

**Dataset\_for\_classification.csv**, giải thích:

- Là dataset của ngành Tiếp thị ngân hàng, cho biết dữ liệu khách hàng của ngân hàng. Mục tiêu của chiến dịch tiếp thị là thuyết phục khách hàng gửi tiền tiết kiệm có kì hạn
- Số lượng data point (số hàng): trên 41K
- Số lượng data features: 20, gồm các cột từ A→ T, giải thích ý nghĩa như sau (lưu ý có 1 số chữ được dịch ra tiếng Việt cho dễ hiểu nhưng sinh viên không cần chuyển qua tiếng Việt trong bài làm)
  1. tuoi: tuổi
  2. nghe\_nghiep: loại công việc (gồm có các loại nghề nghiệp trong tiếng Anh)
  3. hon\_nhan: tình trạng hôn nhân (phân loại: "đã ly dị", "đã kết hôn", "độc thân", "không xác định")
  4. hoc\_van: tình trạng học vấn
  5. co\_the\_tin\_dung: có, hoặc không, hoặc không rõ (unknown)
  6. co\_nha\_o: có nhà ở không? Gồm có: có, không, không rõ
  7. vay\_ca\_nhan: các khoản vay cá nhân, gồm có: có, không, không rõ
  8. kênh\_lien\_lac: liên lạc với khách hàng qua điện thoại di động hoặc dt bàn
  9. thang\_lien\_lac: tháng mà nhân viên ngân hàng đã liên hệ với khách hàng
  10. ngay\_lien\_lac: ngày mà nhân viên ngân hàng đã liên hệ với khách hàng
  11. thoi\_luong\_lien\_lac: thời lượng của lần liên lạc cuối cùng, tính bằng giây
  12. so\_luong\_lien\_lac: số lượng lần liên lạc được thực hiện trong chiến dịch này và cho khách hàng này
  13. ngay: số ngày trôi qua sau khi khách hàng được liên hệ lần cuối từ một chiến dịch tiếp thị trước đó (số; 999 có nghĩa là khách hàng không được liên hệ trước đó)
  14. so\_luong\_lien\_lac\_truoc\_day: số lượng liên hệ được thực hiện cho các chiến dịch tiếp thị lần trước và cho khách hàng này
  15. ket\_qua\_lan\_truoc: kết quả của chiến dịch tiếp thị trước đó (gồm có: "thất bại", "không tồn tại", "thành công")
  16. ti\_le\_thay\_doi\_viec\_lam: tỷ lệ thay đổi việc làm theo quý trong năm
  17. CPI: chỉ số giá tiêu dùng tính theo hàng tháng
  18. CCI: chỉ số niềm tin của người tiêu dùng theo hàng tháng
  19. lai\_suat\_3thang: lãi suất chuẩn liên ngân hàng tại Khu vực sử dụng đồng euro
  20. so\_luong\_nhan\_vien: số lượng nhân viên làm việc cho ngân hàng
- Cột cuối cùng là Label, gồm có 2 phân lớp cần phân loại:
  - o 0: “thất bại”: khách hàng không chịu gửi tiết kiệm có kì hạn
  - o 1: “thành công” : khách hàng đã gửi tiết kiệm có kì hạn
- Các nhóm dùng 1 tập dataset này làm cho tất cả các giải thuật đã học về Classification
- Lưu ý:
  - o Không phải data features nào cũng giúp cho Classification 😊 , do đó mỗi nhóm tùy ý chọn data features nào để dùng cho giải thuật
  - o Nếu muốn vẽ plot phân lớp thì có thể chọn  $\leq 3$  features để plot trên không gian 3 chiều

- Các giá trị của Data features:
  - Có một số giá trị bị thiếu trong một số features, thường được ghi bằng thông tin “unknown”. Sinh viên có thể dùng các giá trị bị thiếu này như là 1 thông tin để dùng cho giải thuật Classification, hoặc có thể hoặc sử dụng các kỹ thuật xóa hoặc cắt bỏ tùy ý
  - Sinh viên xử lý theo cách nào thì ghi vào báo cáo
- Các nhóm tùy ý chia tỷ lệ số lượng data points dành cho training và testing

**Dataset\_for\_clustering.csv**, giải thích:

- Là dataset miêu tả các thông số kỹ thuật của xe hơi, gồm có:
  1. ten\_xe
  2. luong\_hao\_xang
  3. so\_luong\_xi\_lanh
  4. the\_tich\_dong\_co
  5. ma\_luc
  6. ty\_le\_truc\_sau
  7. khoi\_luong\_xe
  8. gia\_toc\_xe
  9. loai\_xy\_lanh\_dong\_co
  10. loai\_truyen\_dong
  11. so\_luong\_banh\_rang
  12. so\_luong\_bo\_che\_hoa\_khi
- Các nhóm dùng 1 tập dataset này làm cho tất cả các giải thuật đã học về Clustering
- Lưu ý:
  - Không phải data features nào cũng giúp cho Clustering 😊, do đó mỗi nhóm tùy ý chọn data features nào để dùng cho giải thuật
  - Nếu muốn vẽ plot phân lớp thì có thể chọn  $\leq 3$  features để plot trên không gian 3 chiều
  - Giải thuật HC yêu cầu vẽ được cây Dendrogram
  - Chia ra bao nhiêu phân cụm là tùy sinh viên lựa chọn
  - Khi nộp report, yêu cầu các nhóm nộp lại dataset này (file CSV) kèm thêm 1 cột cuối cùng:
    - Đặt tên cột là Cluster
    - Với từng dòng thì ghi rõ dòng đó thuộc Cluster nào, ví dụ, 1, 2, 3, 4....