

کلان داده و فناوری هدوپ

سحر جعفری

دانشجوی مهندسی کامپیوتر، دانشگاه صنعتی شیراز، شیراز، ایران

چکیده

کلان داده مربوط به داده‌هایی با حجم زیاد است که به صورت نمایی در حال رشد می‌باشد؛ این داده حجیم با یک سرعت زیاد از منابع مختلف و در انواع مختلف ساختاریافته، بدون ساختار و نیمه ساختاریافته تولید می‌شود که ما می‌توانیم اطلاعات ارزشمندی را از آن استخراج کنیم و در تصمیم‌گیری‌ها از کمک‌های آن بهره‌مند شویم. این پژوهش به منظور آشنایی بیشتر با کلان داده و فناوری هدوپ بوده و یافته‌ها حاکی از آن است که به طور کلی کلان داده توسط سه ویژگی اساسی خود، یعنی حجم (volume)، سرعت (velocity) و تنوع (variety) مشخص می‌شود که این سه ویژگی باید هم‌زمان وجود داشته باشند در غیر این صورت نمی‌توان درباره داده‌های بزرگ صحبت کرد. برخی از محققان برای بیان بهتر کلان داده، ویژگی‌های دیگری را نیز معرفی کرده‌اند از جمله ارزش (value) و صحت (veracity). تحلیل کلان داده با فراهم کردن اطلاعات ارزشمند، می‌تواند در حوزه‌های مختلف پزشکی، تجارت و سیاست بسیار کمک کننده باشد؛ اما استفاده از روش‌های سنتی برای ذخیره سازی و پردازش کلان داده کاری وقت‌گیر و هزینه‌بر است به همین خاطر فناوری‌هایی مانند هدوپ از طریق برقرار کردن امکان ذخیره‌سازی هر نوع داده در یک محیط توزیع‌شده و پردازش آن‌ها به صورت موازی به کمک ما آمده‌اند. آپاچی هدوپ از سه بخش سیستم فایل توزیع‌شده (HDFS)، چهارچوب برنامه‌نویسی نگاشت‌کاهش (MapReduce) و سرویس مدیریت منابع (YARN) تشکیل شده است که به ترتیب به عنوان واحد ذخیره‌سازی، واحد پردازش و واحد مدیریت منابع در هدوپ مورد استفاده قرار می‌گیرند و از این طریق مدیریت کلان داده برای ما میسر می‌شود.

واژه‌های کلیدی: کلان داده، هدوپ، نگاشت‌کاهش، سیستم‌فایل توزیع‌شده، HDFS، YARN، MapReduce.

مقدمه

داده‌ها با پیشرفت تکنولوژی روز به روز در حال گسترش هستند؛ IBM نشان داد که هر روز ۲.۵ اگر بایت داده تولید می‌شود و همچنین ۹۰٪ از داده‌ها در ۵ سال اخیر تولید شده است [1]. گوشی‌های موبایل، حسگرهای محیطی، لاگ نرم‌افزارهای مختلف، دوربین‌ها، ایستگاه‌های هواشناسی، شبکه‌های اجتماعی، اطلاعات پزشکی، اطلاعات سامانه‌های خرید از فروشگاه‌ها و غیره، بخشی از منابع تولید داده در مقیاس کلان هستند که با کمک این داده‌های حجیم می‌توان مسیر حرکت کسب و کار و فرایند چرخش کار سازمان‌ها را مشخص کرد. بنابراین ذخیره‌سازی و پردازش این داده‌ها از اهمیت بالایی در صنعت و تجارت برخوردار است و نیازمند استفاده از سیستم‌ها و چهارچوب‌های مدرنی است که با رویکرد‌های سنتی پردازش داده متفاوت باشد.

بهترین راهکار برای ذخیره‌سازی کلان‌داده استفاده از محیط توزیع‌شده می‌باشد؛ همچنین استفاده از مدل برنامه‌نویسی نگاشت‌کاهش (MapReduce) روشی بهینه برای پردازش موازی و مدیریت کلان‌داده است که توسط دانشمندان علم داده مورد استفاده قرار می‌گیرد. در این تکنیک داده‌ی ما با توجه به یک سری صفت‌ها به دسته‌ای نگاشت می‌شود، سپس تمام داده‌های نگاشت شده با هم ادغام می‌شوند و خروجی نهایی را تولید می‌کنند. تمام این امکانات در یک چهارچوبی به اسم هادوپ فراهم شده است که اجازه می‌دهد مدیریت و پردازش کلان‌داده با سرعت و اثرگذاری بیش‌تری انجام شود.

انواع داده

داده‌ها را می‌توان در سه شکل مختلف ساختاریافته، بدون ساختار و نیمه ساختاریافته تقسیم بندی نمود.

داده‌های ساختاریافته

هر داده‌ای که قابلیت ذخیره شدن، قابلیت دسترسی و پردازش را داشته باشد و به یک فرمت ثابت نیز باشد، به عنوان یک داده ساختار یافته در نظر گرفته می‌شود [2]. داده‌هایی که در پایگاه داده‌های رابطه‌ای ذخیره می‌کنیم از نوع ساختاریافته هستند [3].

داده‌های بدون ساختار

هر داده‌ای با فرمت یا ساختار ناشناخته، به عنوان داده‌های بدون ساختار طبقه‌بندی می‌شود. داده‌های بدون ساختار علاوه بر داشتن اندازه بزرگ، چالش‌های متعددی را از نظر پردازش آن برای استخراج ارزش آن داده‌ها، به وجود می‌آورد [2]. انواع مختلف فایل‌های متنی، تصویری، ویدیویی و صوتی جز داده‌های بدون ساختار هستند [3].

داده‌های نیمه ساختاریافته

داده‌های نیمه ساختار یافته می‌توانند هر دو نوع داده (ساختار یافته و غیرساختاریافته) را شامل شوند. ما می‌توانیم داده‌های نیمه ساختار یافته را به صورت ساختاریافته ببینیم اما این داده‌ی نیمه ساختاریافته، به صورتی که بخواهیم آنها را به عنوان یک جدول رابطه‌ای در نظر بگیریم، در DBMS تعریف نشده است [2]. فایل‌های XML و JSON نمونه‌ای از داده‌های نیمه ساختاریافته می‌باشند [3].

کلان داده

اگر بخواهیم تعریفی از بیگ دیتا یا کلان داده ارائه کنیم میتوانیم آن را مجموعه ای از داده های حجیم بدانیم که توسط منابع مختلف مثل شبکه های اجتماعی، سنسورها، دستگاه های IOT و بسیاری منابع دیگر و در انواع مختلف ساختاریافته، بدون ساختار و نیمه ساختاریافته تولید می شوند، به طوری که نتوان آن ها را با روش های معمول، ذخیره، پردازش و تحلیل کرد [1]. [3]

نمونه هایی از کاربرد کلان داده

در اینجا به برخی از کاربردهای کلان داده در سه حوزه سلامت، تجارت و سیاست اشاره خواهیم کرد.

کاربرد کلان داده در بخش سلامت

با تجزیه و تحلیل اطلاعات در مورد یک بیماری می توان درکی خوب از پیشروی بیماری به دست آورد و از این طریق می توان اقدامات پیشگیرانه ای انجام داد همچنین با تجزیه و تحلیل روی جراحی های قبلی موجود در پرونده بیماران و ترکیب آن با علائم فعلی یا سوابق بیمار می توان بهترین روش های درمان را بر اساس مشخصات بیمار پیدا کرد [4].

کاربرد کلان داده در بخش تجارت

با تجزیه و تحلیل داده های بزرگ در بخش تجارت می توان رفتار کاربران و علایق آن ها را درک کرد؛ همین موضوع باعث افزایش رضایت مشتری، افزایش سود و افزایش رقابت در سازمان ها می شود به طور کلی تحلیل کلان داده، در دستیابی به فرصتهای تجاری و پیش بینی رکود، به سازمان های تجاری کمک میکند [4].

به عنوان مثال در تحقیقی در سنگاپور با استفاده از ۵۰۰ شرکت کننده در صنعت خرده فروشی انجام گردید. نتایج مطالعه نشان داد که در میان تحلیل های مختلف داده های بزرگ مورد استفاده در صنعت خرده فروشی سنگاپور، تحلیل های رسانه های اجتماعی به طور عمده توسط شرکت کنندگان پاسخ داده شده است. تأثیر داده های بزرگ بر عملکرد سازمان مطابق نتایج این مطالعه شامل افزایش بیشتر فروش و کاهش هزینه می باشد، اما بیشتر پاسخ ها به نفع فروش بیشتر بوده است و در صورت تأثیر بر خدمات مشتری، رضایت بالای مشتریان مشاهده شد. این مطالعه عمدتاً بر استراتژی در دو راهی قرار گرفتن به عنوان رویکرد اصلی در مدیریت داده های بزرگ و الگوی VLDB به عنوان روش شناسی اصلی تأکید دارد. از نظر کاهش موانع در زمینه مدیریت داده های بزرگ، استخدام منابع انسانی ماهر در تجزیه و تحلیل های مختلف بوسیله بیشتر شرکت کنندگان انتخاب شده بود. علاوه بر این، برای بستر قابلیت های احتمالی بر بازاریابی سفارشی یا سفارشی سازی انبوه تأکید گردید [5].

کاربرد کلان داده در بخش سیاست

با تحلیل اطلاعات دوربین های نظارت دولتی و خصوصی، نظرات شهروندان در شبکه های اجتماعی، معاملات آنلاین، داده های GPS و ارتباطات تلفن همراه می توان کنترل و نظارت پیوسته، برای محافظت از شهروندان و کاهش جرایم داشت [4].

خصوصیات کلان داده و چالش‌های آن

یکی از چالش‌های بزرگ در حوزه‌های مختلف سیستم محاسباتی، داده‌های بزرگ بوده و این یک موضوع اصلی در کسب و کار هوش داده ای است. دو تکنیک اصلی برای پردازش داده‌های بزرگ وجود دارد: پردازش دسته ای و پردازش جریانی. در پردازش دسته ای، داده‌ها ابتدا در پایگاه داده‌های بزرگی ذخیره میشوند و بعداً پردازش می‌گردند؛ اینکار معمولاً با مدل‌های برنامه نویسی مقیاس پذیری مانند Google's MapReduce انجام میشود. با این حال با اندازه رو به رشد داده ها، هزینه انتقال و ذخیره سازی آنها قابل جلوگیری نیست. بعلاوه در بسیاری از دامنه ها، چیزی که مهم است نگه داشتن داده‌های اولیه نیست بلکه تحلیل آنها در سریعترین زمان ممکن است تا اطلاعات ارزشمندی را ایجاد کنند. سیستم‌های پردازش جریانی برای حل این مسائل، بر واکنش پذیری و تحلیل داده‌ها به محض تولید شدن آنها، تأکید دارند. سالهای اخیر شاهد پیدایش راه حل‌های مختلفی از پردازش جریانی بوده است. بیشتر مطالعات در این خصوص بر روی تأثیر محفظه‌های اجرایی مختلف بر عملکرد یک سیستم پردازش جریانی ارتجاعي تمرکز دارند و در اینگونه تحقیقات سلسله مراتب محفظه‌های اجرایی (ماشین‌ها و فرآیندها) را بررسی میکنند و به این نتیجه رسیده اند که فراهم سازی آنها با هزینه‌های مختلفی انجام می‌شود و مهمتر اینکه فراهم سازی نوع اشتباهی از محفظه‌ها میتواند سبب افت عملکرد گردد [6].

حجم، سرعت و تنوع سه مورد از ویژگی‌های اصلی کلان داده می‌باشند که به عنوان 3V's مطرح شدند [1] اما در ادامه ویژگی‌های بیش‌تری توسط محققان برای کلان داده مطرح شده است:

- حجم (volume): به اندازه داده اشاره دارد که به صورت نمایی در حال رشد است؛ معمولاً کلان داده دارای حجمی بیش‌تر از ترابایت و پتابایت می‌باشد [7].

- سرعت (velocity): داده‌های حجیم توسط منابع مختلف با سرعت بسیار زیاد و به صورت بلادرنگ تولید میشوند [7].

- تنوع (variety): داده‌های حجیم با سرعت زیاد توسط منابع مختلف و در انواع متفاوت تولید میشوند، بیش‌تر حجم داده‌های دنیا بدون ساختار و بسیار متنوع هستند. امروزه بخشی از داده‌ها در بانک‌های اطلاعاتی، بخشی در صفحات وب و بقیه نیز در فایل‌ها با قالب‌های متفاوت ذخیره شده اند که پردازش آن‌ها پیچیده می‌باشد [1].

- ارزش (value): باید بتوان از داده حجیم که در انواع مختلف و با سرعت زیاد تولید می‌شود اطلاعات ارزشمندی را استخراج کرد به عبارت دیگر باید دید که این داده از نظر اطلاعاتی برای تصمیم‌گیری چقدر دارای ارزش است [1].

- صحت (veracity): با توجه به اینکه داده‌ها توسط منابع مختلفی ایجاد می‌شوند؛ این موضوع بر این دلالت دارد که داده‌ها قابل اطمینان و دارای اطلاعات درست باشد [1].

با توجه به خصوصیات کلان داده که در بالا اشاره کردیم؛ در هنگام مواجهه با بیگ‌دیتا با چالش‌هایی روبه‌رو خواهیم شد. اولین چالش مربوط به ذخیره کردن کلان داده می‌باشد با توجه به حجم کلان داده که فراتر از ترابایت و پتابایت است هراندازه هم که واحد ذخیره سازی را بزرگ در نظر بگیریم باز هم جوابگوی آن حجم از داده نخواهد بود [3]. دومین چالش مربوط به پردازش کلان داده می‌باشد که کاری وقت‌گیر و هزینه‌بر است و عملاً روش‌های پردازش سریالی قادر به پردازش این حجم از داده نخواهند بود [3]. سومین چالش عدم توانایی سیستم‌ها در پردازش داده‌های بدون ساختار می‌باشد با توجه به اینکه بیش‌تر حجم داده‌های دنیا بدون ساختار و بسیار متنوع هستند این موضوع بسیار پر اهمیت می‌باشد [3].

این چالش‌ها منجر به ظهور پلتفرم‌های جدیدی مانند Apache Hadoop شده است که می‌تواند مجموعه داده‌های بزرگ را به راحتی مدیریت کند [8].

تعریف هادوپ

با افزایش داده های بزرگ، بنیاد نرم افزار آپاچی در سال ۲۰۰۸ یک چارچوب متن باز به نام Apache Hadoop ایجاد کرد که راه حلی برای تمام مشکلات کلان داده است. هادوپ یک چهارچوب متن باز (open source) می باشد، که به زبان جاوا نوشته شده است، برای ذخیره سازی توزیع شده و پردازش حجم عظیمی از مجموعه داده ها مورد استفاده قرار میگیرد [3]. متن باز بودن به این معنی است که به صورت رایگان در دسترس است و حتی ما می توانیم سورس کد آن را مطابق با نیاز خود تغییر دهیم. هادوپ یک سیستم بسیار مقاوم در برابر خطا و بسیار در دسترس است [3].

مواردی که باعث شده تا هادوپ بعنوان یک ابزار برای مدیریت کلان داده مورد استفاده قرار بگیرد، عبارتند از:

- امکان ذخیره سازی هر نوع داده از جمله ساختاریافته، نیمه ساختاریافته و بدون ساختار بدون هیچ محدودیتی [3].
- توانایی انجام پردازش های پیچیده و سرعت بخشیدن به پردازش ها به این طریق که در هادوپ به جای آنکه داده ها به سمت برنامه ی در حال اجرا منتقل شوند کدهای برنامه ی پردازشی به سمت داده ها ارسال می شوند در حقیقت کدهای برنامه بین ماشین های کلاستر توزیع می شوند. کلاستر (خوشه) به معنای گروهی از سیستم ها است که از طریق یک شبکه ارتباطی به هم متصل شده اند. این موضوع باعث استفاده بهینه از پهنای باند کلاستر شده و از ایجاد بار اضافی ناشی از انتقال داده ها در بین گره های کلاستر جلوگیری می کند و اینگونه سرعت افزایش میابد [3].
- هادوپ می تواند روی یک دستگاه معمولی نصب شود. همین موضوع هادوپ را، هم از نظر اقتصادی به صرفه می کند و هم باعث می شود دسترسی پذیری بیش تری داشته باشد چون وابسته به یک فروشنده خاص نیست [3].

ساختار هادوپ

هادوپ متشکل از سه بخش اصلی HDFS، MapReduce و YARN می باشد [9] که در ادامه هر یک از این بخش ها را با جزییات بیش تر شرح خواهیم داد.

سیستم فایل توزیع شده هادوپ (HDFS)

HDFS واحد ذخیره سازی در هادوپ می باشد [8] که امکان ذخیره سازی حجم عظیمی از داده را در چندین گره از یک کلاستر فراهم می کند، به عبارتی ما داده خود را در یک محیط توزیع شده ذخیره می کنیم. سیستم فایل توزیع شده ی هادوپ بصورتی توزیع شدگی را ایجاد میکند که از حافظه ی اختصاصی هر یک از دستگاه ها برای ذخیره سازی داده ها استفاده می کند [3]. برای ذخیره سازی داده در HDFS، داده ی ما ابتدا به چندین بلاک با سایز پیش فرض 128 MB تقسیم می شود و سپس در گره ها ذخیره می شوند. HDFS از دو بخش name node و data node تشکیل شده است [8]. Name node فقط فراداده ها (داده هایی درباره داده) را ذخیره می کند [9]؛ در حالی که data node گره هایی هستند که شامل داده های حقیقی یا بلاک های داده هستند [8].

نگاشت کاهش (MapReduce)

نگاشت کاهش واحد پردازش در هادوپ می باشد [8]. نگاشت کاهش یک تکنیک برنامه نویسی است که امکان اجرای پردازش موازی بروی مجموعه بزرگی از داده ها در یک محیط توزیع یافته را می دهد [3]. نگاشت/کاهش شامل دو وظیفه مجزا است یعنی نگاشت و کاهش و همانطور که از نام آن مشخص است، بعد از اجرای کامل فاز نگاشت، فاز کاهش اجرا میشود. خروجی

فاز نگاشت، جفت‌های کلید-مقدار می‌باشد که به عنوان ورودی در فاز کاهش در نظر گرفته می‌شوند [8] [3] این ورودی‌ها در تابع کاهش پردازش می‌شوند و سپس به مجموعه‌های کوچک‌تر تبدیل می‌شوند و با یکدیگر ادغام می‌شوند و خروجی نهایی را تولید می‌کنند، خروجی نهایی در HDFS ذخیره می‌شود [3]. با استفاده از این تکنیک درخواست پردازشی ارسال شده توسط کاربر به بخش‌های مستقلی تقسیم می‌شود که هر کدام از آن‌ها بین گره‌هایی که حاوی داده هستند پخش می‌شوند و به صورت موازی و همزمان اجرا می‌شوند [3]. در نتیجه‌ی این کار داده‌های مورد نظر به بخش‌های کوچک‌تر می‌شکنند و یک سری مقادیر به یک سری مقادیر دیگر نگاشت می‌شوند و جفت‌های کلید-مقدار به وجود می‌آیند سپس این جفت‌ها بر اساس کلیدشان طبقه بندی می‌شوند و خروجی‌های این مرحله وارد فاز کاهش می‌شوند و در نهایت یک خروجی تولید می‌شود [8].

سرویس مدیریت منابع (YARN)

YARN واحد مدیریت منابع در هدوپ می‌باشد که زمان‌بندی کارها را نیز انجام می‌دهد [8]. در یک کلاستر با چندین گره مدیریت منابع کاری دشوار است به همین خاطر با کمک YARN منابع به شکل کاملاً کارآمد مدیریت می‌شوند [3]. برای اینکه YARN بتواند هدف مدیریت منابع را محقق کند احتیاج به دو بخش resource manager و node manager دارد [8] [9]. Resource manager وظیفه مدیریت کل منابع موجود در یک خوشه را بر عهده دارد. Node manager وظیفه مدیریت و کنترل اجرای منابع موجود در یک گره را عهده دار است [8]. هر گره خود شامل دو بخش container و application master می‌باشد [8] [9] که در مدیریت منابع به node manager کمک می‌کنند؛ container مجموعه‌ای است که منابع فیزیکی مانند RAM و CPU را در خود دارد که application master درخواست نیاز به منابع را از طرف node manager به سمت container ارسال می‌کند و در آخر node manager نیز گزارشی از منابع خود را به سمت resource manager می‌فرستد [8]. به طور کلی به ازای هر خوشه تنها یک resource manager داریم و به ازای هر گره یک node manager داریم.

به طور کلی هدوپ به سبک MasterSlave (ارباب و برده) کار می‌کند. در یک خوشه هدوپ تنها یک گره ارباب وجود دارد و تعداد زیادی گره برده که این تعداد می‌تواند تا هزاران گره نیز ادامه پیدا کند. گره ارباب، گره‌های برده را مدیریت، حفظ و کنترل می‌کند در حالی که گره‌های برده، مولفه‌های واقعی انجام کار هستند [3]. Name node و resource manager در گره‌های master قرار دارند در حالی که data node و node manager در گره‌های slave قرار دارند [8]. گره‌های برده هر سه ثانیه سیگنال‌هایی به سمت گره‌های ارباب ارسال می‌کنند و گزارشی از وضعیت خود می‌دهند [8].

نتیجه گیری

همانطور که در گفته شد، کلان‌داده توسط سه ویژگی اساسی خود، یعنی حجم، سرعت و تنوع، مشخص می‌شود که این سه ویژگی باید هم‌زمان وجود داشته باشند البته برخی از محققان ویژگی‌های دیگری همچون "ارزش و صحت" را نیز به این ویژگی‌ها اضافه کرده‌اند. کلان‌داده می‌تواند در حوزه‌های مختلف پزشکی، تجارت و سیاست بسیار مفید باشد؛ به همین دلیل بجای استفاده از روش‌های سنتی برای ذخیره سازی و پردازش کلان‌داده از فناوری‌هایی همچون هدوپ استفاده می‌شود. به طور کلی برای مدیریت کلان‌داده می‌توانیم از چهارچوب هدوپ که تحت بنیاد آپاچی است استفاده کنیم و برای مدیریت

مطلوب بر کلان داده‌ها، آپاچی هدوپ از سه بخش سیستم فایل توزیع شده (HDFS)، چهارچوب برنامه‌نویسی نگاشت کاهش (MapReduce) و سرویس مدیریت منابع (YARN) استفاده می‌کند. با استفاده از هدوپ می‌توانیم داده حجیم خود را در یک سیستم توزیع شده ذخیره کنیم و آن را به صورت موازی با استفاده از تکنیک MapReduce در گره‌های محاسباتی موجود در خوشه‌ها پردازش کنیم که همه منابع موجود در خوشه‌ها برای انجام کارهای پردازشی و دیگر کارها به شکل کاملاً کارآمد، توسط YARN مدیریت می‌شوند.

منابع

۱. و. جانی، "داده های حجیم و نقش آن در بهبود تجارت الکترونیک"، در سومین کنفرانس بین‌المللی پژوهش در علوم و تکنولوژی، ۱۳۹۵.
2. "upGrad," [Online]. Available: <https://www.upgrad.com/blog/what-is-big-data-types-characteristics-benefits-and-examples/>
۳. س. ص. م. واهبی، "بررسی سرویس پردازش موازی هادوپ در تحلیل کالان داده"، در پنجمین کنفرانس بین‌المللی مهندسی برق، کامپیوتر و مکانیک، ۱۳۹۹.
۴. ن. ا. پ. م. جوانمرد، "داده های بزرگ، مروری بر فرصت ها، کاربردها و ابزارهای مورد استفاده"، در کنگره بین‌المللی مطالعات میان رشته‌ای در علوم پایه و مهندسی، ۱۳۹۶.
5. S. S. S. A. C. C. J. S. Song Ying, Managing big data in the retail industry of Singapore: Examining the impact on customer satisfaction and organizational performance, European Management Journal, 2020 .
6. V. Marangozova-Martin, Multi-Level Elasticity for Data Stream Processing, Grenoble Alpes, CNRS, LIG, F-38000 Grenoble France: Noël de Palma and Ahmed El Rheddane Univ, 2019 .
7. F.-Z. B. ., A. A. L. ., B. Ahmed Oussous, "Big Data technologies," in Journal of King Saud University Computer and Information Sciences, 2017 .
8. "TechVidvan," [Online]. Available: <https://techvidvan.com/tutorials/how-hadoop-works-internally/>
9. "TechVidvan2," [Online]. Available: <https://techvidvan.com/tutorials/apache-hadoop-tutorials/>

Big Data and Hadoop Technology

Sahar Jafari

Computer engineering student, Shiraz University of Technology, Shiraz, Iran,

Abstract

Big data is related to data with a large volume that is growing exponentially; This massive data is produced at a high speed from different sources and in different types of structured, unstructured and semi-structured, from which we can extract valuable information and benefit from its help in making decisions. In general, big data is characterized by three characteristics. Its basic characteristics, namely volume, velocity and variety, are defined. These three characteristics must exist at the same time, otherwise we cannot talk about big data. In order to better express big data, some researchers include other characteristics have introduced, including value and veracity. By providing valuable information, big data analysis can be very helpful in various fields of medicine, business and politics; But using traditional methods to store and process big data is time-consuming and expensive, that's why technologies like Hadoop have come to our aid by making it possible to store any type of data in a distributed environment and process them in parallel. Apache Hadoop from There are three parts of distributed file system (HDFS), map reduction programming framework (MapReduce) and resource management service (YARN) which are used as storage unit, processing unit and resource management unit respectively in Hadoop and through this management Big data becomes available to us.

Keywords: Big Data, Hadoop, MapReduce, Distributed File System, HDFS, YARN, MapReduce.
