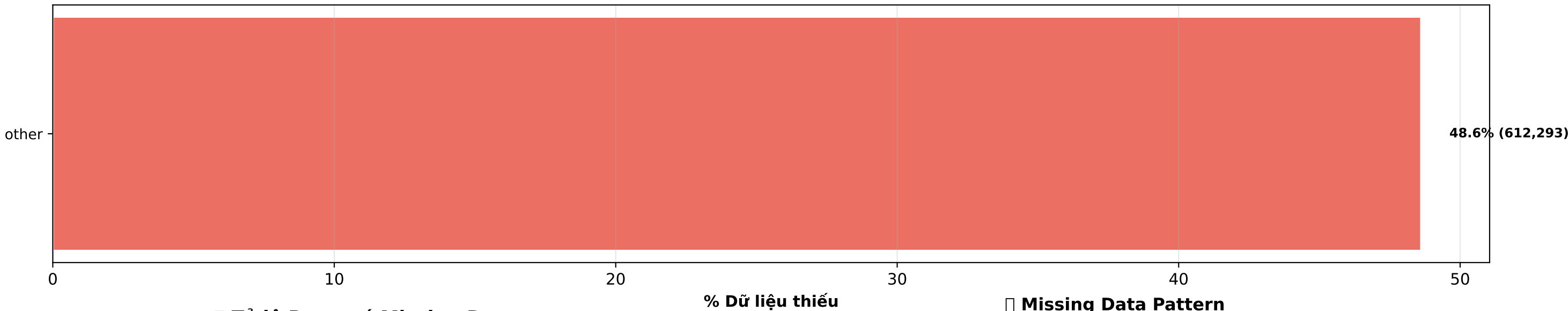


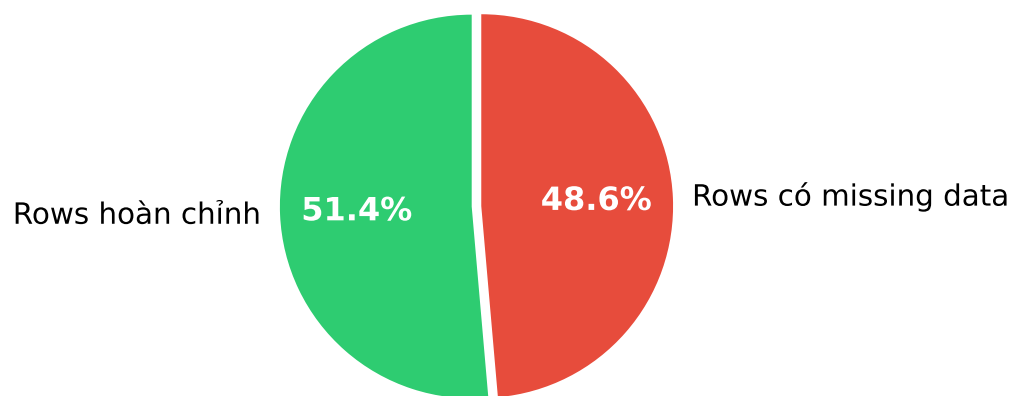
KHÓ KHĂN 3: MISSING DATA (Dữ liệu thiếu)

Phân tích các trường dữ liệu bị null và tác động đến phân tích

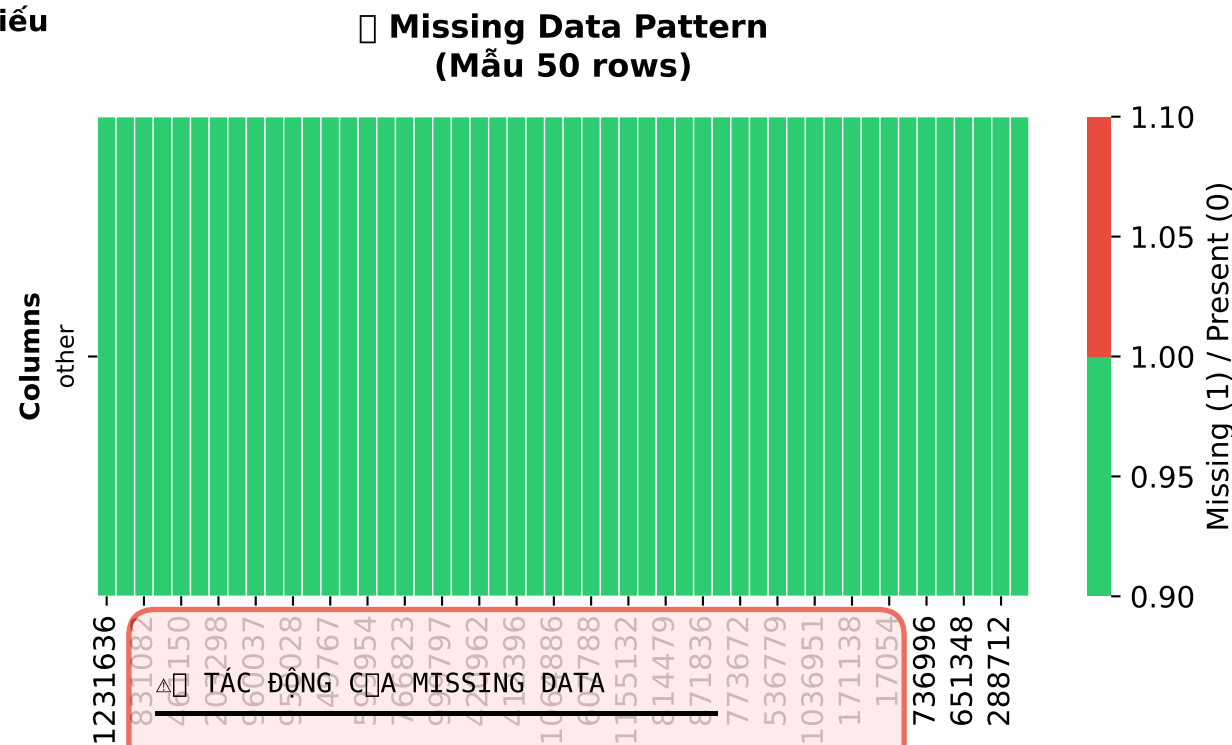
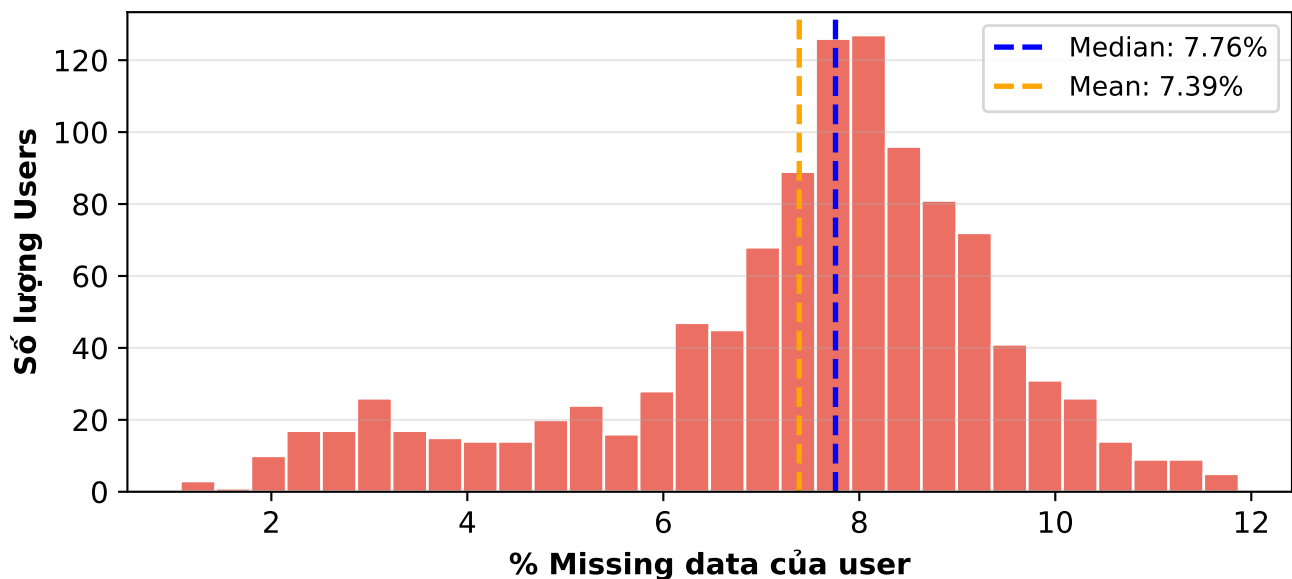
Top Columns có Missing Data
(Tỷ lệ % dữ liệu bị thiếu)



Tỷ lệ Rows có Missing Data



Phân phối Missing Data per User
(Chỉ users có missing)



TÁC ĐỘNG CỦA MISSING DATA

- Thống kê:
- Tổng cells: 10,075,288
 - Missing cells: 612,293 (6.08%)
 - Rows bị ảnh hưởng: 612,293 (48.62%)
 - Columns bị ảnh hưởng: 1/8
- Vấn đề:
- Làm giảm chất lượng feature engineering
 - Nhiều thuật toán ML không xử lý được null
 - Có thể làm sai lệch phân tích thống kê
 - Affect clustering và recommendation
- Giải pháp:
- Phân tích pattern: MCAR, MAR, hay MNAR?
 - Imputation: mean, median, mode, KNN
 - Forward/Backward fill cho time series
 - Loại bỏ columns có >50% missing
 - Tạo indicator variable (is_missing)
 - Sử dụng models chấp nhận null (XGBoost)