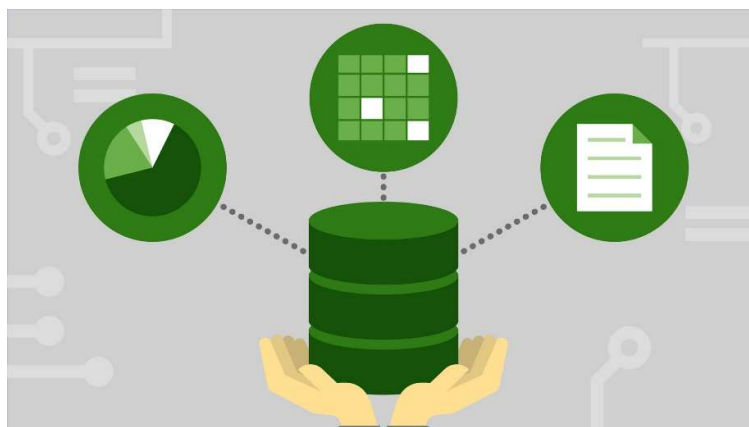


به نام خدا



دانشگاه تهران  
پردیس دانشکده‌های فنی  
دانشکده برق و کامپیوتر



## آزمایشگاه پایگاه داده

دستورکار شماره ۱۰

مهلت تحویل :

۱۴۰۱/۰۳/۲۰

مجتبی بنائی

## دستور کار شماره ۱۰ - کار با آپاچی دروید و کافکا

یکی از دیتابیس‌هایی که اخیراً و به صورت روزانه با آن سروکار دارم و به کارگیری آنرا در بسیاری از شرکت‌های متوسط و بزرگ ایران برای ذخیره و پاسخگویی به حجم عظیم داده‌های ورودی، یک ضرورت میدانم، **دیتابیس تحلیلی آپاچی دروید** است.

این دیتابیس که در رسته بانک‌های اطلاعاتی تحلیلی به عنوان نسل جدید دیتاویزهوس‌ها قرار میگیرد، ویژگی‌های کاربردی بسیار خوبی دارد. به عنوان نمونه:

– یک دیتابیس تحلیلی است و ذخیره داده‌ها به صورت ستونی در آن، هم حجم داده‌ها را بسیار کم میکند و هم سرعت پاسخگویی به کوئری‌های تحلیلی که معمولاً ترکیبی از گروه‌بندی و فیلترینگ داده‌ها هستند را بسیار بهبود می‌بخشد.

– معماری توزیع شده بسیار پیشرفته‌ای دارد که مقیاس‌پذیری و پاسخ‌گویی به هر حجمی از داده‌ها در زمان مناسب (زیر یک ثانیه) تضمین میکند (البته نیاز به تخصیص منابع و پایش دارد)

– یک دیتابیس سری زمانی است و تمامی داده‌ها باید حاوی مهر زمان یا تایم‌استمپ باشند.

– بخش دریافت داده یا اینجسشن بسیار پیشرفته‌ای دارد و کافی است داده‌های خود را به `Kafka/MinIo/HDFS` منتقل کنید و ادامه کار یعنی دریافت از کافکا – و سایر منابع –، فیلترینگ و پردازش اولیه داده‌ها و ذخیره آنها به صورت بلادرنگ را به دروید بسپارید.

– امکان تجمیع داده‌ها به صورت خودکار در آن فراهم شده است یعنی مثلاً داده‌های هر پنج دقیقه را تجمیع کرده، آمار مورد نیاز آن بازه زمانی را به صورت خلاصه برای شما محاسبه میکند.

– الگوریتم‌های تقریبی پیشرفته‌ای دارد که همزمان با تجمیع داده‌ها میتوانید با تقریب بسیار مناسبی، آماره‌هایی راجع به داده‌هایی که در حین تجمیع، حذف می‌شوند را ذخیره کنید.

با این مقدمه، دستور کار آخر را برای آشنایی با این دیتابیس خوش‌آبیه و نیز بروکر (توزیع‌کننده پیام در یک شبکه) کافکا در نظر گرفته‌ام.

توضیح اینکه امروزه کافکا نقش محوری در توسعه سامانه‌های مقیاس‌پذیر اطلاعاتی دارد و نقش واسط بین سامانه‌های مختلف و نیز انتقال توزیع شده اطلاعات در یک شبکه را بر عهده دارد و در این دستور کار هم، داده‌هایی را به صورت تصادفی تولید و در یک کانال کافکا (هر پیام در یک گروه قرار میگیرد که به آن، تاپیک می‌گوئیم) قرار میدهم و به کمک موتور اینجسشن دروید، آنها را به صورت مداوم دریافت و در دیتاسورس‌های متناظر دروید ذخیره می‌کنیم.

آموزش یک ساعته‌ای را برای این دیتابیس و نحوه دریافت اطلاعات از کافکا برای مجموعه آموزشی نیک آموز در یک وینار یک ساعته در انتهای آذر ماه برگزار کرده‌ام. فایل ویدئوی این کارگاه عملی را از آدرس زیر دریافت کنید:

<https://dl6.nikamooz.com/webcast1400/ApachDruid.rar>

نکته: حجم ویدئوها، حدود ۱.۴ گیگابایت است و ممکن است دانلود آنها برای شما هزینه‌بر باشد. می‌توانید ویدئوهای فوق را در سه قسمت در **کانال آپارات مهندس داده**<sup>۱</sup> هم به صورت آنلاین مشاهده کنید (قسمت اول و دوم برای این گزارش کار کافی است) و کیفیت پخش آنرا به دلخواه تنظیم کنید که حجم اینترنت بسیار کمتری برای این موضوع، مصرف شود.

<sup>1</sup> <https://www.aparat.com/playlist/1265528>

فایل‌های کارگاه عملی هم که برای این دستور کار به آنها نیاز خواهید داشت در کنار دستور کار، در قالب یک فایل زیپ بارگزاری شده است اما می‌توانید آنها را از آدرس زیر (پوشه druid) هم دانلود کنید :

[https://gitlab.com/nikamooz\\_bigdata/webinars](https://gitlab.com/nikamooz_bigdata/webinars)

بعد از مشاهده این ویدئو و آشنایی با نحوه نصب و راه اندازی دروید و همینطور کلاستر کافکا (کلاستر : مجموعه نودهایی که با هم، یک هدف واحد را در شبکه بر عهده دارند)، برای دستور کار شماره ۱۰ که ترکیب دروید و کافکا خواهد بود، موارد زیر را انجام و گزارش آنرا آپلود کنید :

- دروید را دانلود و راه اندازی کنید.

- کلاستر کافکا را به کمک داکر راه اندازی نمایید (در خط فرمان، درون پوشه داکر، دستور `docker-compose up` را اجرا کنید)

- برای اینکه بتوانید بین سیستم خودتان و کافکایی که در داکر راه اندازی شده است، ارتباط برقرار کنید و دروید هم بتواند به آن متصل شود نیاز دارید که سطر

127.0.0.1 kafka

را به فایل `hosts` ویندوز یا لینوکس اضافه کنید. (`kafka` نام بروکر یا نود کافکایی است که درون فایل داکر کامپوز تعریف شده است و اگر در فایل داکر کامپوز، چندین بروکر کافکا داشته باشید نام تمام آنها را در این جا باید وارد کنید - هر نام در یک سطر) یک راهنمای ویرایش فایل `hosts` ویندوز را در این آدرس می‌توانید ببینید :

<https://howtotech.ir/what-is-windows-hosts-file-how-to-change-hosts-file>

- چک کنید که واسط کاربری کافکا یعنی `AKHQ` در آدرس `localhost:8000` قابل مشاهده باشد. تایپ `events_topic` با ۳ پارتیشن را در کافکا ایجاد نمایید.

- فایل `user_event_producer.py` موجود در پوشه `codes` را اجرا کنید به گونه ای که درون واسط کاربری `AKHQ` هر نیم ثانیه یک بار، پیامی را مشاهده کنید (می‌توانید تایپیک را انتخاب کنید و از بخش `Live tail` آخرین پیام‌ها را ببینید)

- درون دروید، طبق آموزش ویدئویی فوق، یک اینجسشن برای دریافت داده‌های موجود در تایپیک `events_topic` کافکا ایجاد کنید به گونه ای که در پایان کار، دیتاسورس مربوطه ایجاد شود.

- در مرحله آخر، یک `Rollup` یک ساعته و یک `Rollup` یک روزه با توجه به دیتای تولید شده ایجاد کنید. منظور از رولاپ، جمع‌آوری اطلاعات در یک بازه زمانی است. برای این منظور، مجدداً باید یک اینجسشن بر روی همین تایپیک ایجاد کنید و سپس در بخش `configure schema` گزینه `rollup` را علامت بزنید و سپس فیلدهایی که داده‌ها بر اساس آنها در آن بازه زمانی باید گروه بندی شوند را مشخص کنید (دایمنشن یا ابعاد) و سپس متریک‌های مورد نیاز مانند مجموع یا تعداد یا ماکزیمم و می‌نیمم برخی فیلدها را تعیین کنید. با اینکار، به محض دریافت داده‌ها، بسته به متریک تعیین شده، داده‌ها جمع‌آوری می‌شود. با اینکار جدولی ایجاد می‌کنیم که خلاصه داده‌های یک ساعت یا یک روز در آن ذخیره شده است و نیاز به کوئری گرفتن از جداول اصلی نداریم. از طرفی می‌توانیم فقط این جداول تجمیعی را نگهداری کنیم و خود داده‌های اصلی را حذف و یا به یک دریاچه داده منتقل کنیم. برای تجمیع یکساعته، `query granularity` را برابر یک ساعت و برای تجمیع یک روزه، آنرا برابر روز قرار دهید. با اینکار فیلد `time` براساس این مقدار، در جدول نهایی رند خواهد شد.

با ایجاد یک تجمیع یکروزه ، تمام داده های آنروز در یک رکورد خلاصه شده و در دیتابیس ذخیره خواهند شد. بنابراین ، جدول تجمیعی روزانه باید کمترین میزان رکوردها را داشته باشد و دیتاسورس غیر تجمیعی، بیشترین تعداد رکورد را به خود اختصاص دهد.

- پنج کوثری مختلف بر روی این سه دیتاسورس اجرا و نتیجه را نمایش دهید. کوثری های شما حتما شامل گروه بندی ، بازه زمانی ، تعیین ترتیب و یک شرط عادی باشند.

از مراحل فوق و نیز کوثری ها و نتایج آنها، عکس گرفته و در گزارش کار آخر آنها را ذکر کنید . (در حدی که بیان کننده کار شما باشد )