

CODE BOOK Bài 2a

Thông tin	Nội dung
Tên bộ dữ liệu	Iris Plants Database
Nguồn thu thập và cách thức thu thập	(a) Creator: R.A. Fisher (b) Donor: Michael Marshall (MARSHALL%PLU@io.arc.nasa.gov) (c) Date: July, 1988
Số thuộc tính	5
Thông tin tên các thuộc tính	1. sepal length in cm 2. sepal width in cm 3. petal length in cm 4. petal width in cm 5. class: -- Iris Setosa -- Iris Versicolour -- Iris Virginica
Thông tin tác giả	(a) Creator: R.A. Fisher (b) Donor: Michael Marshall (MARSHALL%PLU@io.arc.nasa.gov)

Raw data: file iris.data

Tidy data: được lưu lại thành file: iris.csv

Instruction list:

CODE R

```
iris.R
1 rm(list=ls())
2
3 fileiris <- read.csv("dataset/iris.data", header = FALSE, sep = ",")
4
5 variables <- c("sepal length", "sepal width", "petal length", "petal width", "class")
6
7 colnames(fileiris) <- variables
8
9 write.csv(fileiris, file = "iris.csv", row.names = FALSE)
```

CODE BOOK Bài 2b

Thông tin	Nội dung
Tên bộ dữ liệu	Bank Marketing
Nguồn thu thập và cách thức thu thập	Paulo Cortez (Univ. Minho) and Sérgio Moro (ISCTE-IUL) @ 2012
Số thuộc tính	16 + thuộc tính đầu ra
Thông tin tên các thuộc tính	<p><i># Dữ liệu khách hàng của ngân hàng:</i></p> <p>1 - age (numeric)</p> <p>2 - job : loại công việc (phân loại: "admin.", "unknown", "unemployed", "management", "housemaid", "entrepreneur", "student", "bluecollar", "selfemployed", "retired", "technician", "services")</p> <p>3 - marital : Tình trạng hôn nhân (phân loại: "married", "divorced", "single"; note: "divorced" means divorced or widowed)</p> <p>4 - education (phân loại: "unknown", "secondary", "primary", "tertiary")</p> <p>5 - default: Có tính dụng trong tình trạng vỡ nợ? (nhị phân: "yes", "no")</p> <p>6 - balance: số dư trung bình hàng năm (numeric)</p> <p>7 - housing: có cho vay mua nhà không? (nhị phân: "yes", "no")</p> <p>8 - loan: có khoản vay cá nhân không? (nhị phân: "yes", "no")</p> <p>9 - contact: phương thức liên lạc (phân loại: "unknown", "telephone", "cellular")</p> <p>10 - day: ngày liên lạc cuối cùng của tháng (numeric)</p> <p>11 - month: tháng liên lạc cuối cùng của năm (phân loại: "jan", "feb", "mar", ..., "nov", "dec")</p> <p>12 - duration: thời lượng liên lạc cuối cùng, tính bằng giây (numeric)</p> <p><i># Thuộc tính khác:</i></p> <p>13 - campaign: số liên hệ được thực hiện trong chiến dịch này và cho khách hàng này (số, bao gồm liên hệ cuối cùng)</p> <p>14 - pdays: số ngày trôi qua sau khi khách hàng được liên hệ lần cuối từ một chiến dịch trước đó (numeric, -1 có nghĩa là khách hàng chưa được liên hệ trước đó)</p> <p>15 - previous: số lượng địa chỉ liên hệ được thực hiện trước chiến dịch này và cho khách hàng này (numeric)</p>

	<p>16 - poutcome: kết quả chiến dịch tiếp thị trước đó (phân loại: "unknown", "other", "failure", "success")</p> <p>Biến đầu ra (mục tiêu mong muốn):</p> <p>17 – y: khách hàng đã đăng kí kì hạn chưa? (nhị phân: "yes", "no")</p>
Thông tin tác giả	[Moro et al., 2014] S. Moro, P. Cortez and P. Rita. A Data-Driven Approach to Predict the Success of Bank Telemarketing. Decision Support Systems, Elsevier, 62:22-31, June 2014

Raw data: file bank-full.csv

Tidy data: được lưu lại thành file: bankfull.csv

Instruction list:

CODE R

```

1 rm(list=ls())
2
3 myfiles <- read.csv("dataset/bank-full.csv", header = FALSE, sep = ";")
4
5 variables <- c("age", "job", "marital", "education", "default", "balance", "housing", "loan",
6               | "contact", "day", "month", "duration", "campaign", "pdays", "previous", "poutcome", "y")
7
8 colnames(myfiles) <- variables
9
10 myfiles <- myfiles[ -c(1),]
11
12 write.csv(myfiles, file = "bankfull.csv", row.names = FALSE )
13
14
15

```

CODE BOOK Bài 2c

Thông tin	Nội dung
Tên bộ dữ liệu	Car Evaluation Database
Nguồn thu thập và cách thức thu thập	<p>(a) Creator: Marko Bohanec</p> <p>(b) Donors: Marko Bohanec (marko.bohanec@ijs.si) Blaz Zupan (blaz.zupan@ijs.si)</p> <p>(c) Date: June, 1997</p>
Số thuộc tính	6
Thông tin tên các thuộc tính	<p>1. buying v-high, high, med, low</p> <p>2. maint v-high, high, med, low</p> <p>3. doors 2, 3, 4, 5-more</p>

	4. persons 2, 4, more 5. lug_boot small, med, big 6. safety low, med, high
Thông tin tác giả	(a) Creator: Marko Bohanec (b) Donors: Marko Bohanec (marko.bohanec@ijs.si) Blaz Zupan (blaz.zupan@ijs.si)

Raw data: file car.data

Tidy data: được lưu lại thành file: car.csv

Instruction list:

CODE R

```

1 rm(list=ls())
2
3 myfiles <- read.csv("dataset/car.data", header = FALSE, sep = ",")
4
5 variables <- c("buying", "maint", "doors", "persons", "lug_boot", "safety")
6
7 colnames(myfiles) <- variables
8
9 myfiles <- myfiles[1:length(myfiles) - 1]
10
11 write.csv(myfiles, file = "car.csv", row.names = FALSE)
12
13 |
14

```

CODE BOOK Bài 2d

Thông tin	Nội dung
Tên bộ dữ liệu	Wine recognition data
Nguồn thu thập và cách thức thu thập	(a) Forina, M. et al, PARVUS - An Extendible Package for Data Exploration, Classification and Correlation. Institute of Pharmaceutical and Food Analysis and Technologies, Via Brigata Salerno, 16147 Genoa, Italy. (b) Stefan Aeberhard, email: stefan@coral.cs.jcu.edu.au (c) July 1991
Số thuộc tính	13
Thông tin tên các thuộc tính	1) Alcohol 2) Malic acid 3) Ash 4) Alcalinity of ash

	5) Magnesium 6) Total phenols 7) Flavanoids 8) Nonflavanoid phenols 9) Proanthocyanins 10) Color intensity 11) Hue 12) OD280/OD315 of diluted wines 13) Proline
Thông tin tác giả	(a) Creator: Marko Bohanec (b) Donors: Marko Bohanec (marko.bohanec@ijs.si) Blaz Zupan (blaz.zupan@ijs.si)

Raw data: file wine.data

Tidy data: được lưu lại thành file: wine.csv

Instruction list:

CODE R

```

1 rm(list=ls())
2
3 myfiles <- read.csv("dataset/wine.data", header = FALSE, sep = ",")
4
5 variables <- c("class identifier", "Alcohol", "Malic acid", "Ash", "Alcalinity of ash",
6               "Magnesium", "Total phenols", "Flavanoids", "Nonflavanoid phenols", "Proanthocyanins",
7               "Color intensity", "Hue", "OD280/OD315 of diluted wines", "Proline")
8
9 colnames(myfiles) <- variables
10
11 write.csv(myfiles, file = "wine.csv", row.names = FALSE)
12
13
14
15
16
17
18

```