

# ĐẠI HỌC QUỐC GIA THÀNH PHỐ HỒ CHÍ MINH



**UIT**  
Trường Đại học  
Công nghệ Thông tin

**Khoa Khoa học  
và Kỹ thuật Thông tin**

## CODEBOOK MÔ TẢ BỘ DỮ LIỆU DIABETES

Môn học: Thu thập và tiền xử lý dữ liệu - DS103.N21

Tên: Nguyễn Thị Huyền Trang

MSSV: 21520488

**TP HỒ CHÍ MINH, 2023**

# MỤC LỤC

<b>CODEBOOK</b>	3
Bộ dữ liệu	3
Codebook	3
Raw data	4
Tidy data	4
Instruction list	4
<b>THAM KHẢO</b>	5

# CODEBOOK

**Bộ dữ liệu:** Diabetes

**Codebook:**

Thông tin	Nội dung
Tên bộ dữ liệu	Diabetes patient records
Nguồn thu thập và cách thức thu thập	<p>Từ 2 nguồn:</p> <ol style="list-style-type: none"><li>Máy đo tự động: có đồng hồ bấm giờ tự động, sẽ ghi lại chính xác thời gian tại lúc đo số liệu.</li><li>Thu thập bằng tay: Sử dụng bản ghi giấy (paper record), giờ được định sẵn vào các khung giờ: sáng (8:00), trưa (12:00), chiều (18:00) và tối (22:00)</li></ol>
Số thuộc tính	4
Thông tin tên các thuộc tính	<p>Date: Ngày thu thập, định dạng: MM-DD-YYYY</p> <p>Time: Giờ thu thập, định dạng: XX:YY (24 giờ).</p> <p>Code: Mã code theo danh sách sau:</p> <ul style="list-style-type: none"><li>33 = Regular insulin dose</li><li>34 = NPH insulin dose</li><li>35 = UltraLente insulin dose</li><li>48 = Unspecified blood glucose measurement</li><li>57 = Unspecified blood glucose measurement</li><li>58 = Pre-breakfast blood glucose measurement</li><li>59 = Post-breakfast blood glucose measurement</li><li>60 = Pre-lunch blood glucose measurement</li><li>61 = Post-lunch blood glucose measurement</li><li>62 = Pre-supper blood glucose measurement</li><li>63 = Post-supper blood glucose measurement</li><li>64 = Pre-snack blood glucose measurement</li><li>65 = Hypoglycemic symptoms</li><li>66 = Typical meal ingestion</li><li>67 = More-than-usual meal ingestion</li><li>68 = Less-than-usual meal ingestion</li><li>69 = Typical exercise activity</li></ul>

	70 = More-than-usual exercise activity 71 = Less-than-usual exercise activity 72 = Unspecified special event Value: Giá trị thu thập được.
Thông tin tác giả	kahn@informatics.WUSTL.EDU (Internet) or 70333,34 (CompuServe)

### Raw data:

Raw data gồm tập hợp các file: data-01, data-02, ... data-70.

### Tidy data:

Tidy data sẽ được lưu lại thành file: diabetes.csv.

### Instruction list:

```
rm(list=ls())
myFiles <- list.files(path="diabetes-data", pattern="data")
data <- read.csv('diabetes-data/data-01', sep='\t',header = FALSE)
k = TRUE
# Tiến hành đọc từng file
for (f in myFiles) {
  if (k==TRUE) {
    file <- read.csv(paste("diabetes-data/", f, sep=""), sep="\t",header = FALSE)
    k = FALSE
  }
  else {
    file <- rbind(file, read.csv(paste("diabetes-data/", f,
                                     sep=""), sep="\t", header = FALSE))
  }
}
dataset <- file
variables <- c("Date", "Time", "Code", "Value")
# Đặt tên cho cột trong bộ dữ liệu
colnames(dataset) <- variables
write.csv(dataset, file = "diabetes.csv")
```

# THAM KHẢO

Diabetes Data Set, [UCI Machine Learning Repository: Diabetes Data Set](#)