

گزارش یک:

RDDs از نظر معماری و سطوح انتزاع یک سطح انتزاعی پایینتر نسبت به DataFrame ها دارند و مستقیماً روی توزیع داده ها و عملیات بر روی آنها کنترل دارند. RDD ها مجموعه ای توزیع شده از اشیاء هستند و با استفاده از توابع تحولی (Transformation) و عملیات (Action) مدیریت میشوند. کاربران میتوانند با استفاده از RDD ها کنترل دقیقی روی توزیع دادهها، تکرارپذیری و پردازش موازی داشته باشند. بهینه سازی عملکرد RDD ها به طور دستی انجام میشود. کاربران باید خودشان بهینه سازیهای لازم را برای عملیاتهای سنگین و توزیع داده ها انجام دهند. RDD ها از بهینه سازیهای خودکار اسپارک کاتالیزور استفاده نمیکند. در مورد موارد استفاده آنها داریم که زمانی که نیاز به کنترل دقیق تری روی داده ها و نحوه توزیع و پردازش آنها باشد، استفاده از RDD ها مناسب تر است. برای پردازش داده هایی که ساختار مشخصی ندارند، RDD ها گزینه مناسبی هستند. DataFrame از نظر معماری و سطوح انتزاعی بالاتر نسبت به RDD ها دارند و به صورت جداول دو بعدی سازماندهی میشود. DataFrame ها به صورت جدول با ستونهای نامگذاری شده نمایش داده میشوند و از API SQL و توابع تحولی سطح بال بهره مند هستند. DataFrame ها از بهینه ساز کاتالیزور (Optimizer Catalyst) استفاده میکنند که به طور خودکار برنامههای اجرایی بهینه را ایجاد میکند. عملیتهای بر روی DataFrame ها میتواند به صورت (Lazy) اجرا شود و بهینه سازی هایی مانند پردازش درون حافظه های و فشرده سازی دادهها را شامل شود. برای تحلیل دادههای ساختاریافته و نیمهساختاریافته که شباهت به جداول دارند، DataFrame ها انتخاب بهتری هستند. برای انجام تحلیلهای پیچیده و نیاز به بهینه سازیهای خودکار، استفاده از DataFrame ها بسیار موثر است.

گزارش دو:

۱. قابلیت های پرس و جو

Spark SQL

Spark SQL از زبان SQL استاندارد استفاده میکند که برای بسیاری از کاربران آشنا و شفاف است. قابلیت های پرس و جوی پیچیده های دارد، شامل joinها، subqueryها، aggregation، window functionها و موارد دیگر. امکان ادغام با Hive برای استفاده از جداول و متادیتاهای موجود.

API DataFrame

DataFrame API با استفاده از کد برنامه نویسی و به صورت توابع قابل استفاده است که میتواند برای توسعه دهندگان نرم افزار جذاب باشد. امکان استفاده از توابع خاص تر و پیچیده تر نسبت به SQL. برخی از عملیات پیچیده و سفارشی را میتوان به سادگی در DataFrame پیاده سازی کرد که در SQL دشوارتر است.

۲. بهینه سازی عملکرد

Spark SQL

از یک بهینه‌ساز داخلی به نام Catalyst استفاده میکند که به صورت خودکار برنامه‌های اجرایی را بهینه‌سازی میکند.

بهینه‌سازیهای مبتنی بر هزینه که بر اساس آمار و داده‌های موجود تصمیم میگیرند.

API DataFrame

API DataFrame نیز از Optimizer Catalyst بهره‌مند است، بنابراین بسیاری از بهینه‌سازیها مشترک هستند.

استفاده از UDF ها ممکن است برخی از بهینه‌سازیهای خودکار را غیرفعال کند، بنابراین عملکرد ممکن است در برخی موارد کاهش یابد.

۳. قابلیت استفاده

Spark SQL

برای کاربرانی که با SQL آشنایی دارند بسیار راحت است. یک زبان استاندارد است که میتواند در سیستمهای مختلف استفاده شود.

بسیاری از ابزارهای Intelligence Business به خوبی با SQL سازگار هستند.

API DataFrame

مناسب برای کسانی که به برنامه‌نویسی تابعی عاقله‌مند هستند. امکان ترکیب با سایر API ها و کتابخانه‌های Spark.

API DataFrame با زبانهای مختلفی مانند Python، Scala و Java سازگار است.

مثالهایی از ترجیح بر دیگری:

استفاده از SQL Spark:

- برای کاربران تجاری که به گزارشهای پیچیده نیاز دارند و با SQL آشنا هستند. استفاده از DataFrame API:

برای توسعه‌دهندگانی که به پردازشهای پیچیده و سفارشی نیاز دارند. هنگامی که نیاز به ادغام کدهای پردازشی با سایر بخشهای برنامه کاربردی است.

گزارش سه:

Partitioning یا تقسیم‌بندی داده‌ها در Spark Apache به معنای توزیع داده‌ها در بخشهای کوچکتر و جداگانه است که میتوانند به صورت موازی پردازش شوند.

در Spark Apache، یک RDD یا DataFrame به چندین بخش کوچکتر به نام پارتیشن تقسیم میشود. هر پارتیشن یک زیرمجموعه از داده‌ها را نگهداری میکند و میتواند به صورت مستقل توسط یک نود (گره) در کالستر پردازش شود.

اهمیت Partitioning در پردازش داده توزیع شده:

1. افزایش کارایی (Efficiency): با توزیع داده‌ها در بین نودهای مختلف، از منابع محاسباتی و حافظه هر نود بهینه‌تر استفاده میشود.

پارتیشن‌ها به صورت موازی پردازش میشوند که این باعث افزایش سرعت پردازش و کاهش زمان کلی اجرا میشود.

2. (Balancing Load): توزیع یکنواخت پارتیشن‌ها بین نودها باعث میشود که هیچ نودی بار کاری بیش از حد نداشته باشد.

3. افزایش انعطافپذیری (Tolerance Fault): اگر یک نود خراب شود، فقط پارتیشن‌های مربوط به آن نود نیاز به بازپردازش دارند، نه کل داده‌ها. این باعث میشود که سیستم در برابر خطاها مقاومتر باشد.

گزارش چهار:

1. پارتیشن‌بندی پیشفرض (Partitioning Default):

اگر تعداد پارتیشن‌ها مشخص نشود، اسپارک به طور خودکار تعداد پارتیشن‌ها را تعیین میکند.

مزایا: برای بسیاری از jobهای ساده مناسب است.

معایب: در برخی موارد، تعداد پارتیشن‌ها ممکن است بهینه نباشد و منجر به پارتیشن‌های بیش از حد بزرگ یا کوچک شود.

2. پارتیشن‌بندی بر اساس کلید (Partitioning Hash):

داده‌ها بر اساس مقدار هش کلیدهای مشخص شده پارتیشن‌بندی میشوند.

مزایا: توزیع یکنواخت داده‌ها در پارتیشن‌ها، مناسب برای عملیات join و aggregation.

معایب: در صورتی که داده‌ها دارای توزیع یکنواخت نباشند، ممکن است منجر به عدم تعادل در بارگذاری پارتیشن‌ها شود.

3. پارتیشن‌بندی بر اساس محدوده (Partitioning Range):

داده‌ها بر اساس محدوده مقادیر کلیدهای مشخص شده پارتیشن‌بندی میشوند.

مزایا: مناسب برای عملیات مرتب‌سازی و جستجو، توزیع یکنواخت در صورت تعریف محدوددهای مناسب.

معایب: تعیین محدوددهای مناسب نیازمند دانش پیشین درباره توزیع داده‌ها است.

4. پارتیشن‌بندی سفارشی (Partitioning Custom):

کاربران میتوانند با پیاده‌سازی یک تابع پارتیشن‌بندی سفارشی، داده‌ها را به صورت دلخواه پارتیشن‌بندی کنند.

مزایا: انعطافپذیری بالا، امکان بهینه‌سازی برای سناریوهای خاص.

معایب: نیاز به پیاده‌سازی و نگهداری پیچیده‌تر، خطای انسانی محتمل.

تاثیر پارتیشن‌بندی بر عملکرد اجرای job در اسپارک: (Balancing Load):

پارتیشن‌های یکنواخت‌تر منجر به توازن بهتر بار و کاهش زمان اجرای jobها میشود.

(Overhead Network): پارتیشن‌بندی بهینه میتواند تعداد و حجم داده‌های منتقل شده بین گره‌ها را کاهش

دهد.

(Utilization CPU and Memory): پارتیشن‌های کوچک‌تر ممکن است به حافظه کمتری نیاز داشته باشند،

اما تعداد زیاد آنها میتواند سربرار پردازشی ایجاد کند. پارتیشن‌های بزرگ‌تر ممکن است منجر به استفاده بهینه‌تر از

پردازنده‌ها شود، اما نیازمند حافظه بیشتری هستند.

مرحله shuffle: عملیات shuffle در اسپارک که داده‌ها را بین گره‌ها منتقل میکند، بسیار هزینه‌بر است.

پارتیشن‌بندی مناسب میتواند هزینه‌های مربوط به shuffle را کاهش دهد.

امتیازی:

narrow و transformations wide به نوع ارتباط و پردازش داده‌ها در **RDD** اشاره دارد. **transformations Narrow** به آن دسته از تبدیلاتی گفته میشود که در آنها هر پارتیشن از **RDD** ورودی حداکثر به یک پارتیشن از **RDD** خروجی نگاشت میشود. این تبدیلات نیازی به جابجایی داده‌ها بین نودهای مختلف ندارند و به صورت **local** پردازش میشوند. مثالی از **narrow** است **mapPartitions** و **map** , **filter** شامل **transformations**

transformations Wide به آن دسته از تبدیلاتی گفته میشود که در آنها داده‌ها باید بین پارتیشن‌ها نیازمند عملیات **shuffling** است. این تبدیلات نیازمند جابجایی داده‌ها بین نودهای جابجا شوند، که معمولاً مختلف هستند و شامل عملیاتی مثل **reduceByKey** , **groupByKey** و **join** میشود. تفاوت‌ها:

محلی بودن پردازش: **transformations Narrow**: پردازش داده‌ها به صورت محلی (**local**) انجام میشود. **transformations Wide**: نیازمند جابجایی داده‌ها بین نودهای مختلف (**shuffling**) هستند. کارایی: **transformations Narrow** معمولاً سریعتر و کارآمدتر هستند زیرا نیازی به جابجایی داده‌ها ندارند.

transformations Wide معمولاً کندتر هستند زیرا شامل عملیات **shuffling** میشوند که به انتقال داده‌ها بین نودها نیاز دارد.

تصاویر:

File Edit View Run Kernel Tabs Settings Help

+

+

+

+

+

Filter files by name

data

spark.ipynb

Name

Last Modified

data

an hour ago

spark.ipynb

seconds ago

spark.ipynb

Python 3 (pykernel)

Markdown

DataFrames

A DataFrame is two-dimensional. Columns can be of different data types. DataFrames accept many data inputs including series and other DataFrames. You can pass indexes (row labels) and columns (column labels). Indexes can be numbers, dates, or strings/tuples.

Load the data into Spark dataframe

```
[48]: schema = StructType([
    StructField("work_year", IntegerType(), True),
    StructField("experience_level", StringType(), True),
    StructField("employment_type", StringType(), True),
    StructField("job_title", StringType(), True),
    StructField("salary", IntegerType(), True),
    StructField("salary_currency", StringType(), True),
    StructField("salary_in_usd", IntegerType(), True),
    StructField("employee_residence", StringType(), True),
    StructField("remote_ratio", IntegerType(), True),
    StructField("company_location", StringType(), True),
    StructField("company_size", StringType(), True)
])

[49]: # Convert Pandas DataFrame to Spark DataFrame with defined schema
salaries_df = spark.createDataFrame(data_engineer_salary.values.tolist(), schema=schema)
# Show the DataFrame schema and some data
salaries_df.printSchema()

root
 |-- work_year: integer (nullable = true)
 |-- experience_level: string (nullable = true)
 |-- employment_type: string (nullable = true)
 |-- job_title: string (nullable = true)
 |-- salary: integer (nullable = true)
 |-- salary_currency: string (nullable = true)
 |-- salary_in_usd: integer (nullable = true)
 |-- employee_residence: string (nullable = true)
 |-- remote_ratio: integer (nullable = true)
 |-- company_location: string (nullable = true)
 |-- company_size: string (nullable = true)
```

Basic data analysis and manipulation

File Edit View Run Kernel Tabs Settings Help

+

+

+

+

+

Filter files by name

data

spark.ipynb

Name

Last Modified

data

an hour ago

spark.ipynb

6 minutes ago

spark.ipynb

Python 3 (pykernel)

Markdown

Spark Context and Spark Session

You will create the Spark Context and initialize the Spark session needed for SparkSQL and DataFrames. SparkContext is the entry point for Spark applications and contains functions to create RDDs such as `parallelize()`. SparkSession is needed for SparkSQL and DataFrame operations.

```
[43]: # Creating a spark context class
sc = SparkContext()

# # Creating a spark session
spark = SparkSession.builder.appName("pyspark-notebook").master("spark://spark-master:7077").config("spark.executor.memory", "1024m").getOrCreate()

[44]: spark

[44]: SparkSession - in-memory

SparkContext

Spark UI

Version          v3.0.0
Master           local[*]
AppName          pyspark-shell

Initialize Spark session

To work with dataframes we just need to verify that the spark session instance has been created.

[45]: if 'spark' in locals() and isinstance(spark, SparkSession):
    print("SparkSession is active and ready to use.")
else:
    print("SparkSession is not active. Please create a SparkSession.")

SparkSession is active and ready to use.
```

DataFrames

File Edit View Run Kernel Tabs Settings Help

Filter files by name

Name Last Modified

data an hour ago

spark.ipynb 5 minutes ago

Basic data analysis and manipulation

In this section, we perform basic data analysis and manipulation. We start with previewing the data and then applying some filtering and columnwise operations.

Task 1:

Show the first 5 records of the DataFrame.

```
[51]: # to-do
salaries_df.show(5)
```

work_year	experience_level	employment_type	job_title	salary	salary_currency	salary_in_usd	employee_residence	remote_ratio	company_location	company_size
2024	SE	FT	AI Engineer	202730	USD	202730	US	0	US	M
2024	SE	FT	AI Engineer	92118	USD	92118	US	0	US	M
2024	SE	FT	Data Engineer	130500	USD	130500	US	0	US	M
2024	SE	FT	Data Engineer	96000	USD	96000	US	0	US	M
2024	SE	FT	Machine Learning ...	190000	USD	190000	US	0	US	M

only showing top 5 rows

Task 2:

Select the `salary` and `job_title` columns from the DataFrame and display the first 5 rows of this columns.

```
[14]: # to-do
salaries_df.select("salary", "job_title").show(5)
```

salary	job_title
202730	AI Engineer
92118	AI Engineer
130500	Data Engineer
96000	Data Engineer
190000	Machine Learning ...

only showing top 5 rows

Task 3:

Display the first five rows of `salary` and `job_title` and `salary_currency` from the DataFrame where the salary is less than 50,000.

File Edit View Run Kernel Tabs Settings Help

Filter files by name

Name Last Modified

data an hour ago

spark.ipynb 2 minutes ago

```
[51]: # to-do
salaries_df.show(5)
```

work_year	experience_level	employment_type	job_title	salary	salary_currency	salary_in_usd	employee_residence	remote_ratio	company_location	company_size
2024	SE	FT	AI Engineer	202730	USD	202730	US	0	US	M
2024	SE	FT	AI Engineer	92118	USD	92118	US	0	US	M
2024	SE	FT	Data Engineer	130500	USD	130500	US	0	US	M
2024	SE	FT	Data Engineer	96000	USD	96000	US	0	US	M
2024	SE	FT	Machine Learning ...	190000	USD	190000	US	0	US	M

only showing top 5 rows

Task 2:

Select the `salary` and `job_title` columns from the DataFrame and display the first 5 rows of this columns.

```
[53]: # to-do
salaries_df.select("salary", "job_title").show(5)
```

salary	job_title
202730	AI Engineer
92118	AI Engineer
130500	Data Engineer
96000	Data Engineer
190000	Machine Learning ...

only showing top 5 rows

Task 3:

Display the first five rows of `salary` and `job_title` and `salary_currency` from the DataFrame where the salary is less than 50,000.

```
[15]: # to-do
filtered_data = salaries_df.filter(salaries_df.salary < 50000).select("job_title", "salary", "salary_currency").show()
```

job_title	salary	salary_currency
Data Analyst	45864	USD
Data Architect	45000	GBP
Data Analyst	30000	EUR
Data Analyst	25000	EUR
Data Specialist	45607	GBP
Data Specialist	30921	GBP
Data Analyst	45000	GBP
W. Engineer	25000	USD
Data Analyst	47500	USD
Data Analyst	46000	GBP

```
File Edit View Run Kernel Tabs Settings Help

Filter files by name

Name Last Modified
data an hour ago
spark.pyynb 2 minutes ago

Task 3:
Display the first five rows of salary and job_title and salary_currency from the DataFrame where the salary is less than 50,000.

[55]: # to-do
filtered_data = salaries_df.filter(salaries_df.salary < 50000).select("job_title", "salary", "salary_currency").show()

+-----+
| job_title[salary[salary_currency]] |
+-----+
| Data Analyst[45864 USD] |
| Data Architect[45800 GBP] |
| Data Analyst[30000 EUR] |
| Data Analyst[25800 EUR] |
| Data Specialist[45607 GBP] |
| Data Specialist[39921 GBP] |
| Data Analyst[45800 GBP] |
| ML Engineer[25800 USD] |
| Data Analyst[47500 USD] |
| Data Analyst[44800 GBP] |
| Data Scientist[40000 EUR] |
| ML Engineer[25800 USD] |
| Data Analyst[46500 USD] |
| Data Specialist[41418 GBP] |
| Data Specialist[38223 GBP] |
| Data Science[33946 GBP] |
| Data Science[31396 GBP] |
| Data Scientist[45800 USD] |
| Data Scientist[30000 USD] |
| Data Engineer[35800 EUR] |
+-----+
only showing top 20 rows

Task 4:
Multiply the salary column by 0.5 and display the salary, job_title, and the newly calculated half_salary for the first five rows.

[16]: # to-do
salaries_df = salaries_df.withColumn("half_salary", col("salary") * 0.5)
salaries_df.select("job_title", "salary", "half_salary").show(5)

+-----+
| job_title[salary][half_salary] |
+-----+
| AI Engineer[202730] 101365.0 |
| AI Engineer[92118] 46059.0 |
| Data Engineer[130500] 65250.0 |
| Data Engineer[95000] 48000.0 |
| Machine Learning ...[100000] 50000.0 |
+-----+
only showing top 5 rows
```

```
File Edit View Run Kernel Tabs Settings Help

Filter files by name

Name Last Modified
data an hour ago
spark.pyynb 3 minutes ago

Task 4:
Multiply the salary column by 0.5 and display the salary, job_title, and the newly calculated half_salary for the first five rows.

[56]: # to-do
salaries_df = salaries_df.withColumn("half_salary", col("salary") * 0.5)
salaries_df.select("job_title", "salary", "half_salary").show(5)

+-----+
| job_title[salary][half_salary] |
+-----+
| AI Engineer[202730] 101365.0 |
| AI Engineer[92118] 46059.0 |
| Data Engineer[130500] 65250.0 |
| Data Engineer[95000] 48000.0 |
| Machine Learning ...[100000] 50000.0 |
+-----+
only showing top 5 rows

Task 5:
Join these two dataframes on emp_id.

[17]: # DataFrame 1
data = [("A101", "John"), ("A102", "Peter"), ("A103", "Charlie")]
columns = ["emp_id", "emp_name"]
dataframe_1 = spark.createDataFrame(data, columns)

[18]: # DataFrame 2
data = [("A101", 1000), ("A102", 2000), ("A103", 3000)]
columns = ["emp_id", "salary"]
dataframe_2 = spark.createDataFrame(data, columns)
```

```
File Edit View Run Kernel Tabs Settings Help

Filter files by name

Name Last Modified
data an hour ago
spark.pyynb 3 minutes ago

Task 5:
Join these two dataframes on emp_id.

[58]: # DataFrame 1
data = [("A101", "John"), ("A102", "Peter"), ("A103", "Charlie")]
columns = ["emp_id", "emp_name"]
dataframe_1 = spark.createDataFrame(data, columns)

[59]: # DataFrame 2
data = [("A101", 1000), ("A102", 2000), ("A103", 3000)]
columns = ["emp_id", "salary"]
dataframe_2 = spark.createDataFrame(data, columns)

[60]: # to-do
# create a new DataFrame, "combined_df" by performing an inner join
joined_dataframe = dataframe_1.join(dataframe_2, "emp_id", "inner").show()

+-----+
| emp_id[emp_name][salary] |
+-----+
| A103 Charlie[3000] |
| A102 Peter[2000] |
| A101 John[1000] |
+-----+

Task 6:
```

```
File Edit View Run Kernel Tabs Settings Help

+ - + - C
Filter files by name
/
Name Last Modified
data an hour ago
spark.pyrb 4 minutes ago

spark.pyrb
data = [{"A101", "John"}, ("A102", "Peter"), ("A103", "Charlie")]
columns = ["emp_id", "emp_name"]
dataframe_1 = spark.createDataFrame(data, columns)

[59]: # DataFrame 2
data = [{"A101", 1000}, ("A102", 2000), ("A103", 3000)]
columns = ["emp_id", "salary"]
dataframe_2 = spark.createDataFrame(data, columns)

[61]: # to-do
# create a new DataFrame, "combined_df" by performing an inner join
joined_dataframe = dataframe_1.join(dataframe_2, "emp_id", "inner").show()

+-----+
|emp_id|emp_name|salary|
+-----+
| A101| Charlie| 3000|
| A102| Peter| 2000|
| A101| John| 1000|
+-----+

Task6:
Filter the DataFrame to include only entries with experience levels "SE" (Senior Engineer) and "MI" (Mid-Level Engineer), then calculate and display the average salary for each experience level.

[62]: # to-do
filtered_df = salaries_df.filter((salaries_df["experience_level"] == "SE") | (salaries_df["experience_level"] == "MI"))
average_salary_df = filtered_df.groupBy("experience_level").agg(avg("salary").alias("average_salary"))
average_salary_df.show()

+-----+
|experience_level| average_salary|
+-----+
| MI|157116.97672114008|
| SE|168852.49343955013|
+-----+

Task7
```

```
File Edit View Run Kernel Tabs Settings Help

+ - + - C
Filter files by name
/
Name Last Modified
data an hour ago
spark.pyrb 4 minutes ago

spark.pyrb
Filter the DataFrame to include only entries with experience levels "SE" (Senior Engineer) and "MI" (Mid-Level Engineer), then calculate and display the average salary for each experience level.

[63]: # to-do
filtered_df = salaries_df.filter((salaries_df["experience_level"] == "SE") | (salaries_df["experience_level"] == "MI"))
average_salary_df = filtered_df.groupBy("experience_level").agg(avg("salary").alias("average_salary"))
average_salary_df.show()

+-----+
|experience_level| average_salary|
+-----+
| MI|157116.97672114008|
| SE|168852.49343955013|
+-----+

Task7
Write dataframe in the HDFS and load

[64]: # to-do
salaries_df.write.csv("hdfs://hadoop-namenode:9000/user/root/output_csv2")
loaded_df = spark.read.csv("hdfs://hadoop-namenode:9000/user/root/output_csv2", schema=schema)

[65]: loaded_df.show(5)

+-----+
|work_year|experience_level|employment_type| job_title|salary|salary_currency|salary_in_usd|employee_residence|remote_ratio|company_location|company_size|
+-----+
| 2023| SE| FT|Data Architect| 85000| GBP| 104584| GB| 0| GB| M|
| 2023| SE| FT|Data Architect| 75000| GBP| 92280| GB| 0| GB| M|
| 2023| MI| FT| Data Analyst|125000| USD| 125000| US| 0| US| M|
| 2023| MI| FT| Data Analyst| 86400| USD| 86400| US| 0| US| M|
| 2023| SE| FT| Data Engineer|232760| USD| 232760| US| 100| US| M|
+-----+
only showing top 5 rows

Introduction to SparksSQL

Spark SQL is a Spark module for structured data processing. It is used to query structured data inside Spark programs, using either SQL or a familiar DataFrame API.

Create a Table View
```

```
File Edit View Run Kernel Tabs Settings Help

+ - + - C
Filter files by name
/
Name Last Modified
data an hour ago
spark.pyrb 2 minutes ago

spark.pyrb
example we create a temporary view using the createTempView() function. Once we have a table view, we can run queries similar to querying a SQL table.

[67]: salaries_df.createTempView("salaries")

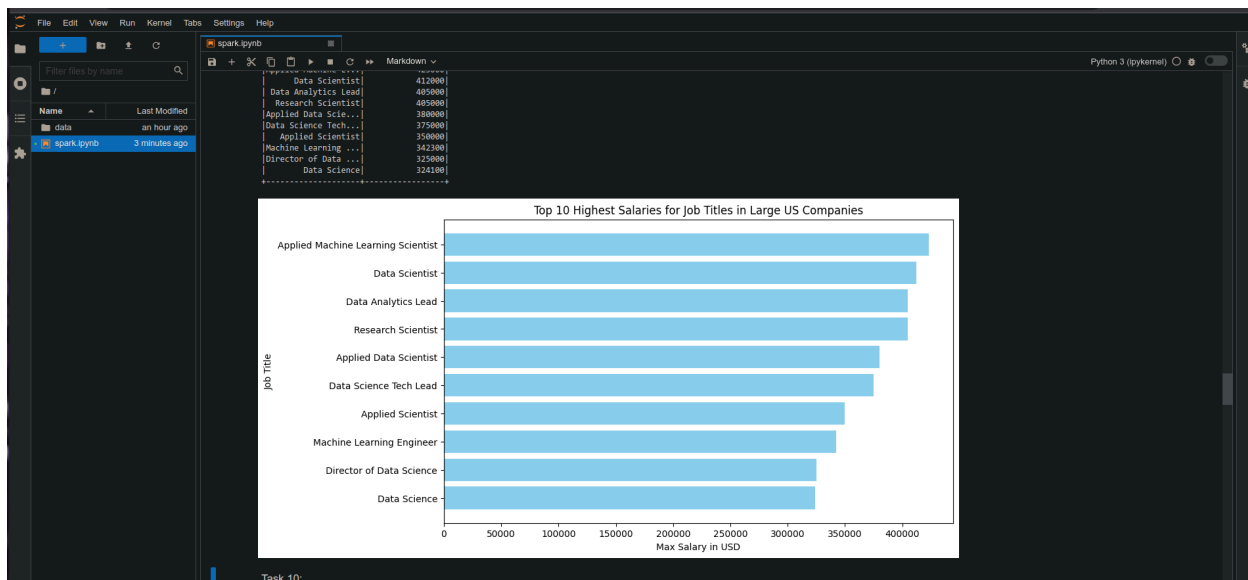
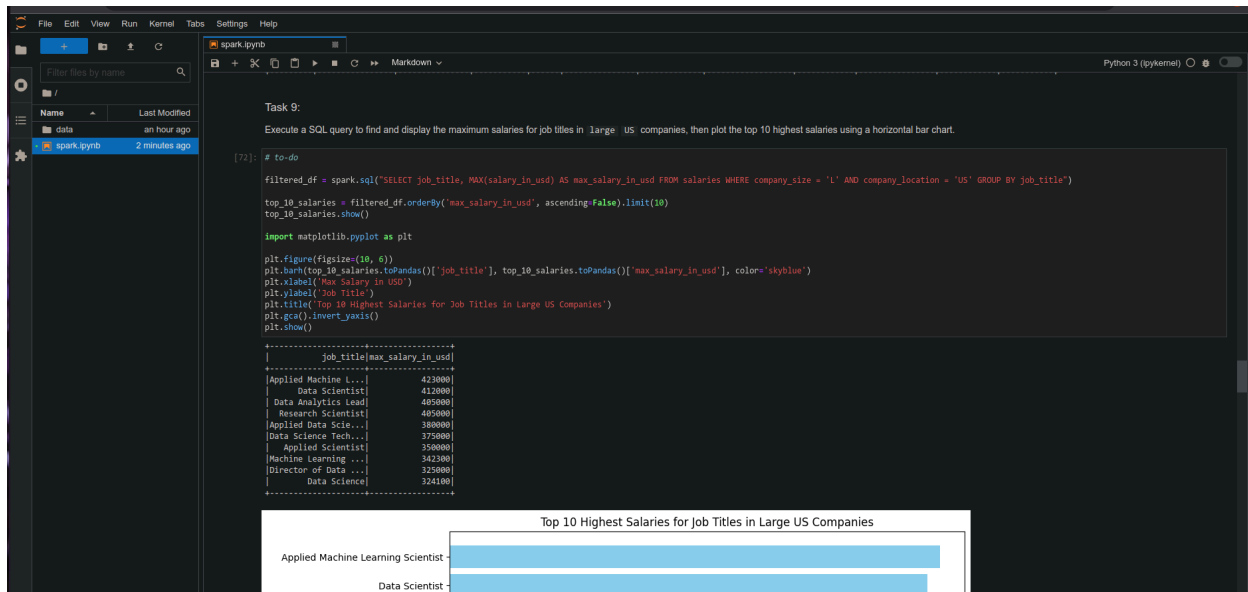
Running SQL queries and aggregating data

Task 8:
Display the first five records of AI Engineers with salaries over 10,000 with Spark SQL.

[71]: # to-do
query = """
SELECT *
FROM salaries
WHERE job_title = 'AI Engineer' AND salary > 10000
LIMIT 5
"""
result = spark.sql(query)
result.show()

+-----+
|work_year|experience_level|employment_type| job_title|salary|salary_currency|salary_in_usd|employee_residence|remote_ratio|company_location|company_size|half_salary|
+-----+
| 2024| SE| FT|AI Engineer|202730| USD| 202730| US| 0| US| M| 101365.0|
| 2024| SE| FT|AI Engineer| 92118| USD| 92118| US| 0| US| M| 46059.0|
| 2024| MI| FT|AI Engineer|150000| USD| 150000| US| 100| US| M| 75000.0|
| 2024| MI| FT|AI Engineer| 90200| USD| 90200| US| 100| US| M| 45100.0|
| 2024| SE| FT|AI Engineer|236872| USD| 236872| US| 0| US| M| 118436.0|
+-----+

Task 9:
Execute a SQL query to find and display the maximum salaries for job titles in large US companies, then plot the top 10 highest salaries using a horizontal bar chart.
```




```
File Edit View Run Kernel Tabs Settings Help

spark.pytb
total Salary given by Large Companies: 528/40143

Task 14:
Filter and display the top five highest-paying fully remote job titles based on salary.

[84]: # to do
fully_remote_rdd = rdd.filter(lambda x: x['remote_ratio'] == 100)
# Sort the RDD based on salary in descending order
sorted_rdd = fully_remote_rdd.sortBy(lambda x: x['salary'], ascending=False)
# Display the top five highest-paying job titles and salaries
top_five = sorted_rdd.map(lambda x: (x['job_title'], x['salary'])).take(5)
# Print the results
for title, salary in top_five:
    print(f"Job Title: {title}, Salary: {salary}")

Job Title: Data Scientist, Salary: 30400000
Job Title: Data Scientist, Salary: 6600000
Job Title: Data Engineer, Salary: 4450000
Job Title: Data Scientist, Salary: 4200000
Job Title: Data Scientist, Salary: 4000000
```

```
File Edit View Run Kernel Tabs Settings Help

spark.pytb
Implement a Pandas UDF in PySpark named categorize_salary that accepts a pandas Series containing salaries and categorizes each salary into "Low" (less than 50,000), "Medium" (50,000 to 149,999), or "High" (150,000 and above) by returning a new Series with these categories.

[86]: # to do
pandas.udf(StringType())
def categorize_salary(salaries):
    categories = []
    for salary in salaries:
        if salary < 50000:
            categories.append("Low")
        elif 50000 <= salary < 149999:
            categories.append("Medium")
        else:
            categories.append("High")
    return pd.Series(categories)

# Register the UDF
spark.udf.register("categorize_salary", categorize_salary)
# Apply the UDF to categorize salaries
categorized_salaries_df = spark.sql("SELECT *, categorize_salary(salary) AS salary_category FROM salaries")
categorized_salaries_df.show()

work_year|experience_level|employment_type|job_title|salary|salary_currency|salary_in_usd|employee_residence|remote_ratio|company_location|company_size|half_salary|salary_category|
-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| 2024|SE|FT|AI Engineer|202730|USD|202730|US|0|US|M|101365.0|High|
| 2024|SE|FT|AI Engineer|92118|USD|92118|US|0|US|M|46059.0|Medium|
| 2024|SE|FT|Data Engineer|138500|USD|138500|US|0|US|M|69250.0|Medium|
| 2024|SE|FT|Data Engineer|96000|USD|96000|US|0|US|M|48000.0|Medium|
| 2024|SE|FT|Machine Learning ...|190000|USD|190000|US|0|US|M|95000.0|High|
| 2024|SE|FT|Machine Learning ...|160000|USD|160000|US|0|US|M|80000.0|High|
| 2024|MI|FT|ML Engineer|400000|USD|400000|US|0|US|M|200000.0|High|
| 2024|MI|FT|ML Engineer|65000|USD|65000|US|0|US|M|32500.0|Medium|
| 2024|EN|FT|Data Analyst|101520|USD|101520|US|0|US|M|50760.0|Medium|
| 2024|EN|FT|Data Analyst|45864|USD|45864|US|0|US|M|22932.0|Low|
| 2024|SE|FT|Data Engineer|172469|USD|172469|US|0|US|M|86234.5|High|
| 2024|SE|FT|Data Engineer|114945|USD|114945|US|0|US|M|57472.5|Medium|
| 2024|EX|FT|MLP Engineer|200000|USD|200000|US|0|US|M|100000.0|High|
| 2024|EX|FT|MLP Engineer|150000|USD|150000|US|0|US|M|75000.0|High|
| 2024|MI|FT|Data Scientist|156450|USD|156450|US|100|US|M|78225.0|High|
| 2024|MI|FT|Data Scientist|119200|USD|119200|US|100|US|M|59600.0|Medium|
| 2024|SE|FT|Data Analyst|170000|USD|170000|US|0|US|M|85000.0|High|
| 2024|SE|FT|Data Analyst|130000|USD|130000|US|0|US|M|65000.0|Medium|
| 2024|SE|FT|Applied Scientist|222200|USD|222200|US|0|US|L|111100.0|High|
| 2024|SE|FT|Applied Scientist|136000|USD|136000|US|0|US|L|68000.0|Medium|
-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
only showing top 20 rows
```

```
File Edit View Run Kernel Tabs Settings Help

spark.pytb
[ 2024|SE|FT|Data Analyst|130000|USD|130000|US|0|US|M|65000.0|Medium|
| 2024|SE|FT|Applied Scientist|222200|USD|222200|US|0|US|L|111100.0|High|
| 2024|SE|FT|Applied Scientist|136000|USD|136000|US|0|US|L|68000.0|Medium|
-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
only showing top 20 rows

Task 16:
Integrate the categorize_salary UDF into a PySpark DataFrame by adding a new column salary_category which classifies each salary. Subsequently, display the original salary and its corresponding category to verify the successful application of the UDF.

[88]: # to do
categorized_salaries_df.select("salary", "salary_category").show()

[salary|salary_category]
-----|-----|
| 202730|High|
| 92118|Medium|
| 138500|Medium|
| 96000|Medium|
| 190000|High|
| 160000|High|
| 400000|High|
| 65000|Medium|
| 101520|Medium|
| 45864|Low|
| 172469|High|
| 114945|Medium|
| 200000|High|
| 150000|High|
| 156450|High|
| 119200|Medium|
| 170000|High|
| 130000|Medium|
| 222200|High|
| 136000|Medium|
-----|-----|
only showing top 20 rows

[49]: spark.stop()
```