

3. IMPUTER

Mục tiêu IMPUTER: được sử dụng để xử lý các giá trị thiếu trong dữ liệu, tức là các ô dữ liệu mà không có giá trị (có thể là NaN, null, hoặc các giá trị không có ý nghĩa khác). Khi dữ liệu bị thiếu, việc sử dụng Imputer giúp điền các giá trị thay thế vào những vị trí thiếu này để có thể tiếp tục phân tích dữ liệu một cách chính xác.

- [Thực hiện trên hệ thống theo các bước hướng dẫn sau:](#)

Thực hiện trên hệ thống theo các bước hướng dẫn sau:

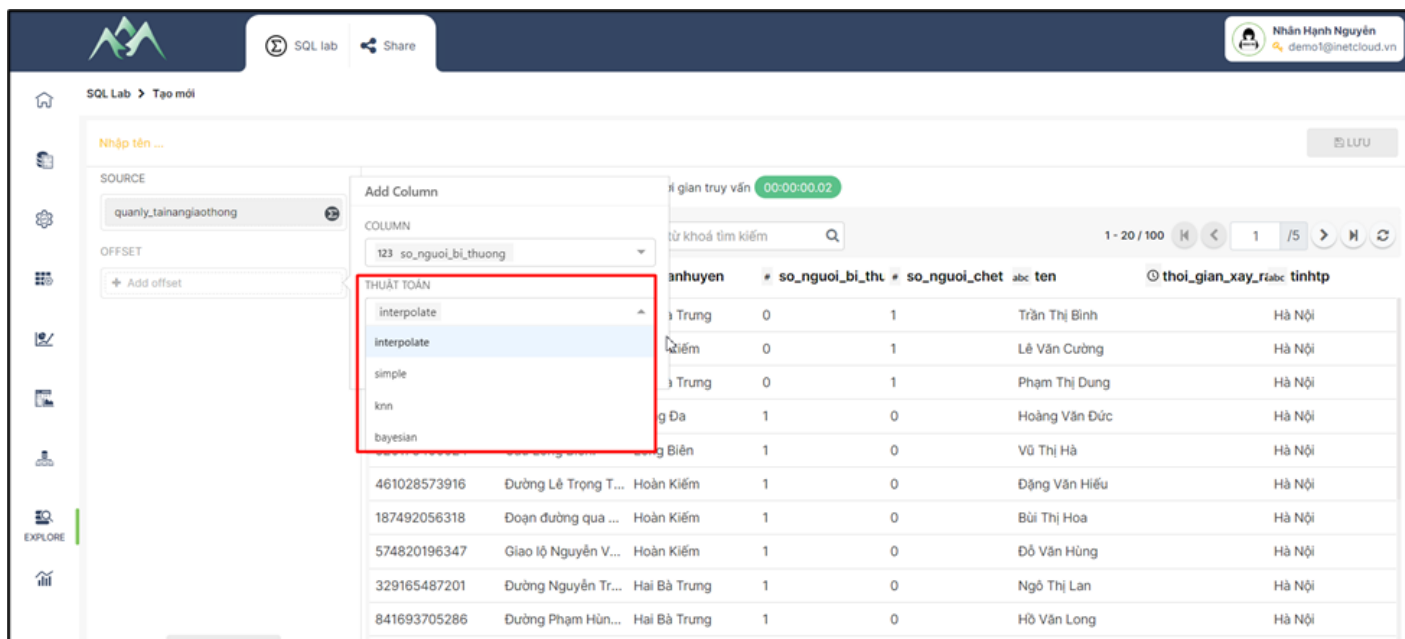
Bước 1: Vào hệ thống chọn module Explore, chọn tập dữ liệu cần tổng hợp và tạo mới một tập dữ liệu tổng hợp từ IMPUTER

The screenshot shows the SQL Lab interface. On the left, the 'DATASET' list includes 'storagehub.hanh_chinh.quanly_t...'. The main area displays a table with columns: 'ccccc', 'địa_diem_tai_nha', 'quan_huyen', 'so_nguoi_bi_thu', 'so_nguoi_chet', 'ten', and 'thoi_gian_xay'. The 'MORE TOOLS' menu is open, and the 'IMPUTER' option is highlighted.

Bước 2: Chọn cột dữ liệu để thuật toán Imputer giúp điền các giá trị thay thế vào những vị trí thiếu.

The screenshot shows the 'Add Column' dialog box in the SQL Lab interface. The 'COLUMN' list includes '123_so_nguoi_bi_thuong'. The 'Add offset' button is highlighted.

Bước 3: Chọn thuật toán xử lý cho cột dữ liệu đã chọn.



The screenshot shows the SQL Lab interface. In the center, an 'Add Column' dialog box is open, displaying a dropdown menu for selecting a mathematical method. The dropdown menu is highlighted with a red box and contains the following options: 'interpolate', 'interpolate', 'simple', 'knn', and 'bayesian'. The background shows a data table with columns like 'so_nguoi_bi_thu' and 'so_nguoi_chet'.

Trong đó:

- o Thuật toán nội suy (interpolation) được sử dụng để ước lượng giá trị tại các điểm trung gian dựa trên các giá trị đã biết của một tập dữ liệu.
- o Thuật toán Simple được dùng để xử lý các bài toán như tìm kiếm tuyến tính, sắp xếp nổi bọt, hoặc tính giai thừa.
- o Thuật toán KNN là một thuật toán dựa trên khoảng cách để phân loại hoặc dự đoán giá trị của một điểm dữ liệu mới dựa trên các điểm dữ liệu đã biết trong không gian đa chiều.
- o Thuật toán Payesian là một trong những thuật toán phân loại phổ biến nhất trong học máy, đặc biệt là đối với các bài toán phân loại văn bản và lọc thư rác.

Bước 4: “RUN” chạy test dữ liệu được tổng hợp.

