

DP-203^{Q&As}

Data Engineering on Microsoft Azure

Pass Microsoft DP-203 Exam with 100% Guarantee

Free Download Real Questions & Answers PDF and VCE file from:

https://www.pass2lead.com/dp-203.html

100% Passing Guarantee 100% Money Back Assurance

Following Questions and Answers are all new published by Microsoft
Official Exam Center

- Instant Download After Purchase
- 100% Money Back Guarantee
- 365 Days Free Update
- 800,000+ Satisfied Customers





QUESTION 1

You are designing a fact table named FactPurchase in an Azure Synapse Analytics dedicated SQL pool. The table contains purchases from suppliers for a retail store. FactPurchase will contain the following columns.

Name	Data type	Nullable
PurchaseKey	Bigint	No
DateKey	Int	No
SupplierKey	Int	No
StockItemKey	Int	No
PurchaseOrderID	Int	Yes
OrderedQuantity	Int	No
OrderedOuters	Int	No
ReceivedOuters	Int	No
Package	Nvarchar(50)	No
IsOrderFinalized	Bit	No
LineageKey	Int	No

FactPurchase will have 1 million rows of data added daily and will contain three years of data. Transact-SQL queries similar to the following query will be executed daily.

SELECT

SupplierKey, StockItemKey, IsOrderFinalized, COUNT(*)FROM FactPurchase WHERE DateKey >= 20210101

AND DateKey 1 unique values while others may end with zero values.

2.

Does not have NULLs, or has only a few NULLs.

3.

Is not a date column. Incorrect Answers:

C: Round-robin tables are useful for improving loading speed.

Reference: https://docs.microsoft.com/en-us/azure/synapse-analytics/sql-data-warehouse/sql-data-warehouse-tables-distribute

QUESTION 2

DRAG DROP

2024 Latest pass2lead DP-203 PDF and VCE dumps Download

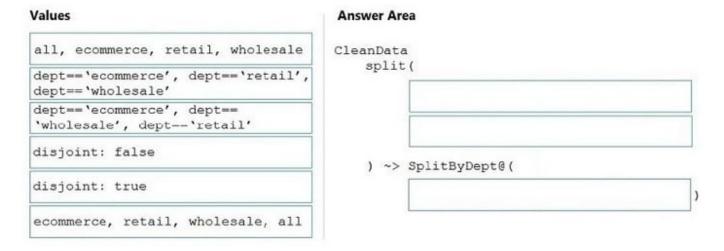
You need to create an Azure Data Factory pipeline to process data for the following three departments at your company: Ecommerce, retail, and wholesale. The solution must ensure that data can also be processed for the entire company.

How should you complete the Data Factory data flow script? To answer, drag the appropriate values to the correct targets. Each value may be used once, more than once, or not at all. You may need to drag the split bar between panes or

scroll to view content.

NOTE: Each correct selection is worth one point.

Select and Place:



Correct Answer:

```
Answer Area

all, ecommerce, retail, wholesale

cleanData
split(

dept=='ecommerce', dept=='retail',
dept=='wholesale'

disjoint: false

) ~> SplitByDept@(

disjoint: true

disjoint: true
```

The conditional split transformation routes data rows to different streams based on matching conditions. The conditional split transformation is similar to a CASE decision structure in a programming language. The transformation evaluates expressions, and based on the results, directs the data row to the specified stream.

2024 Latest pass2lead DP-203 PDF and VCE dumps Download

```
CleanData
split(

dept=='ecommerce', dept=='retail',
dept=='wholesale'

disjoint: false

) ~> SplitByDept@( ecommerce, retail, wholesale, all )
```

Box 1: dept==\\'ecommerce\\', dept==\\'wholesale\\' First we put the condition. The order must match the stream labeling we define in Box 3.

Syntax:

split(

disjoint: {true | false}
) ~> @(stream1, stream2, ...,)

Box 2: discount : false

disjoint is false because the data goes to the first matching condition. All remaining rows matching the third condition go to output stream all.

Box 3: ecommerce, retail, wholesale, all

Label the streams

QUESTION 3

You have an Azure subscription that contains an Azure Synapse Analytics dedicated SQL pool named Pool1. Pool1 receives new data once every 24 hours. You have the following function.



```
create function dbo.udfFtoC(F decimal)
return decimal
as
begin
return (F - 32) * 5.0 / 9
end
```

You have the following query.

```
select avg_date, sensorid, avg_f, dbo.udfFtoC(avg_temperature) as avg_c from SensorTemps where avg_date = @parameter
```

The query is executed once every 15 minutes and the @parameter value is set to the current date.

You need to minimize the time it takes for the query to return results.

Which two actions should you perform? Each correct answer presents part of the solution.

NOTE: Each correct selection is worth one point.

- A. Create an index on the avg_f column.
- B. Convert the avg_c column into a calculated column.
- C. Create an index on the sensorid column.
- D. Enable result set caching.
- E. Change the table distribution to replicate.

Correct Answer: BD

D: When result set caching is enabled, dedicated SQL pool automatically caches query results in the user database for repetitive use. This allows subsequent query executions to get results directly from the persisted cache so recomputation is not needed. Result set caching improves query performance and reduces compute resource usage. In addition, queries using cached results set do not use any concurrency slots and thus do not count against existing concurrency limits.

Incorrect:

Not A, not C: No joins so index not helpful.

Not E: What is a replicated table?

A replicated table has a full copy of the table accessible on each Compute node. Replicating a table removes the need to transfer data among Compute nodes before a join or aggregation. Since the table has multiple copies, replicated tables

work best when the table size is less than 2 GB compressed. 2 GB is not a hard limit. If the data is static and does not

2024 Latest pass2lead DP-203 PDF and VCE dumps Download

change, you can replicate larger tables.

Reference: https://docs.microsoft.com/en-us/azure/synapse-analytics/sql-data-warehouse/performance-tuning-result-set-caching

QUESTION 4

DRAG DROP

You are responsible for providing access to an Azure Data Lake Storage Gen2 account.

Your user account has contributor access to the storage account, and you have the application ID and access key.

You plan to use PolyBase to load data into an enterprise data warehouse in Azure Synapse Analytics.

You need to configure PolyBase to connect the data warehouse to storage account.

Which three components should you create in sequence? To answer, move the appropriate components from the list of components to the answer area and arrange them in the correct order.

Select and Place:

Answer Area

Correct Answer:

a database encryption key an asymmetric key an external data source an external file format

Step 1: a database scoped credential

To access your Data Lake Storage account, you will need to create a Database Master Key to encrypt your credential secret used in the next step. You then create a database scoped credential.

Step 2: an external data source Create the external data source. Use the CREATE EXTERNAL DATA SOURCE command to store the location of the data. Provide the credential created in the previous step.

Step 3: an external file format Configure data format: To import the data from Data Lake Storage, you need to specify the External File Format. This object defines how the files are written in

Data Lake Storage.

References:

https://docs.microsoft.com/en-us/azure/sql-data-warehouse/sql-data-warehouse-load-from-azure-data-lake-store

QUESTION 5

HOTSPOT

The following code segment is used to create an Azure Databricks cluster.



```
{
    "num workers": null,
    "autoscale": {
         "min workers": 2,
         "max workers": 8
    },
    "cluster name": "MyCluster",
    "spark version": "latest-stable-scala2.11",
    "spark conf": {
         "spark.databricks.cluster.profile": "serverless",
         "spark.databricks.repl.allowedLanguages": "sql,python,r"
    },
    "node type id": "Standard DS13 v2",
    "ssh_public_keys": [],
    "custom tags": {
         "ResourceClass": "Serverless"
    "spark env vars": {
         "PYSPARK PYTHON": "/databricks/python3/bin/python3"
    "autotermination minutes": 90,
    "enable elastic disk": true,
    "init scripts": []
}
```

For each of the following statements, select Yes if the statement is true. Otherwise, select No.

NOTE: Each correct selection is worth one point.

Hot Area:

Answer Area

Statements		No
The Databricks cluster supports multiple concurrent users.	0	0
The Databricks cluster minimizes costs when running scheduled jobs that execute notebooks.	0	0
The Databricks cluster supports the creation of a Delta Lake table.	0	0

2024 Latest pass2lead DP-203 PDF and VCE dumps Download

Correct Answer:

Answer Area

Statements		No
The Databricks cluster supports multiple concurrent users.	0	0
The Databricks cluster minimizes costs when running scheduled jobs that execute notebooks.	0	0
The Databricks cluster supports the creation of a Delta Lake table.	0	0

Box 1: Yes

A cluster mode of 'High Concurrency' is selected, unlike all the others which are 'Standard'. This results in a worker type of Standard DS13 v2.

Box 2: No

When you run a job on a new cluster, the job is treated as a data engineering (job) workload subject to the job workload pricing. When you run a job on an existing cluster, the job is treated as a data analytics (all-purpose) workload subject to

all-purpose workload pricing.

Box 3: Yes

Delta Lake on Databricks allows you to configure Delta Lake based on your workload patterns.

Reference:

https://adatis.co.uk/databricks-cluster-sizing/

https://docs.microsoft.com/en-us/azure/databricks/jobs

https://docs.databricks.com/administration-guide/capacity-planning/cmbp.html

https://docs.databricks.com/delta/index.html

DP-203 PDF Dumps

DP-203 VCE Dumps

DP-203 Practice Test