ساعت رسید! در سال‌های اخیر حجم داده‌های تولید شده هر دو روز دو برابر می‌شود! در سال ۲۰۰۸ تعداد دستگاه‌های متصل به اینترنت از تعداد آدمهای روی کره زمین پیشی گرفت و به یک نقطه عطف در هوشمندی دنیا رسید! شرکت‌های بزرگ دنیا هر روز آماری از تعداد دستگاه‌های هوشمند که تا سال ۲۰۲۰ قرار است وجود داشته باشند، ارائه می‌دهند. بر اساس پیش‌بینی شرکت سیسکو، تعداد دستگاه‌های متصل به اینترنت تا سال ۲۰۲۰ به بیش از ۲۶ میلیارد خواهد رسید!

همه عناوین بالا، عناوین اصلی اخباری هستند که ما هر روز در گوش و کنار مجله‌ها و سایت‌های معتبر مرتبط با علوم کامپیوتر مشاهده می‌کنیم. همین امر ایجاب می‌کند که ما در مورد این فناوری‌ها به اندازه کافی مطالعه داشته باشیم تا بتوانیم دانش خود را به روز نگه داشته و خود را با قطار سریع السیر دنیای فناوری، که هر

روز با سرعت بیشتری پیش می‌رود، هم جهت نموده و از این فناوری‌ها در امور کشوری، مدیریت بحران‌ها، ساده‌سازی زندگی افراد جامعه و غیره استفاده نماییم. در این فصل نگاهی به فناوری‌های دنیای امروز، اهمیت آن‌ها و چرایی وجود آن‌ها خواهیم داشت.

۲-۱ کلان داده‌ها

مفهوم کلان داده‌ها^۱ برای اولین بار در دهه اول قرن جاری به علت پیشی گرفتن سرعت و نرخ تولید داده از تحلیل و فضای ذخیره سازی موجود، به طور جدی مطرح گردید. دقیقاً زمانی که اولین جرقه‌های مفهومی فراتر از داده مطرح گردید، سازمان‌ها و شرکت‌های بزرگ با ذخیره‌سازی و تحلیل داده به عنوان چالشی اساسی مواجه شدند.

به دلیل نوپا بودن و وجود سازمان‌های مختلف با حوزه‌های کاری مختلف، تاکنون برای کلان داده‌ها استاندارد واحدی تعریف نشده است. در ادامه، تعریف‌های متعدد و رایجی که برای این مفهوم وجود دارد را ارائه می‌دهیم.

موسسه گارتنر

کلان داده‌ها به داده‌هایی گفته می‌شوند که دارای حجم^۲ بالا، تنوع^۳ بالا و نرخ تولید^۴ بالا باشند.

شرکت اوراکل

کلان داده‌ها به داده‌هایی گفته می‌شوند که دارای چهار ویژگی حجم بالا، تنوع بالا، نرخ تولید بالا و همچنین صحت و درستی باشند.

مک‌کینزی^۱

مجموعه داده‌هایی که اندازه آنها از قابلیت ابزارهای پایگاه داده سنتی برای دریافت، ذخیره‌سازی، مدیریت و تحلیل خارج است.

موسسه ملی استاندارد و فناوری

مجموعه داده‌هایی هستند که در آن‌ها حجم داده‌ها، سرعت تولید و نمایش داده‌ها، توانایی ما را برای تحلیل کارآمد و موثر با استفاده از روش‌های سنتی محدود می‌سازد و برای تجزیه و تحلیل و پردازش کارآمد این داده‌ها، استفاده از روش‌های نوین و مقیاس‌پذیر ضروری می‌باشد.

تعريف ساده^۲

کلان داده‌ها به داده‌هایی اطلاق می‌گردد که نتوان آن‌ها را توسط یک سیستم عادی ذخیره‌سازی و پردازش کرد. بدین معنی که یک سیستم عادی توان پردازشی و ذخیره‌سازی لازم برای تجزیه و تحلیل این داده‌ها را نداشته باشد.

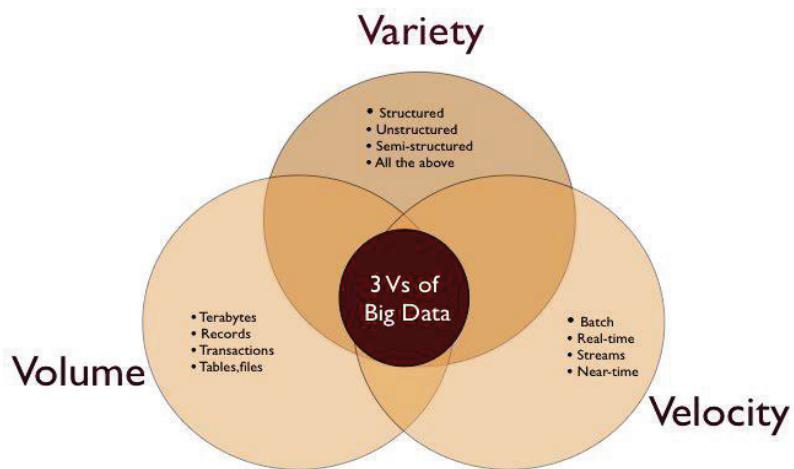
۱-۲-۱- ویژگی‌های کلان داده‌ها

حال این سوال پیش می‌آید که کلان داده‌ها چه ویژگی‌هایی دارند و ما چه مجموعه از داده‌ها را کلان داده می‌نامیم؟ آیا تمام داده‌های دنیای امروز کلان داده محسوب می‌شوند؟ باید بگوییم خیر؛ اینگونه نیست و دانشمندان داده، مجموعه‌ای از ویژگی‌ها را برای کلان داده‌ها در نظر گرفته‌اند که اگر مجموعه داده ما چنین ویژگی‌هایی را دارا باشد می‌توان آن را کلان داده در نظر گرفت. برای کلان داده سه ویژگی اصلی شامل حجم، تنوع و نرخ تولید در نظر گرفته می‌شود [۱]. در دنیای کلان داده‌ها، این سه ویژگی به یک استاندارد

McKinsey^۱

^۲ این تعریف توسط نویسنده‌گان برداشت شده است.

تبدیل شده است (شکل ۱-۱). ویژگی‌های دیگری همچون صحت^۱، اعتبار^۲، بصری‌سازی^۳، نوسان^۴، ارزش^۵ و آسیب‌پذیری^۶ هم وجود دارند [۲] که جنبه کاربردی آن‌ها بیشتر مربوط به سازمان‌های خاص می‌باشد و کمتر مورد توجه و اهمیت عام قرار می‌گیرند. در ادامه، به توصیف هر یک از این ویژگی‌ها می‌پردازیم.



شکل ۱-۱ ویژگی‌های اصلی کلان داده‌ها

حجم

حجم یکی از ویژگی‌های بسیار مهم کلان داده‌هاست؛ به ویژه زمانی که گفته می‌شود ۹۰ درصد داده‌های دنیا در دو سال اخیر تولید شده‌اند. اما حجم در کلان داده مفهوم مبهمی محسوب می‌شود و در واقع نمی‌توان چندان به ویژگی حجم برای تشخیص کلان داده بودن یک مجموعه اکتفا نمود. در واقع ما نمی‌دانیم که دقیقاً چه حجمی را باید به عنوان حجم پایه برای کلان داده در نظر بگیریم. به عنوان مثال اگر یک تراباتیت داده را حجم پایه‌ای انتخاب نماییم، آیا ۹۹۰ گیگابایت داده، کلان داده محسوب نمی‌شود؟ اگرچه گفته می‌شود که

Veracity ^۱
Validity ^۲
Visualization ^۳
Volatility ^۴
Value ^۵
Vulnerability ^۶

حجم کلان داده‌ها در دنیای امروز در حدود ترابایت و پتابایت می‌باشد ولی ما در مورد ویژگی حجم تنها به این نکته بسته می‌کنیم: مجموعه داده‌ای کلان داده است که حجم آن توسط پایگاه داده‌های سنتی قابل مدیریت نباشد. حجم به دو عامل سرعت رشد داده و سرعت توسعه دیسک‌های ذخیره‌سازی بستگی دارد.

تنوع

تنوع باعث می‌شود که کلان داده‌ها واقعاً کلان شوند. تنوع در کلان داده‌ها به این معنی است که مجموعه داده‌های موجود در کلان داده شامل انواع داده‌های ساختاریافته^۱ (داده‌هایی که در پایگاه داده‌های رابطه‌ای ذخیره می‌شوند؛ مانند مجموعه نمرات دانشجویان یک کلاس)، داده‌های نیمه‌ساختاریافته^۲ (داده‌هایی که دارای قالب نسبتاً مشخصی هستند ولی امکان دارد که خصوصیت‌های مختلفی داشته باشند؛ مانند فرمت‌های JSON و XML) و داده‌های غیر‌ساختاریافته^۳ (داده‌هایی که هیچ‌گونه ساختار منظمی در آن‌ها وجود ندارد؛ مانند تصویر، فایل‌های چندرسانه‌ای و داده‌های جمع آوری شده از شبکه‌های اجتماعی) است. منابع تولید داده می‌تواند هم ماشین باشد و هم انسان.

نرخ تولید

نرخ تولید یا سرعت تولید تداعی کننده این امر است که مجموعه داده‌ها در کلان داده در مدت زمانی کم به میزان بالایی افزایش یابد؛ یعنی منابع تولید داده با سرعت بالایی داده تولید می‌کنند و حجم داده با آهنگ بالایی رشد می‌نماید. داده‌های بلاذرنگ^۴ نیز در این دسته قرار می‌گیرند؛ چون در مورد داده‌های بلاذرنگ، سرعت تولید و پردازش داده اهمیت بالاتری از حجم داده دارد. به عنوان مثال شرکت فیسبوک روزانه نزدیک به ۵۰۰ میلیون تصویر را در پایگاه داده‌های خود ذخیره می‌نماید و حجم پایگاه داده‌ی آن با یک آهنگ نمائی رشد می‌کند.

Structured^۱

Semi-structured^۲

Unstructured^۳

Real-time^۴