

Snowflake DSA-C03 Dumps—Best Option For Preparation



2026 Latest RealValidExam DSA-C03 PDF Dumps and DSA-C03 Exam Engine Free Share: <https://drive.google.com/open?id=1whnUBr7shhqHoZyxWc6kYgSIunopQsM9>

Download the free DSA-C03 demo of whatever product you want and check its quality and relevance by comparing it with other available study contents within your access. RealValidExam's study guides and DSA-C03 Dump will prove their worth and excellence. Check also the feedback of our clients to know how our products proved helpful in passing the exam.

For the challenging Snowflake DSA-C03 exam, they make an effort to locate reputable and recent Snowflake DSA-C03 practice questions. The high anxiety and demanding workload the candidate must face being qualified for the Snowflake DSA-C03 Certification are more difficult than only passing the Snowflake DSA-C03 exam.

>> DSA-C03 Dumps Collection <<

Snowflake DSA-C03 Test Sample Questions | DSA-C03 Latest Test Format

Our study materials have enough confidence to provide the best DSA-C03 exam torrent for your study to pass it. With many years work experience, we have fast reaction speed to market change and need. In this way, we have the latest DSA-C03 guide torrent. You don't worry about that how to keep up with the market trend, just follow us. We can say that our DSA-C03 Test Questions are the most suitable for examinee to pass the exam, you will never regret to buy it.

Snowflake SnowPro Advanced: Data Scientist Certification Exam Sample Questions (Q152-Q157):

NEW QUESTION # 152

You are preparing a dataset in Snowflake for a K-means clustering algorithm. The dataset includes features like 'age', 'income' (in USD), and 'number_of_transactions'. 'Income' has significantly larger values than 'age' and 'number_of_transactions'. To ensure that all features contribute equally to the distance calculations in K-means, which of the following scaling approaches should you consider, and why? Select all that apply:

- A. Apply PowerTransformer to transform income and StandardScaler to other features to handle skewness.
- B. Apply RobustScaler to handle outliers and then StandardScaler or MinMaxScaler to further scale the features.
- C. Apply MinMaxScaler to all three features to scale them to a range between 0 and 1 .
- D. Apply StandardScaler to all three features ('age', 'income', 'number_of_transactions') to center the data around zero and scale it to unit variance.
- E. Do not scale the data, as K-means is robust to differences in feature scales.

Answer: B,C,D

Explanation:

K-means clustering is sensitive to the scale of the features because it relies on distance calculations. Features with larger values will have a disproportionate influence on the clustering results. StandardScaler centers the data around zero and scales it to unit variance, which ensures that all features have a similar range and variance. MinMaxScaler scales the features to a range between 0 and 1, which also addresses the issue of different scales. RobustScaler handles outliers which will then use the other two scaling techniques. Therefore A, B and D are the appropriate scaling techniques. C is not correct as K-means relies on distance calculations and not scaling the data could give some feature a larger weight which isn't the desired outcome. Option E: Using PowerTransformer on 'income' to reduce skewness and StandardScaler on the other features can be a valid approach, but it depends on the distribution of 'income' and the presence of outliers. If 'income' is highly skewed and/or contains outliers, this combination might be more effective than using StandardScaler or MinMaxScaler alone.

NEW QUESTION # 153

You're a data scientist analyzing sensor data from industrial equipment stored in a Snowflake table named 'SENSOR READINGS'. The table includes 'TIMESTAMP', 'SENSOR ID', 'TEMPERATURE', 'PRESSURE', and 'VIBRATION'. You need to identify malfunctioning sensors based on outlier readings in 'TEMPERATURE', 'PRESSURE', and 'VIBRATION'. You want to create a dashboard to visualize these outliers and present a business case to invest in predictive maintenance. Select ALL of the actions that are essential for both effectively identifying sensor outliers within Snowflake and visualizing the data for a business presentation. (Multiple Correct Answers)

- A. Implement a clustering algorithm (e.g., DBSCAN) within Snowflake using Snowpark Python to group similar sensor readings, identifying outliers as points that do not belong to any cluster or belong to very small clusters.
- B. Calculate basic statistical summaries (mean, standard deviation, min, max) for each sensor and each variable C TEMPERATURE, 'PRESSURE, and 'VIBRATION') and use that information to filter down to the most important sensor, prior to using the other techniques.
- C. Calculate Z-scores for 'TEMPERATURE, 'PRESSURE, and 'VIBRATION' for each 'SENSOR_ID within a rolling window of the last 24 hours using Snowflake's window functions. Define outliers as readings with Z-scores exceeding a threshold (e.g., 3).
- D. Directly connect the 'SENSOR_READINGS' table to a visualization tool and create a 3D scatter plot with 'TEMPERATURE, 'PRESSURE, and 'VIBRATION' on the axes, without any pre-processing or outlier detection in Snowflake.
- E. Create a Snowflake stored procedure to automatically flag outlier readings in a new column 'IS OUTLIER based on a predefined rule set (e.g., IQR method or Z-score threshold), and then use this column to filter data for visualization in a dashboard.

Answer: A,B,C,E

Explanation:

Options A, C, D, and E are essential. A (Z-score calculation with rolling window) provides a dynamic measure of how unusual a reading is relative to recent history for each sensor. C (DBSCAN clustering) helps identify outliers based on density; points far from any cluster are likely outliers. D (Stored procedure with outlier flagging) automates the outlier detection process and makes it easy to filter and visualize outliers in a dashboard, with a business ready explanation. Option E allows you to focus on the right data, allowing you to have a more useful visualisation. Option B (direct 3D scatter plot without pre-processing) is not effective because it will be difficult to identify outliers visually in a high- density scatter plot without any outlier detection or data reduction. The direct scatter plot becomes overwhelming very quickly with sensor data.

NEW QUESTION # 154

A data scientist is analyzing website click-through rates (CTR) for two different ad campaigns. Campaign A ran for two weeks and had 10,000 impressions with 500 clicks. Campaign B also ran for two weeks with 12,000 impressions and 660 clicks. The data scientist wants to determine if there's a statistically significant difference in CTR between the two campaigns. Assume the population standard deviation is unknown and unequal for the two campaigns. Which statistical test is most appropriate to use, and what

Snowflake SQL code would be used to approximate the p-value for this test (assume 'clicks_b', and are already defined Snowflake variables)?

- A. An independent samples t-test, because we are comparing the means of two independent samples. Snowflake code: `SELECT`

```
SELECT ARRAY_CONSTRUCT(CAST(clicks_a/IMPRESSIONS_A AS ARRAY<CONSTRUCT(CAST(impressions_b) AS VAR EQUAL=FALSE))
```

- B. A one-sample t-test, because we are comparing the sample mean of campaign A to the sample mean of campaign Snowflake code: 'SELECT t_test_lsamp(clicks_a/impressions_a - clicks_b/impressions_b, 0)'
- C. An independent samples t-test (Welch's t-test), because we are comparing the means of two independent samples with unequal variances. Snowflake code (approximation using UDF - assuming UDF 'p_value_from_t_stat' exists that calculates p-value from t-statistic and degrees of freedom):

```
CREATE OR REPLACE FUNCTION calculate_welch_t_test(n1 FLOAT, mean1 FLOAT, var1 FLOAT, n2 FLOAT, mean2 FLOAT, var2 FLOAT) RETURNS FLOAT LANGUAGE JAVASCRIPT AS $$
var t_stat = (mean1 - mean2) / Math.sqrt((var1 / n1) + (var2 / n2));
var df = Math.pow((var1 / n1) + (var2 / n2), 2) / ((Math.pow(var1 / n1, 2) / (n1 - 1)) + (Math.pow(var2 / n2, 2) / (n2 - 1)));
return t_stat;
$$
SELECT p_value from t_stat(calculate_welch_t_test(10000, 500/10000, VAR(SELECT clicks FROM campaign_a_data), 12000, 660/12000, VAR(SELECT clicks FROM campaign_b_data)), ...); //fill in degrees of freedom calculation in UDF.
```

- D. A paired t-test, because we are comparing two related samples over time. Snowflake code: 'SELECT t_test_ind(clicks_a/impressions_a, 'VAR EQUAL=TRUE)'
- E. Az-test, because we know the population standard deviation. Snowflake code: 'SELECT normcdf(clicks_a/impressions_a - clicks_b/impressions_b, 0, 1)'

Answer: A

Explanation:

The correct answer is E. Since we're comparing the means of two independent samples (Campaign A and Campaign B) and the population standard deviations are unknown, an independent samples t-test is appropriate. Because the problem stated that the variances are unequal, Welch's t-test provides a more accurate p-value and confidence intervals. The Snowflake function handles independent samples and the 'VAR_EQUAL=FALSE' parameter specifies that the variances should not be assumed to be equal. The other options are incorrect because they use inappropriate tests given the problem conditions. The z-test is not appropriate because the population standard deviations are unknown. A paired t-test is for related samples. A one sample test is to test one mean against a constant not another mean.

NEW QUESTION # 155

You are tasked with estimating the 95% confidence interval for the median annual income of Snowflake customers. Due to the non-normal distribution of incomes and a relatively small sample size (n=50), you decide to use bootstrapping. You have a Snowflake table named 'customer_income' with a column 'annual_income'. Which of the following SQL code snippets, when correctly implemented within a Python script interacting with Snowflake, would most accurately achieve this using bootstrapping with 1000 resamples and properly calculate the confidence interval?

- A.

```
python import snowflake.connector import numpy as np # Snowflake connection details (replace with your actual credentials) conn = snowflake.connector.connect( user='your_user', password='your_password', account='your_account' ) cs = conn.cursor() cs.execute('SELECT annual_income FROM customer_income') data = [row[0] for row in cs.fetchall()] bootstrap_medians = [] for i in range(1000): resample = np.random.choice(data, size=len(data), replace=True) bootstrap_medians.append(np.median(resample)) lower_bound = np.percentile(bootstrap_medians, 2.5) upper_bound = np.percentile(bootstrap_medians, 97.5) print(f'95% Confidence Interval: ({lower_bound}, {upper_bound})') conn.close()
```

- B.

```
python import snowflake.connector import random # Snowflake connection details (replace with your actual credentials) conn = snowflake.connector.connect( user='your_user', password='your_password', account='your_account' ) cs = conn.cursor() cs.execute('SELECT annual_income FROM customer_income') data = [row[0] for row in cs.fetchall()] bootstrap_means = [] for i in range(1000): resample = random.sample(data, len(data)) bootstrap_means.append(sum(resample) / len(resample)) lower_bound = min(bootstrap_means) upper_bound = max(bootstrap_means) print(f'95% Confidence Interval: ({lower_bound}, {upper_bound})') conn.close()
```

- C.

```
python import snowflake.connector import numpy as np # Snowflake connection details (replace with your actual credentials) conn = snowflake.connector.connect( user='your_user', password='your_password', account='your_account' ) cs = conn.cursor() cs.execute('SELECT annual_income FROM customer_income') data = [row[0] for row in cs.fetchall()] bootstrap_medians = [] for i in range(1000): resample = np.random.choice(data, size=len(data), replace=True) bootstrap_medians.append(np.mean(resample)) lower_bound = np.percentile(bootstrap_medians, 0.025) upper_bound = np.percentile(bootstrap_medians, 0.975) print(f'95% Confidence Interval: ({lower_bound}, {upper_bound})') conn.close()
```

- D.

```
python import snowflake.connector import numpy as np # Snowflake connection details (replace with your actual credentials) conn = snowflake.connector.connect( user='your_user', password='your_password', account='your_account' ) cs = conn.cursor() cs.execute('SELECT annual_income FROM customer_income') data = [row[0] for row in cs.fetchall()] bootstrap_medians = [] for i in range(1000): resample = np.random.choice(data, size=len(data), replace=True) bootstrap_medians.append(np.median(resample)) lower_bound = np.percentile(bootstrap_medians, 5) upper_bound = np.percentile(bootstrap_medians, 95) print(f'95% Confidence Interval: ({lower_bound}, {upper_bound})') conn.close()
```

- E.

```

○ """python import snowflake.connector # Snowflake connection details (replace with your actual credentials) conn = snowflake.connector.connect(
user='your_user', password='your_password', account='your_account' ) cs = conn.cursor() cs.execute("SELECT annual_income FROM customer_income") data =
[row[0] for row in cs.fetchall()] # Incorrect Bootstrap implementation bootstrap_values = [sum(data) / len(data) for _ in range(1000)] lower_bound =
min(bootstrap_values) upper_bound = max(bootstrap_values) print(f'95% Confidence Interval: ({lower_bound}, {upper_bound})') conn.close() """

```

Answer: A

Explanation:

Option A is the correct answer. It accurately implements bootstrapping by: (1) Resampling with replacement from the original data. (2) Calculating the median of each resample. (3) Computing the 2.5th and 97.5th percentiles of the bootstrap medians to obtain the 95% confidence interval. Option B calculates the mean instead of the median, and uses 'random.sample' without replacement, which is incorrect for bootstrapping. Option C doesn't resample at all, just calculates the mean of the original data repeatedly. Option D calculates the mean instead of the median. Option E calculates 90% confidence interval instead of 95%.

NEW QUESTION # 156

You are building a data science pipeline in Snowflake to predict customer churn. The pipeline involves extracting data, transforming it using Dynamic Tables, training a model using Snowpark ML, and deploying the model for inference. The raw data arrives in a Snowflake stage daily as Parquet files. You want to optimize the pipeline for cost and performance. Which of the following strategies are MOST effective, considering resource utilization and potential data staleness?

- **A. Implement a series of smaller Dynamic Tables, each responsible for a specific transformation step, with well-defined refresh intervals tailored to the data's volatility and the downstream model's requirements.**
- B. Load all data into traditional Snowflake tables and use scheduled tasks with stored procedures written in Python to perform the transformations and model training.
- **C. Use a combination of Dynamic Tables for feature engineering and Snowpark ML for model training and deployment, ensuring proper dependency management and refresh intervals for each Dynamic Table based on data freshness requirements.**
- D. Use a single, large Dynamic Table to perform all transformations in one step, relying on Snowflake's optimization to handle dependencies and incremental updates.
- E. Schedule all data transformations and model training as a single large Snowpark Python script executed by a Snowflake task, ignoring data freshness requirements.

Answer: A,C

Explanation:

Option B is correct because breaking down the transformations into smaller Dynamic Tables with tailored refresh intervals ensures that only necessary data is recomputed, minimizing cost and resource usage. Option D is also correct because combining Dynamic Tables for feature engineering with Snowpark ML allows for efficient model training and deployment within Snowflake, leveraging the platform's scalability and security features. Furthermore, specifying refresh intervals based on data freshness guarantees that the model is trained on up-to-date information. Option A could lead to unnecessary computations if the entire table needs to be refreshed even for minor changes. Option C relies on traditional tables and stored procedures, which may be less efficient and harder to manage than Dynamic Tables. Option E ignores data freshness, which can significantly impact model accuracy.

NEW QUESTION # 157

.....

You can easily get SnowPro Advanced: Data Scientist Certification Exam (DSA-C03) certified if you prepare with our Snowflake DSA-C03 questions. Our product contains everything you need to ace the DSA-C03 certification exam and become a certified IT professional. So what are you waiting for? Purchase this updated SnowPro Advanced: Data Scientist Certification Exam (DSA-C03) exam practice material today and start your journey to a shining career.

DSA-C03 Test Sample Questions: <https://www.realvalidexam.com/DSA-C03-real-exam-dumps.html>

Snowflake DSA-C03 Dumps Collection With it, you will get a different life, You can successfully prepare for the DSA-C03 exam in a short time with the help of our latest exam questions, RealValidExam has come up with real Snowflake DSA-C03 Questions for students so they can pass SnowPro Advanced: Data Scientist Certification Exam (DSA-C03) exam in a single try and get to their destination, Snowflake DSA-C03 Dumps Collection Our expert team guarantees that each answer and question is useful and valuable.

Advanced database design topics, Ade Miller DSA-C03 is currently a Principal Software Architect at Microsoft Studios, With it, you will get a different life, You can successfully prepare for the DSA-C03 Exam in a short time with the help of our latest exam

