

Vce DSA-C03 Format | DSA-C03 Study Group



2026 Latest ExamsLabs DSA-C03 PDF Dumps and DSA-C03 Exam Engine Free Share: https://drive.google.com/open?id=17kD6agegdVx4KqkPkQUWJ_f-X62tgmbR

As we all know, the influence of DSA-C03 exam guides even have been extended to all professions and trades in recent years. Passing the DSA-C03 exam is not only for obtaining a paper certification, but also for a proof of your ability. Most people regard Snowflake certification as a threshold in this industry, therefore, for your convenience, we are fully equipped with a professional team with specialized experts to study and design the most applicable DSA-C03 Exam prepare. We have organized a team to research and DSA-C03 study question patterns pointing towards various learners.

Once you compare our DSA-C03 study materials with the annual real exam questions, you will find that our DSA-C03 exam questions are highly similar to the real exam questions. We have strong strengths to assist you to pass the exam. All in all, we hope that you are brave enough to challenge yourself. Our DSA-C03 learning prep will live up to your expectations. It will be your great loss to miss our DSA-C03 practice engine.

>> Vce DSA-C03 Format <<

Snowflake - DSA-C03 - SnowPro Advanced: Data Scientist Certification Exam –Valid Vce Format

The Snowflake DSA-C03 certification will further demonstrate your expertise in your profession and remove any room for ambiguity on the hiring committee's part. People need to increase their level by getting the Snowflake DSA-C03 Certification. You can choose flexible timings for the learning Snowflake DSA-C03 exam questions online and practice with Snowflake DSA-C03 exam dumps any time.

Snowflake SnowPro Advanced: Data Scientist Certification Exam Sample Questions (Q150-Q155):

NEW QUESTION # 150

You are using Snowpark Python to process a large dataset of website user activity logs stored in a Snowflake table named 'WEB ACTIVITY'. The table contains columns such as 'USER ID', 'TIMESTAMP', 'PAGE URL', 'BROWSER', and 'IP ADDRESS'.

You need to remove irrelevant data to improve model performance. Which of the following actions, either alone or in combination, would be the MOST effective for removing irrelevant data for a model predicting user conversion rates, and which Snowpark Python code snippets demonstrate these actions? Assume that conversion depends on page interaction and a model will only leverage session id and session duration.

- Remove rows where `PAGE_URL` contains common bot patterns like '/robots.txt'. Directly drop the `IP_ADDRESS` column. Example:
 

```
df = df.filter(~col('PAGE_URL').contains('/robots.txt'))
df = df.drop('IP_ADDRESS')
```
- Remove all rows where `BROWSER` is 'Internet Explorer'. Directly drop the `BROWSER` column. Example:


```
df = df.filter(col('BROWSER') != 'Internet Explorer')
df = df.drop('BROWSER')
```
- Group activities by `USER_ID` and `SESSION_ID` (assuming you can determine session IDs), calculate session duration, and keep only those sessions with duration greater than 5 seconds. Drop `PAGE_URL`, `BROWSER` and `IP_ADDRESS`. Example assuming session IDs are defined elsewhere:


```
from snowflake.snowpark.functions import datediff, to_timestamp, lit
session_df = df.groupby(['USER_ID', 'SESSION_ID']).agg(min(col('TIMESTAMP')).alias('SESSION_START'), max(col('TIMESTAMP')).alias('SESSION_END'))
session_df = session_df.withColumn('SESSION_DURATION', datediff(lit('second'), col('SESSION_START'), col('SESSION_END')))
filtered_df = df.join(session_df.filter(col('SESSION_DURATION') > 5), ['USER_ID', 'SESSION_ID'])
filtered_df = filtered_df.drop('PAGE_URL', 'BROWSER', 'IP_ADDRESS')
```
- Keep only the `USER_ID` and count the number of page views. Drop all other columns. Example:


```
from snowflake.snowpark.functions import count
df = df.groupby('USER_ID').agg(count(' ').alias('PAGE_VIEW_COUNT'))
```
- Randomly sample 10% of the rows and drop the `PAGE_URL`, `BROWSER`, and `IP_ADDRESS` columns. Example:


```
df = df.sample(0.1)
df = df.drop('PAGE_URL', 'BROWSER', 'IP_ADDRESS')
```

- A. Option E
- B. Option B
- C. Option C
- D. Option A
- E. Option D

Answer: C

Explanation:

Option C is the most effective for this scenario. Focusing on sessions and their durations provides a more meaningful feature for predicting conversion rates. Removing bot traffic (A) might be a useful preprocessing step but doesn't fundamentally address session-level relevance. Option B's logic is flawed removing all Internet Explorer traffic isn't inherently removing irrelevant data. Option D oversimplifies the data, losing valuable information about user behavior within sessions. Option E introduces bias by randomly sampling and removing potentially important patterns, plus it is too simplistic. The code example in C demonstrates how to calculate session duration using Snowpark functions, join the filtered session data back to the original data, and then drop the irrelevant columns.

NEW QUESTION # 151

You are using Snowflake ML to train a binary classification model. After training, you need to evaluate the model's performance. Which of the following metrics are most appropriate to evaluate your trained model, and how do they differ in their interpretation, especially when dealing with imbalanced datasets?

- A. AUC-ROC: Measures the ability of the model to distinguish between classes. It is less sensitive to class imbalance than accuracy. Log Loss: Measures the performance of a classification model where the prediction input is a probability value between 0 and 1.
- B. Mean Squared Error (MSE): The average squared difference between the predicted and actual values. R-squared: Represents the proportion of variance in the dependent variable that is predictable from the independent variables. These are great for regression tasks.
- C. Confusion Matrix: A table that describes the performance of a classification model by showing the counts of true positive, true negative, false positive, and false negative predictions. This isn't a metric but representation of the metrics.
- D. Accuracy: It measures the overall correctness of the model. Precision: It measures the proportion of positive identifications that were actually correct. Recall: It measures the proportion of actual positives that were identified correctly. F1-score: It is the harmonic mean of precision and recall.
- E. Precision, Recall, F1-score, AUC-ROC, and Log Loss: Precision focuses on the accuracy of positive predictions; Recall focuses on the completeness of positive predictions; F1-score balances Precision and Recall; AUC-ROC evaluates the separability of classes and Log Loss quantifies the accuracy of probabilities, especially valuable for imbalanced datasets because they provide a more nuanced view of performance than accuracy alone.

Answer: E

Explanation:

Option E correctly identifies the most appropriate metrics (Precision, Recall, F1-score, AUC-ROC, and Log Loss) for evaluating a binary classification model, especially in the context of imbalanced datasets. It also correctly describes the focus of each metric. Accuracy can be misleading with imbalanced datasets. MSE and R-squared are for regression problems (Option B). Confusion Matrix is a table, and Options D, contains incorrect statement.

NEW QUESTION # 152

Consider the following Python UDF intended to train a simple linear regression model using scikit-learn within Snowflake. The UDF takes feature columns and a target column as input and returns the model's coefficients and intercept as a JSON string. You are encountering an error during the CREATE OR REPLACE FUNCTION statement because of the incorrect deployment of the package during runtime. What would be the right way to fix this deployment and execute your model?

- A. The package 'scikit-learn' needs to be included in the import statement and deployed while creation of the 'Create or Replace function' statement, by including parameter. Also the correct code is to ensure the model can be trained and return the coefficients and intercept of the model.
- B. The package 'scikit-learn' needs to be included in the import statement and deployed while creation of the 'Create or Replace function' statement, by including parameter. Also the correct code is to ensure the model can be trained and return the coefficients and intercept of the model.
- C. The code works seamlessly without modification as Snowflake automatically resolves all the dependencies and ensures the execution of code within the create or replace function statement.
- D. The package 'scikit-learn' needs to be included in the import statement and deployed while creation of the 'Create or Replace function' statement, by including parameter. Also the correct code is to ensure the model can be trained and return the coefficients and intercept of the model.
- E. The required packages 'scikit-learn' is not present. The correct way to create UDF is by including the import statement within the function along with the deployment.

Answer: A

Explanation:

Option E is the correct option and provides explanation for deploying the packages and ensuring that model executes successfully.

NEW QUESTION # 153

A data scientist is tasked with building a predictive maintenance model for industrial equipment. The data is collected from IoT sensors and stored in Snowflake. The raw sensor data is voluminous and contains noise, outliers, and missing values. Which of the following code snippets, executed within a Snowflake environment, demonstrates the MOST efficient and robust approach to cleaning and transforming this sensor data during the data collection phase, specifically addressing outlier removal and missing value imputation using robust statistics? Assume necessary libraries like numpy and pandas are available via Snowpark.

- A.

```
import snowflake.snowpark.functions as F
import numpy as np
from snowflake.snowpark.functions import col

def clean_sensor_data(df):
    #Outlier capping based on 1st and 99th percentile values
    p1 = df.approx_quantile("sensor_value", 0.01)
    p99 = df.approx_quantile("sensor_value", 0.99)

    df = df.with_column("sensor_value", F.when(col("sensor_value") < p1[0], p1[0]).when(col("sensor_value") > p99[0], p99[0]).otherwise(col("sensor_value")))

    # Impute missing values using median
    median_val = df.approx_quantile("sensor_value", 0.5)
    df = df.fillna(median_val[0], subset=["sensor_value"])

    return df
```

```

import snowflake.snowpark.functions as F

def clean_sensor_data(df):
    # Remove outliers using a fixed threshold
    df = df[df["sensor_value"] < 1000] # Assuming sensor value should be less than 1000

    # Impute missing values with 0
    df["sensor_value"] = df["sensor_value"].fillna(0)

```

- B. `return df`
- C.

```

import snowflake.snowpark.functions as F
import numpy as np

def clean_sensor_data(df):
    # Outlier removal using interquartile range (IQR)
    Q1 = df["sensor_value"].quantile(0.25)
    Q3 = df["sensor_value"].quantile(0.75)
    IQR = Q3 - Q1
    df = df[(df["sensor_value"] >= (Q1 - 1.5 * IQR)) & (df["sensor_value"] <= (Q3 + 1.5 * IQR))]

    # Missing value imputation using median
    df["sensor_value"] = df["sensor_value"].fillna(df["sensor_value"].median())
    return df

```

- D.

```

import snowflake.snowpark.functions as F

def clean_sensor_data(df):
    # Simple outlier removal using z-score
    z_scores = F.abs((df["sensor_value"] - df["sensor_value"].mean()) / df["sensor_value"].std())
    df = df[z_scores < 3]

    # Simple mean imputation for missing values
    df["sensor_value"] = df["sensor_value"].fillna(df["sensor_value"].mean())
    return df

```

- E.

```

import snowflake.snowpark.functions as F

def clean_sensor_data(df):
    # Do nothing - skip outlier removal and missing value imputation
    return df

```

Answer: A

Explanation:

Option E is the MOST robust and efficient. It uses the interquartile range (IQR) method, which is less sensitive to extreme outliers than the z-score method in Option A. It also utilizes 'approx_quantile' and is therefore more optimized for Snowflake large datasets. The median is also a more robust measure of central tendency for imputation than the mean when dealing with outliers. Option C uses a hard-coded threshold for outlier removal and imputes with 0, which is not adaptive or robust. Option D skips data cleaning altogether. Option A uses z-score which may work however, since IoT has continuous streaming data quantile based outlier removal is better. It is more optimised for large dataset and better at handling streaming datasets.

NEW QUESTION # 154

You're analyzing the performance of two different AIB testing variants of an advertisement. You've collected the following data over

a period of one week: Variant A: 1000 impressions, 50 conversions Variant B: 1100 impressions, 66 conversions Which of the following statements are TRUE regarding confidence intervals and statistical significance in this scenario?

- A. Constructing a 95% confidence interval for the difference in conversion rates between Variant B and Variant A will allow you to assess if there is a statistically significant difference at the 5% significance level. If the confidence interval contains zero, there is no statistically significant difference.
- B. Calculating separate confidence intervals for conversion rates A and B, and noting overlap, is an invalid method to infer statistical significance. One must construct confidence interval for the difference in means.
- C. Increasing the sample size (number of impressions for each variant) will generally widen the confidence interval, making it more likely to contain zero.
- D. If the 95% confidence interval for the conversion rate of Variant A is entirely above the 95% confidence interval for the conversion rate of Variant B, then Variant A is statistically better than Variant B.
- E. A narrower confidence interval for the difference in conversion rates implies a higher degree of certainty about the estimated difference.

Answer: A,B,E

Explanation:

Options A, B, and E are correct. Option A correctly explains the relationship between confidence intervals and statistical significance at a given significance level. Option B is correct because narrower interval correctly infers higher certainty. Option E is correct since you need a single measure of difference not each variable measured separately. Option C is incorrect: increasing the sample size will generally narrow the confidence interval, making it less likely to contain zero. Option D is incorrect. You cannot conclude statistical superiority by comparing if one confidence interval is entirely above other. You must construct a difference interval to compare. There is more to overlap than just that.

NEW QUESTION # 155

.....

The high quality and high efficiency of DSA-C03 study guide make it stand out in the products of the same industry. Our study materials have always been considered for the users. If you choose our DSA-C03 exam questions, you will become a better self. DSA-C03 actual exam want to contribute to your brilliant future. Our study materials are constantly improving themselves. If you have any good ideas, our study materials are very happy to accept them. DSA-C03 Exam Materials are looking forward to having more partners to join this family. We will progress together and become better ourselves.

DSA-C03 Study Group: <https://www.examslabs.com/Snowflake/SnowPro-Advanced/best-DSA-C03-exam-dumps.html>

Snowflake Vce DSA-C03 Format If you have an existing PayPal account, you can log in using your user data to confirm the payment, As you know, there are so many users of our DSA-C03 guide questions, Nowadays, some corporation and employer attach much importance on the Snowflake DSA-C03 certification, Snowflake Vce DSA-C03 Format With the development of scientific and technological progress, being qualified by some certifications plays an increasingly important role in our life.

Moreover, to write the Up-to-date DSA-C03 practice braindumps, they never stop the pace of being better, This lesson is where you learn how to protect your Exchange mailboxes from spam and viruses.

Snowflake DSA-C03 Exam Questions Come With Free 12 Months Updates

If you have an existing PayPal account, you can log in using your user data to confirm the payment, As you know, there are so many users of our DSA-C03 Guide questions.

Nowadays, some corporation and employer attach much importance on the Snowflake DSA-C03 certification, With the development of scientific and technological progress, being DSA-C03 qualified by some certifications plays an increasingly important role in our life.

Some of the test data on the site is free, but more importantly is that it provides a realistic simulation exercises that can help you to pass the Snowflake DSA-C03 exam

- Pass Guaranteed 2026 Perfect Snowflake DSA-C03: Vce SnowPro Advanced: Data Scientist Certification Exam Format
 Search for ➔ DSA-C03 on ➤ www.examcollectionpass.com ↵ immediately to obtain a free download DSA-C03 Latest Exam Dumps
- DSA-C03 Latest Exam Vce Simulated DSA-C03 Test DSA-C03 Reliable Practice Questions Search for ↗ DSA-C03 ↘ and download exam materials for free through (www.pdfvce.com) Answers DSA-C03 Real

Questions

DOWNLOAD the newest ExamLabs DSA-C03 PDF dumps from Cloud Storage for free: https://drive.google.com/open?id=17kD6agegdVx4KqkPkQUWJ_f-X62tgrmB