

به نام خدا



دانشگاه تهران
پردیس دانشکده‌های فنی
دانشکده برق و کامپیوتر



درس کلان داده

تمرین شماره ۳

Spark

خرداد ماه ۱۳۹۹

مقدمه

هدف از این تمرین آشنایی با Spark به عنوان یکی از اصلی ترین فریمورک های حال حاضر کار با کلان داده در جامعه جهانی است که در بسیاری از شرکتها و کاربردها به صورت روزانه مورد استفاده قرار می گیرد.

در این تمرین ابتدا با اصول اولیه اسپارک و اجرای دستورات پایه ای آن آشنا خواهید شد و سپس با دو کتابخانه جانبی و اصلی آن یعنی Spark Graph و Spark SQL کار خواهید کرد. کار با بخش پردازش جریان در اسپارک را در پروژه نهایی این درس انجام خواهید داد.

با توجه به مشکلاتی که در نصب و راه اندازی هدوپ در دو تمرین گذشته با آن مواجه بوده اید و با هدف کار با محیط های آنلاین پردازش داده، توصیه می شود این تمرین را به کمک محیط رایانش ابری شرکت ¹ Databricks به آدرس <https://community.cloud.databricks.com> انجام دهید.

توصیه می کنیم قبل از شروع کار با این محیط، این آموزش ساده و کاربردی اسپارک را که با تمرکز بر این محیط توسعه آنلاین نوشته شده است را مطالعه کرده و دستورات آنرا به عنوان دست گرمی انجام دهید:

<https://bit.ly/SparkUT>

دیتاست های مورد نیاز هر تمرین هم همراه با تمرین آپلود شده است .

برای هر سوال، یک کتابچه پایتون (Python Notebooks) ایجاد کنید و در انتهای کار، کتابچه ها را دانلود کرده، زیپ نموده و همراه گزارش توضیحات تمرین به صورت تک نفره، آپلود نمایید.

¹ <https://databricks.com/>

سوال اول - دستورات پایه

بخش اول

در این قسمت با استفاده از تابع نگاشت-کاهش (map_reduce) تعداد لغات فایل Input.txt را شمارش کرده و نمایش دهید. همچنین گزارش کنید که هر کلمه چند بار تکرار شده است و خروجی را در یک فایل txt ذخیره کنید. در این گام تنها علائم نقطه گذاری (علامت تعجب، سوال، نقطه و ...) را حذف کنید و پیش پردازش دیگری لازم نیست.

بخش دوم

حال تمام در این گام با استفاده از تابع نگاشت-کاهش (map_reduce) تعداد تمامی کلماتی که با حرف (M) آغاز می شوند را بیابید کنید. (مستقل از کوچک و بزرگ بودن M)

بخش سوم

در این بخش نیز همانند دو گام قبلی، با استفاده از تابع نگاشت-کاهش (map_reduce) تعداد لغات 5 حرفی موجود در فایل words.txt را یافته، لغاتی که با حروف صدادار شروع می شوند را از خروجی حذف کنید و نتیجه نهایی را به صورت مرتب نمایش دهید.

بخش چهارم

به کمک مراحل قبلی، ایست واژه ها (stop words) را بیابید. کلمه ای را ایست واژه در نظر بگیرید که جزء ده درصد کلمات پرتکرار این فایل قرار بگیرد. سپس تابعی بنویسید که یک خط را گرفته، تمام حروف غیر الفبایی و ایست واژه های آنرا حذف کند. این تابع را روی تمام خطوط اعمال کرده، نتیجه را در یک فایل، ذخیره کنید.

بخش پنجم

تعداد دو کلمه ای هایی که بیشتر از یک بار در فایل اصلی (input.txt) کنار هم آمده اند را به ترتیب فرکانس، یافته و نمایش دهید. منظور از دو کلمه ای (bigram)، دو لغتی هستند که پشت سر هم به کار رفته اند.

سوال دوم - بررسی یک فایل لاگ وب سرور

فایل لاگ پیوست این تمرین با نام "Log" که مربوط به درخواست های HTTP است. با استفاده از این فایل به سوال زیر پاسخ دهید (برای این بخش از دستورات پایه اسپارک استفاده کنید):

بخش اول

چند Hostد در این لاگ فایل وجود دارد؟

بخش دوم

متوسط تعداد درخواست های روزانه برای هر میزبان منحصر به فرد (آی پی یا نام دامنه) چقدر است؟ ابتدا متوسط تعداد درخواست های هر دامنه در هر روز را به دست آورید و سپس، متوسط نهایی را برای هر دامنه یا آی پی، تعیین کنید.

بخش سوم

تعداد فایل های گیف درخواست شده در این فایل لاگ چقدر است؟

بخش چهارم

دامنه های پرتقاضا (بیش از ۳ بار) را یافته، آنها را به صورت مرتب شده نمایش دهید. آی پی ها را جزء این دامنه ها در نظر نگیرید. سپس دامنه پرتقاضا به ازای هر روز را پیدا کنید (دامنه ای با بیشترین تعداد درخواست در یک روز).

بخش پنجم

خطاهای HTTP (غیر از کد ۲۰۰، بقیه را همه خطا در نظر بگیرید.) را یافته، تعداد تکرار آنها در یک نمودار ستونی نمایش دهید.

سوال سوم - کار با دیتافریم ها / Spark SQL

با توجه به دیتاست stock.csv به سوالات زیر پاسخ دهید . این دیتاست، داده های بورس مربوط به یکی از کمپانی های بزرگ از سال ۲۰۱۲ تا ۲۰۱۷ می باشد.

برای انجام این تمرین از دو روش استفاده کنید یعنی برای هر بخش، خروجی مورد نظر را با هر کدام از دو روش زیر به صورت جداگانه به دست آورید :

1. Spark DataFrames - با توابع دیتافریم (DataFrame Operations such as min,avg,...)

2. Spark SQL - با دستورات SQL (spark.sql)

بخش اول:

یک ستون اطلاعات جدید با ستونی به نام HV ایجاد کنید که این نسبت بالاترین قیمت بر حجم سهام معامله شده برای یک روز است.

بخش دوم:

بیک بالاترین قیمت، برای چه روزی بوده است؟

بخش سوم:

میانگین ستون، Close، چه مقدار است؟

بخش چهارم:

مقدار ماکزیمم و مینیمم ستون Volum را مشخص کنید.

بخش پنجم:

چند روز ستون Close کمتر از 60 دلار بوده است؟

بخش ششم:

Pearson correlation بین ستون های High و Volum چقدر است؟

بخش هفتم:

ماکزیمم ستون High در هر سال چقدر است؟

سوال چهارم - Spark GraphX

فایل پیوست `edgs.txt` یال ها و فایل پیوست `vertex.txt` درجه های یک گراف هستند. گراف مورد نظر ما از مقالات ویکی پدیا استخراج شده اند. هر گره یک مقاله ویکی پدیا و یال از مقاله A به مقاله B نشان دهنده این است که مقاله A به مقاله B ارجاع داده است.

نکته: می توانید برای کار با گراف در اسپارک از `GraphFrames`¹ استفاده کنید.

بخش اول:

با استفاده از فایل یال ها و گره ها، این گراف را ایجاد کنید.

بخش دوم:

بیشترین درجه ورودی در این گراف چقدر است؟ بیشترین درجه خروجی (مقاله ای که احتمالا Survey بوده و شامل لینک زیادی به سایر مقالات است.) چند است؟

بخش سوم:

سایز هر کدام از `ConnectedComponent` ها چقدر است؟

بخش چهارم:

ده تا از مقالات برتر را بیابید (مقالاتی که بیشترین درجه ورودی را داشته اند).

بخش پنجم (نمره اضافی)

آیا می توانید گراف فوق را به صورت بصری نمایش دهید؟

¹ graphframes.github.io/graphframes/