





SQL and Pandas Revision Notes

Introduction to the Class

The session primarily focused on explaining SQL Window functions and their equivalent implementations in Python using Pandas. The instructor initiated the session by addressing the agenda and interacting with the learners to tailor the lecture according to their background knowledge.

SQL Window Functions Overview

SQL Window functions are pivotal in performing operations across a set of table rows that are somehow related to the current row.

Popular SQL Window Functions include:

- **ROW_NUMBER()**: Assigns a unique sequential integer to rows within a partition.
- **RANK()**: Provides a rank to rows, with gaps in ranking for ties.
- **DENSE_RANK()**: Similar to RANK(), but without gaps between ranked positions for ties.
- **NTILE()**: Distributes the rows in an ordered partition into a specified number of sets.
- **LAG()**: Accesses data from a previous row in the same result set without needing a self-join.
- **LEAD()**: Provides access to a subsequent row at a specified physical offset.

Other window functions also include functions for calculating cumulative distribution, percentages, and more [【8:0+source】](#) [【8:11+source】](#).

Pandas Equivalents for SQL Window Functions

The class explained how these SQL functionalities can be implemented using Python's Pandas library. Key equivalent



- **Ranking:** Achieved using `df.rank()` function for RANK and `method='min'` or `method='dense'` for DENSE_RANK.
- **Lag and Lead Operations:** Performed using `df.shift(+1)` for lead and `df.shift(-1)` for lag **【8:2+source】 【8:16+source】**.

Example: Calculating Median

To calculate the median in Pandas without a direct median function:

1. **Sorting:** Sort the data based on a specific column.
2. **Row Number Calculation:** Use a row numbering method through ranks or custom logic to identify median positions.
3. **Conditional Filtering:** Based on odd or even number of rows, derive the median either directly or by averaging two middle values **【8:8+source】**.

Problem Solving

Throughout the session, several exercises were discussed:

- **Query-Based Approaches vs. Programmatic Approaches:** Highlighted the decision-making process when choosing SQL versus Pandas based on the task needs.
- **Join Operations and Data Merging:** Explained with examples how to merge tables and handle joins both in SQL (`JOIN` operations) and in Pandas (`pd.merge` with `how` parameter).
- **Practical Applications:** Problems involved calculating distances, swapping values based on conditions, and more complex operations involving both SQL and Pandas **【8:4+source】 【8:9+source】 【8:19+source】**.

Additional Topics Discussed

- **Efficiency of Code:** Emphasized understanding logic over memorizing syntax, as technical interviews focus on logical problem-solving.
- **Practical Applications:** Mentioned exercises such as calculating variances using a window frame in Pandas, which is akin to calculating over partitions in SQL **【8:14+source】**.

Conclusion



further resources if needed, and an encouragement to solve practice problems across different platforms to improve technical skills.

These notes capture the essence of the class, focusing on the conversion between SQL and Pandas, practical problem-solving, and the technical details shared by the instructor.