

به نام خدا



دانشگاه صنعتی امیرکبیر  
( پلی تکنیک تهران )

تمرین سوم (فاز ۲) درس مبانی رایانش ابری

استاد درس: دکتر جوادی

نام : محمد حسین شهبازی

شماره دانشجویی: ۴۰۰۳۱۰۸۸

## سوال ۱

مقایسه از نظر سطوح انتزاع:

**RDD:** سطح پایینی از انتزاع را ارائه داده و مجموعه‌های از اشیا توزیع شده هستند که به طور موازی در سرتاسر کلاستر، ذخیره میشوند؛ همچنین امکان انجام عملیاتی ساده‌ای مانند filter را ارائه میدهند.

**DataFrames:** سطح بالاتری از انتزاع را ارائه داده و مشابه جداول در یک پایگاه داده رابطهای هستند که به صورت توزیع شده و موازی ذخیره میشوند؛ همچنین امکان انجام عملیاتی مانند filter ، group by ، select و ... را فراهم میسازند.

## مقایسه از نظر جنبه های عملکرد:

به دلیل سطح پایینتر انتزاع، کنترل بیشتری بر روی داده ها و عملیاتها دارند؛ همچنین عملیات بر روی RDD ها به صورت lazy evaluation انجام میشود، یعنی محاسبات تنها زمانی اجرا میشوند که اکشنها فراخوانی شوند.

*DataFrames*: از بهینه سازیهای داخلی مانند optimizer catalyst بهره مند میشوند که باعث بهبود عملکرد خودکار میشود و در بیشتر موارد، سریعتر از RDD ها هستند، زیرا بهینه سازیهای داخلی بیشتری برای کاهش زمان اجرا دارند.

## مقایسه از نظر موارد استفاده:

*RDDs*: برای عملیاتی پیچیده و سفارشی که نیاز به کنترل دقیق دارند و الگوریتمهای سنگینی مانند پردازش گراف و یادگیری ماشین، مناسب بوده و معمولاً وقتی استفاده میشوند که لازم باشد تا به داده ها به صورت مستقیم و بدون استفاده از API سطح بالاتر دسترسی پیدا نمود.

*DataFrames*: برای عملیات تحلیل داده های متداول که میتوانند به صورت SQL بیان شوند، مناسب بوده و برای بارگذاری، پاکسازی و تحلیل داده های ساختاریافته و نیمه ساختاریافته استفاده میشوند.

## سوال ۲

از نظر قابلیت پرس و جو، SQL Spark به دلیل استفاده از زبان SQL که يك زبان استاندارد براي پایگاه داده است، کوئریهای خواناتری را ارائه میدهد و به راحتی میتوان از آن براي کوئری زدن به منابع داده مختلف داده استفاده نمود؛ همچنین براي عملیاتی پیچیده تحلیل داده نیز مناسبتر میباشد؛ اما API DataFrame قابلیت‌های بیشتری براي ترکیب عملیاتی پیچیده داشته و این امکان را ارائه میدهد تا عملیاتی سفارشی تري انجام گیرند، همچنین در زبانهای برنامه نویسی مختلفی در دسترس بوده و براي برنامه نویسان، استفاده از آن راحت تر میباشد. از نظر بهینه سازی عملکرد، هردو از optimizer catalyst براي بهینه سازی کوئریها و از execution tungsten engine براي بهبود عملکرد حافظه و CPU استفاده میکنند؛ با این تفاوت که در SQL Spark، بهینه سازی به صورت خودکار انجام میشود ولي در API DataFrame، کاربران میتوانند بهینه سازیهای بیشتری را به صورت دستی انجام دهند. از نظر قابلیت استفاده از API، SQL Spark سادگی و خوانایی بیشتری داشته و چون زبان SQL استاندارد است، احتمال خطا کاهش مییابد؛ اما API DataFrame انعطاف پذیری و قابلیت ترکیب بیشتری دارد و میتوان از توابع و ماژول های مختلف، براي ایجاد کوئریهای قابل استفاده مجدد و سفارشی استفاده نمود. به طور کلی SQL Spark مناسب براي تحلیل داده و کسانی است که با SQL کار میکنند و براي کوئریهای ساده و سریع، ایده آل میباشد؛ اما API DataFrame براي توسعه دهندگانی که به عملیات های پیچیده تر نیاز دارند مناسب بوده و در پروژه هایی که نیاز به ترکیب چندین عملیات پیچیده است، استفاده میشود.

### سوال ۳

Partitioning در Spark Apache به فرآیند تقسیم داده‌ها به بخش‌های کوچکتر به نام پارتیشن‌ها اشاره دارد. هر پارتیشن یک بخش منطقی از داده است که میتواند به صورت مستقل پردازش شود؛ این پارتیشن‌ها به صورت توزیع شده در میان نودهای مختلف یک کلاستر ذخیره میشوند و پردازش موازی روی آنها انجام میشود Partitioning. باعث افزایش کارایی و سرعت پردازش، بالانس بار و توزیع متوازن job ها، بهینه‌سازی حافظه و منابع، کاهش هزینه‌های ارتباطی و انعطاف پذیری و قابلیت تحمل خطای بیشتری میشود.

## سوال ۴

### پارتیشن بندی پیش فرض (Default Partitioning)

Spark به صورت خودکار داده‌ها را بر اساس حجم و تعداد نودهای موجود در کلاستر پارتیشن بندی میکند؛ این استراتژی برای بسیاری از کاربردها کافی است، اما در برخی موارد ممکن است نیاز به تنظیمات دستی باشد. در بسیاری از موارد، این پارتیشن‌بندی عملکرد خوبی دارد، اما ممکن است در مواردی که داده‌ها به صورت نامتوازن توزیع شده‌اند یا حجم داده‌ها بسیار زیاد است، منجر به توزیع نامناسب بار شود.

### پارتیشن بندی بر پایه کلید (Key-Based Partitioning)

در این روش، داده‌ها بر اساس کلید خاصی پارتیشن بندی میشوند؛ این استراتژی معمولاً در عملیات‌های join استفاده میشود زیرا داده‌هایی که دارای کلیدهای مشابه هستند، در یک پارتیشن قرار میگیرند. این روش میتواند بهبود عملکرد قابل توجهی در عملیات‌های join و by group داشته باشد، زیرا نیاز به انتقال داده‌ها بین نودها کاهش مییابد.

## پارتیشن‌بندی سفارشی (Custom Partitioning)

پارتیشن‌بندی سفارشی این امکان را به کاربران می‌دهد تا پارتیشن‌بندی خاصی را بر اساس نیازهای خاص خود پیاده‌سازی کنند. این روش می‌تواند منجر به بهبود عملکرد در سناریوهای خاص شود، زیرا کاربران می‌توانند داده‌ها را به نحوی پارتیشن‌بندی کنند که منجر به کاهش انتقال داده‌ها و بهینه‌سازی استفاده از منابع شود.

## بخش امتیازی ۱:

مقایسه زمان محاسبه مجموع salary کمپانی ها با اندازه Large ، با caching و بدون caching به صورت زیر میباشد.:

Bonus 1:

1. Calculating sum of salary with company size 'L' without caching.

```
] rdd = salaries_df.rdd

start_time_no_cache = time.time()
dummy = rdd.filter(lambda row: row['company_size'] == 'L').aggregate(0,
                                                                    lambda acc, row: acc + row['salary'],
                                                                    lambda acc1, acc2: acc1 + acc2)
end_time_no_cache = time.time()
end_time_no_cache - start_time_no_cache
```

```
] 0.2176377773284912
```

2. Calculating sum of salary with company size 'L' with caching.

```
] rdd_cached = rdd.cache()

start_time_cache = time.time()
dummy = rdd_cached.filter(lambda row: row['company_size'] == 'L').aggregate(0,
                                                                              lambda acc, row: acc + row['salary'],
                                                                              lambda acc1, acc2: acc1 + acc2)
end_time_cache = time.time()
end_time_cache - start_time_cache
```

```
] 0.3096776008605957
```

مشاهده میشود که با استفاده از caching ، سرعت پردازش بهتری به دست آمده است زیرا در پردازش دوم، به جای اینکه داده ها مستقیماً از منبع اصلی خوانده شوند، داده ها از حافظه خوانده شده اند.



## بخش امتیازی ۲:

Narrow transformation ها شامل عملیات هایی هستند که در آنها هر پارتیشن از RDD خروجی فقط به یک پارتیشن از RDD ورودی وابسته است؛ این تغییرات نیاز به جابهجایی داده ها بین نودها ندارند و به راحتی در محل انجام میشوند.

wide transformation ها شامل عملیاتی هستند که در آنها داده ها بین نودهای مختلف جابه جا میشوند، به عبارت دیگر نیاز به shuffle میباشد؛ این عملیاتها به دلیل نیاز به بازتوزیع داده ها، باعث کاهش عملکرد میشوند زیرا نیاز به تبادل داده بین نودها افزایش می یابد. مقایسه ای بین این دو به صورت زیر میباشد:

Bonus 2:

1.Narrow transforamtion example.

```
[97]: rdd = sc.parallelize(range(1, 1000000))

start_time_narrow = time.time()
narrow_result = rdd.map(lambda x: x * 2).collect()
end_time_narrow = time.time()
end_time_narrow - start_time_narrow
```

[97]: 0.42276501655578613

2.Wide transforamtion example.

```
[98]: rdd_map = rdd.map(lambda x: (x % 100, x))

start_time_wide = time.time()
wide_result = rdd_map.groupByKey().mapValues(list).collect()
end_time_wide = time.time()
end_time_wide - start_time_wide
```

[98]: 0.953167200088501

مشاهده میشود که در transformation narrow زمان اجرا به دلیل استفاده از map و عدم نیاز به جابهجایی داده ها کوتاه است؛ ولی زمان اجرای transformation wide به دلیل استفاده از key by group و نیاز به shuffle و تبادل داده ها بین نودها، بیشتر میباشد.