

Exam 3 (MW) - Results

X

Attempt 1 of 1

Written Dec 3, 2025 9:59 AM - Dec 3, 2025 11:29 AM

Attempt Score 56 / 100 - 56 %

Overall Grade (Highest Attempt) 56 / 100 - 56 %

Multiple Choices, 40 points

Question 1

2 / 2 points

Which of the following is NOT among the supported parameters in REST API of Firebase realtime database?

- orderBy
- limitToFirst
- startAt
- groupBy

Question 2

0 / 2 points

Which of the following data formats is most structured?

- JSON
- XML
- Relation
- Text

Question 3

0 / 2 points

Which of the following MongoDB operators is used to join two MongoDB collections?

- \$group
- \$match
- \$lookup
- \$unwind

Question 4

0 / 2 points

Which of the following is most likely to be the page size of SSD?

- 4B
- 4KB
- 4MB
- 4GB

Question 5

0 / 2 points

Which operations in SSD are conducted page by page?

- Read only
- Read and write
- Write only
- Erase only

Question 6

3 / 3 points

Which of the following is the correct order of tasks being executed in Hadoop MapReduce?

- Map, shuffle, reduce
- Map, reduce, shuffle
- Shuffle, map, reduce
- Reduce, map, shuffle

Question 7

3 / 3 points

Consider joining two relations R and S with the amount of memory M (pages). Sizes (i.e., the number of blocks) of R and S are denoted as B(R) and B(S) respectively. Suppose both B(R) and B(S) is larger than M. Which of the following statements about B(R) and B(S) is correct?

- Nested-loop join algorithm always has higher cost than sort-merge join algorithm.
- Nested-loop join algorithm always has higher cost than partitioned-hash join algorithm.
- Sort-merge join algorithm always has the same cost as the partitioned hash join algorithm.
- If B(R) < B(S), then placing R in the outer loop of the nested-loop join will incur a lower cost.

Question 8

0 / 3 points

Which of the following transformations in Spark involves shuffling?

- map
- filter
- flatMap
- reduceByKey

Question 9

3 / 3 points

Suppose there are some empty author elements in a book XML document. Which of the following XPath expressions correctly finds such elements?

- //author[node()]
- //author[*]
- //author[(not(node()))]
- //author[(not(text()))]

Question 10

3 / 3 points

Consider an XML element about book: <book isbn=1234> <author>John</author> <author>123</author> <isbn>5678</isbn> <title>DataScience</title> </book> What syntax error does the above book element have?

- It has isbn as both attribute and sub-element.
- The value of the isbn attribute (1234) needs to be quoted.
- It can not have multiple author subelements.
- None of above.

Question 11

0 / 3 points

Suppose the actual bandwidth of reading 100 and 200 random blocks of data from a hard drive is B1 and B2 respectively. Which of the following statements about the relationship between B1 and B2 is most accurate?

- B1 < B2
- B1 = B2
- B1 = B2
- B2 = 2*B1

Question 12

3 / 3 points

How many HDFS blocks are needed to store 1GB of data?

- 2
- 4
- 8
- 10

Question 13

3 / 3 points

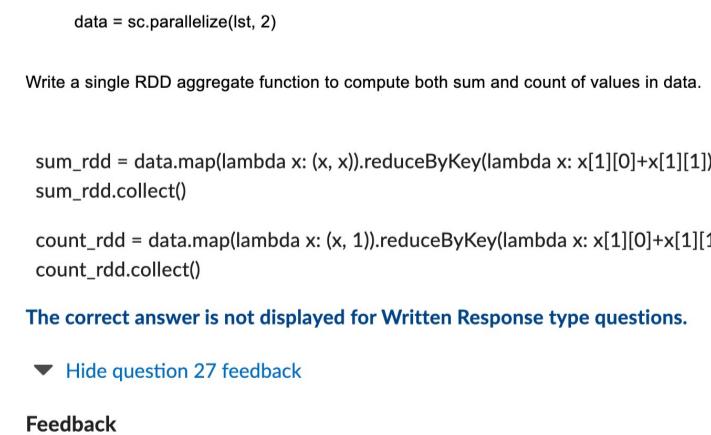
Consider two relations Product int primary key, title varchar(20); SalesId int primary key, qty float, foreign key(productId) references Product(id); Suppose Product has 100 rows and Sales has 1000 rows. How many rows will the following query produce: Select * from Product natural join Sales;

- 100
- 1000
- 900
- 1100

Question 14

3 / 3 points

The following diagram you have seen in class illustrates the shuffling process in Hadoop MapReduce. Recall that flowcharts in boxes marked as map and reduce tasks describe the steps involved in the tasks. How many map tasks and reduce tasks does the MapReduce job illustrated by the diagram have?



- 3 map and 4 reduce tasks
- 4 map and 3 reduce tasks
- 3 map and 3 reduce tasks
- 4 map and 4 reduce tasks

Question 15

3 / 3 points

Continue from the above question. Which tasks will be merging runs?

- Map task only
- Reduce task only
- Both map and reduce tasks
- Neither map nor reduce tasks

Spark DataFrame, 15 points

Consider the countrylanguage data you have seen in the homework. Suppose the data is now stored in a CSV file, countrylanguage.csv. Every row of the file has values for CountryCode, Language, IsOfficial (either 'T' or 'F'), and Percentage (assumed to be an Integer).

For example,

CAN,English,T,60

CAN,French,T,23

USA,English,T,86

USA,French,F,5

MEX,Spanish,T,92

Further suppose the country data is stored in country.csv and it has only two columns: Code and Name. For example,

CAN,Canda

USA,United States

MEX,Mexico

Suppose a Spark DataFrame called "df" has been created for you:

df = spark.read.csv("countrylanguage.csv",...)

country = spark.read.csv("country.csv",...)

The DataFrame also has four columns with data and format corresponding to that in the CSV file, that is, the first attributes are strings and the last attribute is an integer.

Also the following import statement has been executed.

import pyspark.sql.functions as fc

Using the df (and fc), write a Spark DataFrame script for each of the following SQL queries. Note your script should not return extra information that SQL queries are not asking for, and should return the attributes with names stated in the SQL queries.

Question 16

3 / 6 points

Select Language, max(Percentage) as max_p

From countrylanguage

Where IsOfficial = "T"

Having count(*) >= 10;

```
result = df.filter(fc.col("IsOfficial") == "T")
            .groupBy(fc.col("Language"))
            .filter(fc.col("count(*)") >= 10)
            .select(
                "Language",
                fc.max(fc.col("Percentage")).alias("max_p")
            )
            .show()
```

The correct answer is not displayed for Written Response type questions.

▼ Hide question 16 feedback

Feedback

- 1 incorrect usage of groupby
- 1 missing aggregation
- 1 missing count in aggregation

Question 17

4 / 5 points

Select countryName, Language

From country natural join countrylanguage

Where IsOfficial = "T";

```
joined = country.join(countrylanguage, "CountryCode") ==> cl[("Code", "inner")]
result = joined.filter(fc.col("IsOfficial") == "T")
            .select(cl["countryName"], cl["Language"])
result.show()
```

The correct answer is not displayed for Written Response type questions.

▼ Hide question 17 feedback

Feedback

- 1 incorrect join, wrong column names for both dataframes

Question 18

3 / 4 points

(select Language

from countrylanguage

where IsOfficial = "F" and CountryCode = "CAN")

except

(select Language

from countrylanguage

where IsOfficial = "F" and CountryCode = "USA");

```
nonofficial_lang_in_can = \
    df.filter(fc.col("IsOfficial") == "F") & (fc.col("CountryCode") == "CAN") \
    .select("Language")
```

```
nonofficial_lang_in_usa = \
    df.filter(fc.col("IsOfficial") == "F") & (fc.col("CountryCode") == "USA") \
    .select("Language")
```

The correct answer is not displayed for Written Response type questions.

▼ Hide question 18 feedback

Feedback

- 1 missing subtraction

Query execution, 15 points

Now suppose the country and countrylanguage data are stored in an RDBMS server. Consider the following SQL query:

Select *

From country join countrylanguage

on country.Code = countrylanguage.CountryCode

Where IsOfficial = "T";

joined = country.join(countrylanguage, "CountryCode") ==> cl[("Code", "inner")]

result = joined.filter(fc.col("IsOfficial") == "T")

result.show()

The correct answer is not displayed for Written Response type questions.

▼ Hide question 19 feedback

Feedback

- 1 missing subtraction

Question 20

2 / 2 points

Select * from country where id = 100

join countrylanguage on countrylanguage.id = country.id

having count(*) >= 100;

```
def count(*):
    if val == 100:
        output = []
    else:
        output.append(val)
    return output
```

if condition(e1, e2):
 joined.append(count(e1 + e2)) # here '+' means concatenation

return joined

The correct answer is not displayed for Written Response type questions.