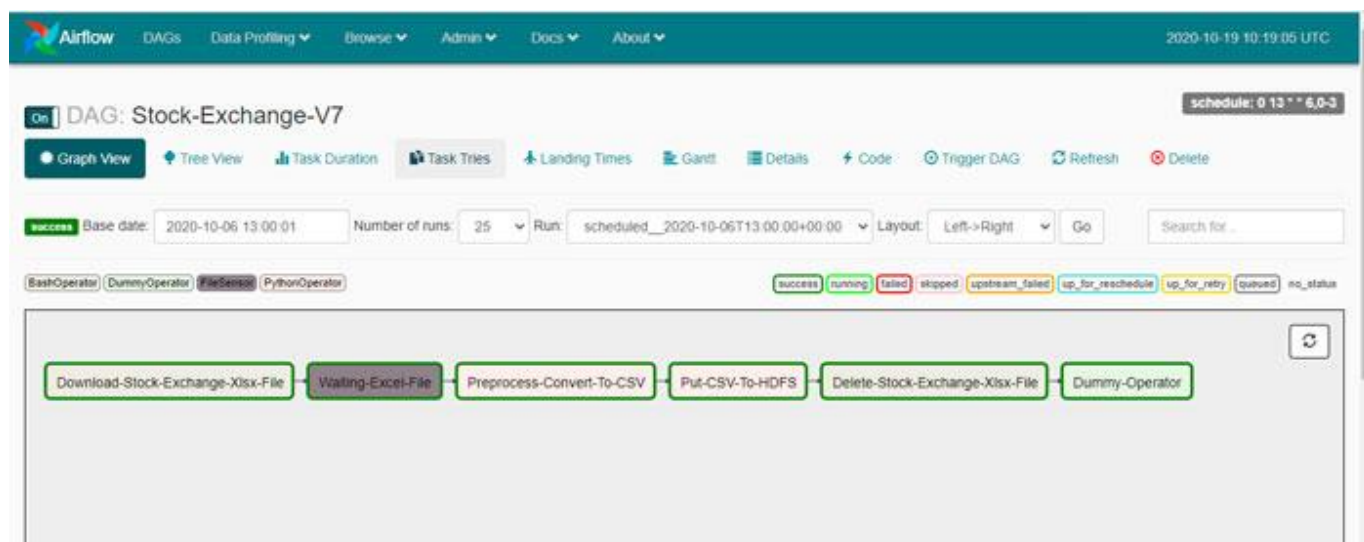


سرفصل‌های این وبکست آموزشی

- آشنایی با تعاریف و مفاهیم پایه مهندسی داده (۳۰ دقیقه)
- انجام مثال عملی پردازش داده‌های روزانه بورس ایران (حدود یک ساعت و ربع)

در این مثال کاربردی، داده‌های روزانه بورس ایران از سایت tsetmc به صورت خودکار و روزانه دانلود شده، بعد از پردازش اولیه و تبدیل قالب آن از اکسل به CSV، درون HDFS (سیستم فایل هدوپ) ذخیره می‌گردد. سپس با تعریف جداولی در Hive و با استفاده از Hue، انواع کوئری‌های تحلیلی مانند بررسی شرکت‌هایی که نرخ خرید حقوقی به حقیقی آنها بیشتر از ۲ است، شرکت‌هایی که حجم معاملات روزانه‌شان از میانگین سه ماهه آنها بیشتر است و ... را روی این داده اجرا خواهیم کرد. بنابراین در مثال اول، با فناوریهای زیر آشنا خواهیم شد:

- Apache Airflow به عنوان محور کار و هماهنگ کننده اصلی تسک‌های روزانه
 - Hadoop برای ذخیره و اجرای کوئری‌های Hive
 - Hive به عنوان یک موتور پردازشی SQL بر روی هدوپ و داده‌های حجیم.
 - Hue به عنوان یک محیط گرافیکی کار با Hive و اجرای کوئری‌ها
- اسکرین‌شات‌هایی از مراحل مختلف انجام این مثال عملی را در زیر می‌توانید مشاهده کنید:



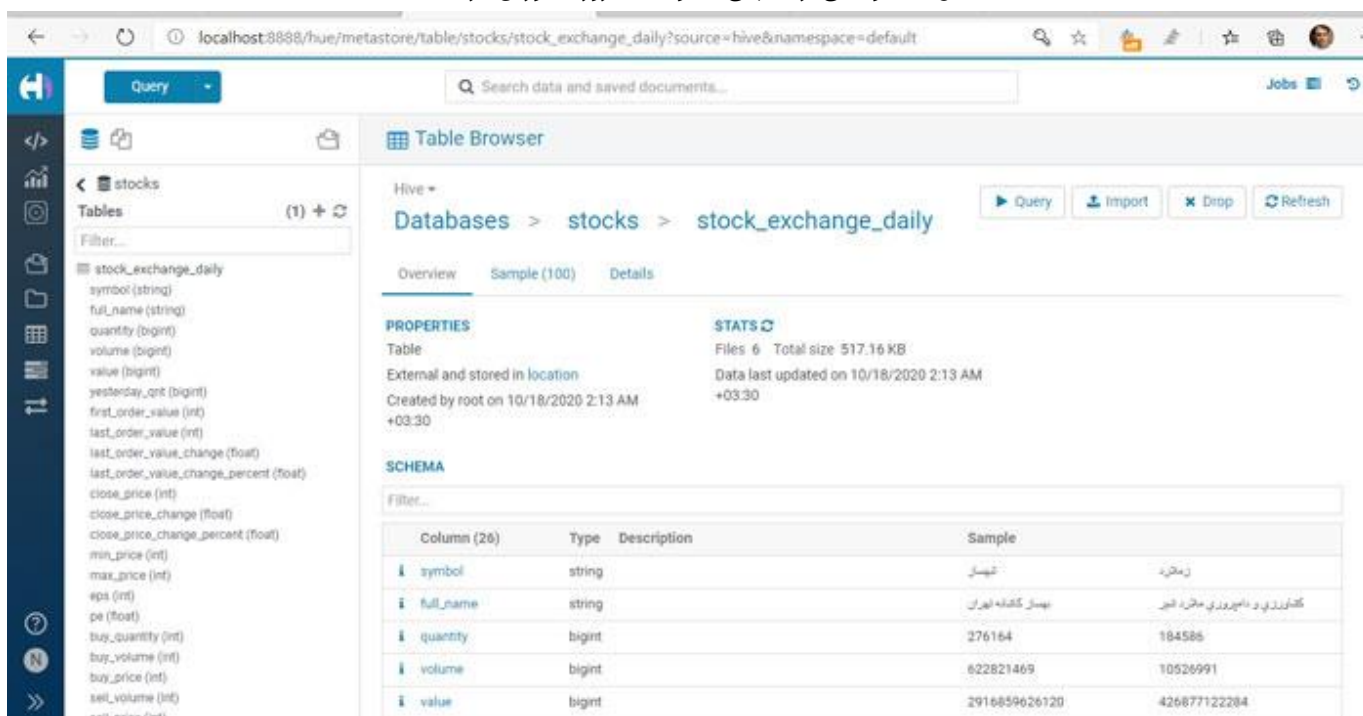
شکل ۱- نمونه‌ای از خط پردازش داده ایجاد شده در Airflow



Browse Directory

Permission	Owner	Group	Size	Last Modified	Replication	Block Size	Name
-rwxr-xr-x	airflow	supergroup	97.29 KB	10/18/2020, 12:14:16 PM	3	128 MB	daily_trades_2020-09-19.csv
-rwxr-xr-x	airflow	supergroup	93 KB	10/18/2020, 12:14:16 PM	3	128 MB	daily_trades_2020-09-20.csv
-rwxr-xr-x	airflow	supergroup	95.94 KB	10/18/2020, 12:14:41 PM	3	128 MB	daily_trades_2020-09-21.csv
-rwxr-xr-x	airflow	supergroup	94.74 KB	10/18/2020, 12:14:41 PM	3	128 MB	daily_trades_2020-09-22.csv
-rwxr-xr-x	airflow	supergroup	96.36 KB	10/18/2020, 12:15:15 PM	3	128 MB	daily_trades_2020-09-23.csv
-rwxr-xr-x	airflow	supergroup	96.03 KB	10/18/2020, 12:15:26 PM	3	128 MB	daily_trades_2020-09-26.csv
-rwxr-xr-x	airflow	supergroup	93.91 KB	10/18/2020, 12:15:26 PM	3	128 MB	daily_trades_2020-09-27.csv

شکل ۲- نمونه ای از فایل‌های دانلود شده روزانه بورس در HDFS



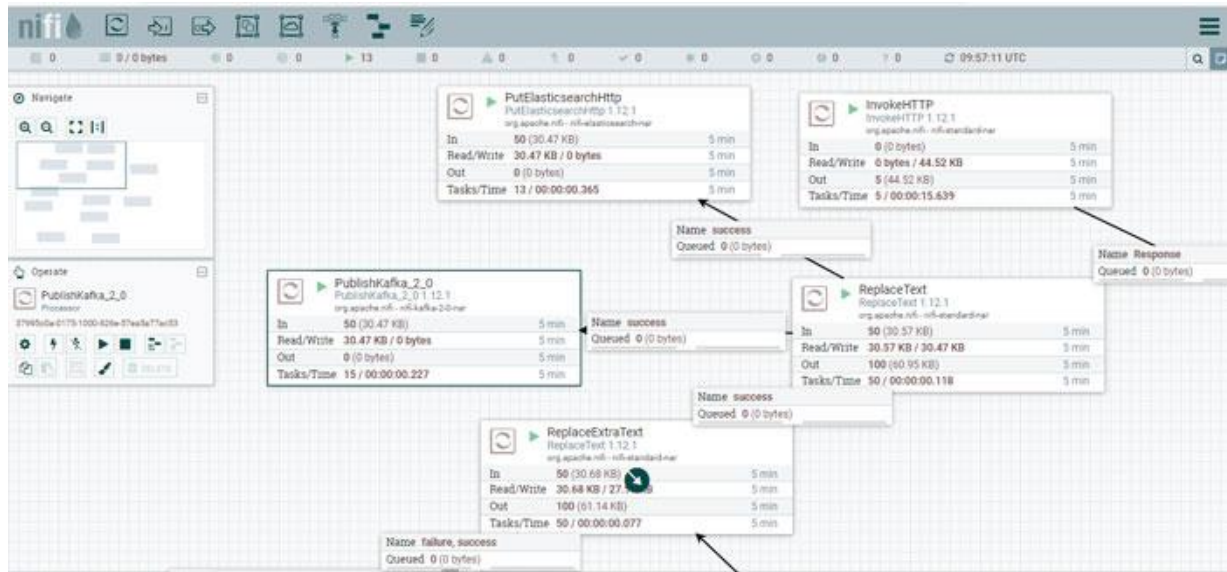
شکل ۳ - تصویری از محیط Hue که جدول ایجاد شده در Hive را نشان می‌دهد

- انجام مثال عملی پردازش توئیت‌های لحظه‌ای سایت سهامیاب (حدود یک ساعت و ربع)

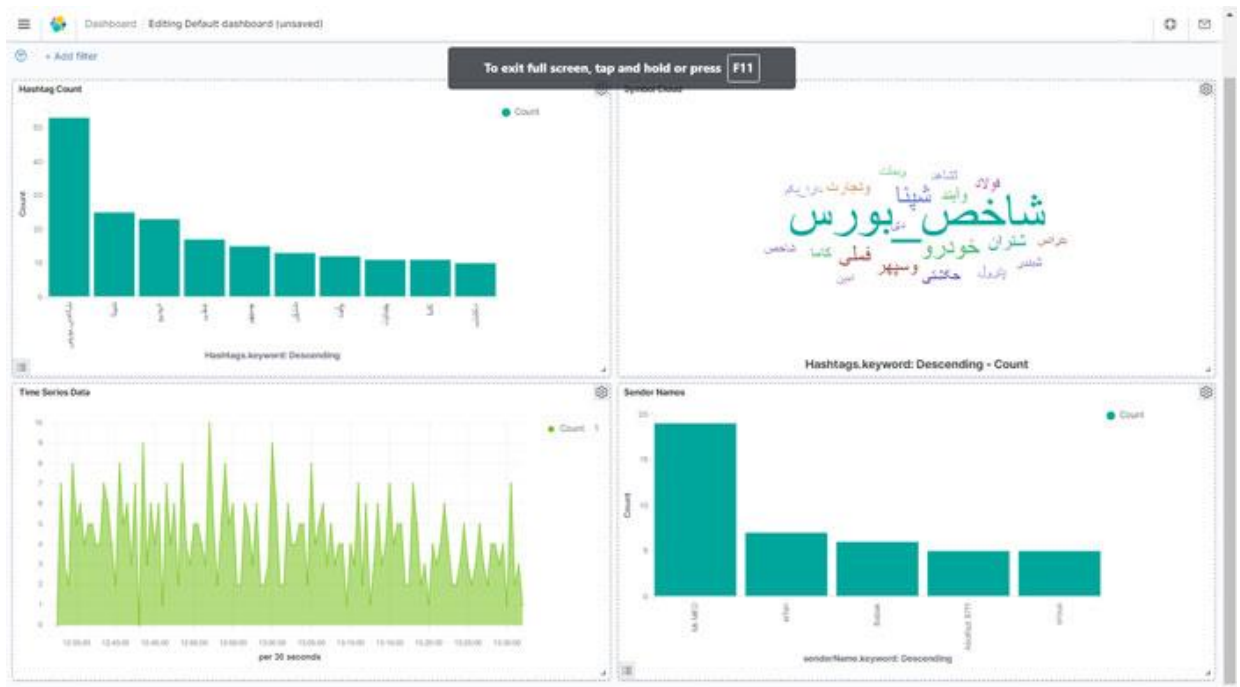
سایت سهامیاب به عنوان یکی از مراجع اصلی کاربران بورس ایران، روزانه توئیت‌ها و مطالب مختلفی را از کاربران دریافت می‌کند. در مثال دوم این کارگاه عملی، سعی می‌کنیم این توئیت‌ها را دریافت کرده، پردازش نماییم و نهایتاً آنها را ذخیره کنیم و علاوه بر ساخت داشبوردهای تحلیلی با الاستیک سرچ، به کمک اسپارک و امکانات پردازش جریان آن، اطلاعات مفیدی راجع به هر نماد بورسی در لحظه استخراج نماییم. در این مثال با ابزار زیر کار خواهیم کرد

:

- Apache Nifi به عنوان محور اصلی کار که طراحی خط پردازش داده، ما بر اساس آن انجام خواهد شد.
 - Kafka به عنوان صف توزیع توئیت‌ها در شبکه
 - Elastic Search برای ذخیره توئیت‌ها
 - Kibana برای ساخت داشبوردهای لحظه‌ای
 - Spark برای پردازش توئیت‌ها و استخراج آمار مورد نیاز که از نسخه جدید اسپارک (نسخه ۳) و رهیافت نوین Spark Structured Streaming استفاده خواهد شد.
- در زیر اسکرین‌شاتهایی از ابزار مورد استفاده و خروجی‌ها می‌توانید مشاهده کنید:



شکل ۴- تصویری از محیط آپاچی نایفای برای پردازش لحظه‌ای توئیت‌های بورس



شکل ۵- داشبورد اولیه طراحی شده برای پایش توئیت‌های لحظه‌ای بورس به کمک کیبانا

