

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:

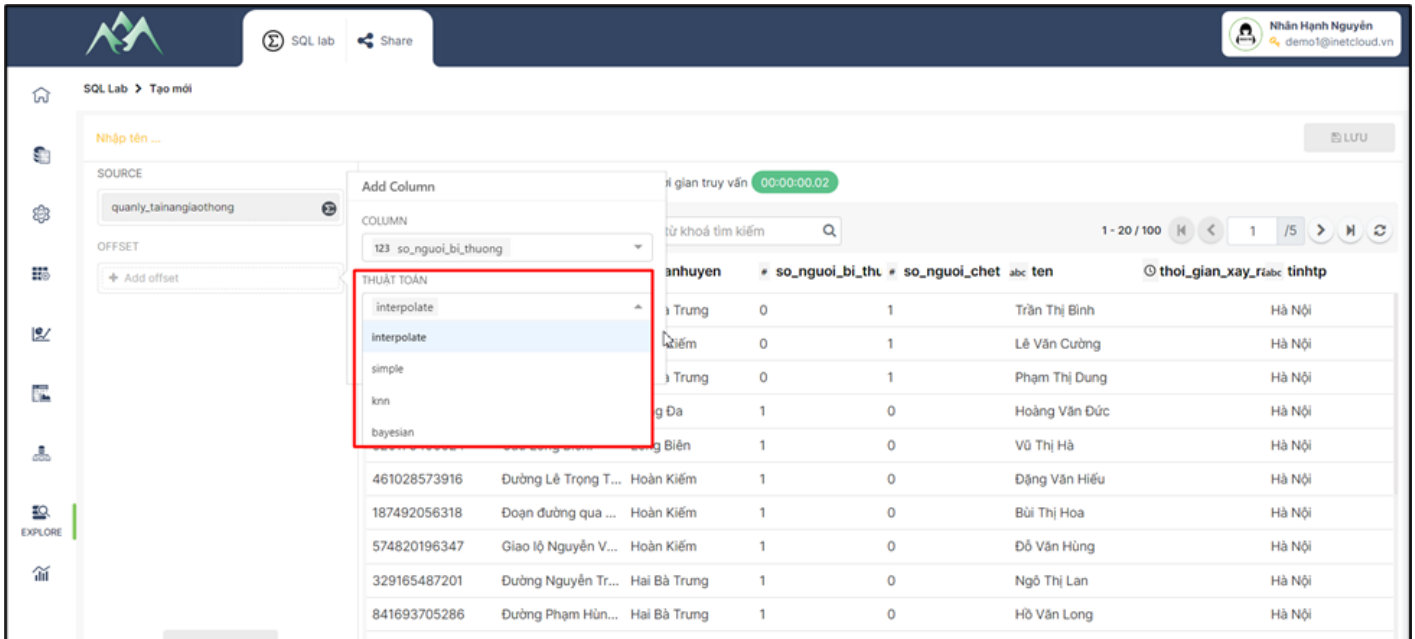
Bước 1: Vào hệ thống chọn module Explore, chọn tập dữ liệu cần tổng hợp và tạo mới một tập dữ liệu tổng hợp từ IMPUTER

The screenshot shows the SQL Lab interface. On the left, the 'DATASET' list includes 'storagehub.hanh_chinh.quanly_t...'. The 'DỮ LIỆU' table displays data with columns: 'abc', 'ccccc', 'địa điểm tại nạn', 'quanhuyen', 'so_nguoi_bi_thu', 'so_nguoi_chet', 'ten', and 'thoi_gian_xay'. The 'IMPUTER' button is highlighted in the top right corner.

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 SQL Lab interface with the 'Add Column' dialog box open. The 'COLUMN' list includes '123 so_nguoi_bi_thuong'. The 'Add offset' button is highlighted in the bottom left corner.

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



The screenshot shows the SQL Lab interface. In the center, the 'Add Column' dialog box is open. The 'COLUMN' dropdown is set to '123 so_nguoi_bi_thuong'. The 'THUẬT TOÁN' (Algorithm) dropdown is open, showing a list of algorithms: 'interpolate', 'interpolate', 'simple', 'knn', and 'bayesian'. The first 'interpolate' option is selected. The background shows a table with columns: 'anhuynh', '# so_nguoi_bi_thu', '# so_nguoi_chet', 'abc ten', and 'thoi_gian_xay_ra'. The table contains data for various locations and counts.

anhuynh	# so_nguoi_bi_thu	# so_nguoi_chet	abc ten	thoi_gian_xay_ra
Trung	0	1	Trần Thị Bình	Hà Nội
Hoàn Kiếm	0	1	Lê Văn Cường	Hà Nội
Trung	0	1	Phạm Thị Dung	Hà Nội
Đa	1	0	Hoàng Văn Đức	Hà Nội
Đa	1	0	Vũ Thị Hà	Hà Nội
Đa	1	0	Đặng Văn Hiếu	Hà Nội
Đa	1	0	Bùi Thị Hoa	Hà Nội
Đa	1	0	Đỗ Văn Hùng	Hà Nội
Đa	1	0	Ngô Thị Lan	Hà Nội
Đa	1	0	Hồ Văn Long	Hà Nội

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.

