

VIETNAM NATIONAL UNIVERSITY, HO CHI MINH CITY  
UNIVERSITY OF TECHNOLOGY  
FACULTY OF COMPUTER SCIENCE & ENGINEERING



## PROBABILITY AND STATISTICS

---

### THE EVOLUTION OF COMPUTER PROCESSORS: A STATISTIC OF COMMON PROPERTIES

---

Supervisor: Nguyen Thi Mong Ngoc, PhD  
Students: Chau Dang Minh - 2013748

Ho Chi Minh City, April 2024



## EVALUATION

N.O.	Student	ID	Works	Completed
1	Chau Dang Minh	2212287	Dataset overview Preprocessing	100%
2	Ha Khoi Nguyen	2212287	Descriptive statistics	100%
3	Nguyen Thi Mai Anh	2210103	Theories Slides and Presentation	
4			Inferential statistics	
5			Inferential statistics	



# Table of Contents

<b>I</b>	<b>Introduction</b>	<b>3</b>
<b>II</b>	<b>Overview of the Dataset</b>	<b>4</b>
<b>III</b>	<b>Preprocessing</b>	<b>5</b>
1	Data Cleaning . . . . .	5
2	Data Pre-computation . . . . .	7

# Chapter I

## Introduction

Phenomena that are meaningful to humans appear not to be stochastic. In the same sense, datasets produced by humans, or nature in time circulations have insights to be analyzed, which is accounted by Statistics. Thanks to Dr. Nguyen Thi Mong Ngoc's supervision in Probability and Statistics course, we have a chance to study basic statistics within an assignment with a tiny dataset. We organized our report in the following structure

1. Overview of the dataset. In this chapter, we carefully describe in details as much as possible the dataset, specifically the properties of each instance. We also notice which features to be used for later statistical tasks.
2. Preprocessing. We process data cleaning and some computations.
3. Descriptive statistics. We calculate some qualitative features of the dataset.
4. Inferential statistics. Our problems are explicitly stated and solved.

## Chapter II

# Overview of the Dataset

As Computer Science students, we are assigned to analyze a [dataset about computer processors](#), namely CPUs and GPUs. Our dataset is credited to Intel, Game-Debate, and the companies involved in producing the part. Information of CPUs and GPUs are collected separately into two files, namely `Intel_CPUs.csv` and `All_GPUs.csv`.

# Chapter III

## Preprocessing

### 1 Data Cleaning

With RStudio, the working directory is automatically determined. Otherwise, it can be indicated by `here` library.

Listing 3.1: Required libraries and working directory setup

```
1 # Libraries and options
2 library(dplyr)
3 library(here)
4 library(knitr)
5 library(kableExtra)
6
7 # Self-defined functions
8 source("utils.R")
9
10 # Working directory
11 setwd(here())
```

Now our working directory have been explicated, we can use relative paths to read the data. With RMarkdown, we can prettify the rendering.

Listing 3.2: RStudio data object initialization

```
1 # Read the CSV file into a data frame
2 cpu_data <- read.csv("dataset/Intel_CPUs.csv")
3 gpu_data <- read.csv("dataset/All_GPUs.csv")
4
5 # Inspect the CPU data
6 kable(head(cpu_data), format = "html") %>%
7   kable_styling()
```

Invalid cells may contain NA, an empty string, or other values showing us that this cell's data was not collecting correctly. At the very first step, we want to selected only columns whose the percentage of valid cells exceeds our predefined value. Then we filter out all instances with invalid features. Note that careful column selection possibly remains more instances for later tasks.

Listing 3.3: Cleaning functions

```
1 # Check if a cell has a valid value
2 is_valid <- function(value) {
3   return(!is.na(value)
4     & !is.null(value)
5     & !value == ""
6     & !value == "N/A"
7     & !value == "-")
```

Product_Collection	Vertical_Segment	Processor_Number	Status	Launch_Date	Lithography	Recommended_Customer_Price	nb_of_Cores	nb_of_Threads	Pro
7th Generation Intel® Core™ i7 Processors	Mobile	i7-7Y75	Launched	Q3'16	14 nm	\$393.00	2	4	1.3
8th Generation Intel® Core™ i5 Processors	Mobile	i5-8250U	Launched	Q3'17	14 nm	\$297.00	4	8	1.6
8th Generation Intel® Core™ i7 Processors	Mobile	i7-8550U	Launched	Q3'17	14 nm	\$409.00	4	8	1.8
Intel® Core™ X-series Processors	Desktop	i7-3820	End of Life	Q1'12	32 nm	\$305.00	4	8	3.6
7th Generation Intel® Core™ i5 Processors	Mobile	i5-7Y57	Launched	Q1'17	14 nm	\$281.00	2	4	1.2
Intel® Celeron® Processor 3000 Series	Mobile	3205U	Launched	Q1'15	14 nm	\$107.00	2	2	1.5

Figure III.1: First instances of CPUs data

```

8      & !value == "missing"
9      & !value == "unknown")
10     # Add your criteria
11   }
12
13   # Select columns with enough valid cells
14   filtered_data <- function(data, valid_percentage=0.8) {
15     selected_columns <- character(0)
16
17     for (col in colnames(data)) {
18       valid_count <- sum(is_valid(data[[col]]))
19       total_instances <- length(data[[col]])
20
21       if ((valid_count / total_instances) >= fill) {
22         selected_columns <- c(selected_columns, col)
23       }
24     }
25
26     return(data[selected_columns])
27   }

```

Listing 3.4: Cleaned data and selected features

```

1   filtered_cpu_data <- filtered_data(cpu_data, 0.4)
2
3   processed_cpu_data <-
4     filtered_cpu_data[
5       apply(filtered_cpu_data, 1, function(row) all(sapply(row, is_valid))), ]
6
7   selected_cpu_data <- processed_cpu_data[, c("Recommended_Customer_Price",
8       "Product_Collection",
9       "Launch_Date",
10      "nb_of_Cores",
11      "nb_of_Threads",
12      "Processor_Base_Frequency",
13      "Bus_Speed")]
14
15   # Adjust selected columns for your later needs
16
17   kable(head(selected_cpu_data), format = "html") %>%
18   kable_styling()

```

	Recommended_Customer_Price	Product_Collection	Launch_Date	nb_of_Cores	nb_of_Threads	Processor_Base_Frequency	Bus_Speed
1	\$393.00	7th Generation Intel® Core™ i7 Processors	Q3'16	2	4	1.30 GHz	4 GT/s OPI
2	\$297.00	8th Generation Intel® Core™ i5 Processors	Q3'17	4	8	1.60 GHz	4 GT/s OPI
3	\$409.00	8th Generation Intel® Core™ i7 Processors	Q3'17	4	8	1.80 GHz	4 GT/s OPI
5	\$281.00	7th Generation Intel® Core™ i5 Processors	Q1'17	2	4	1.20 GHz	4 GT/s OPI
423	\$281.00	7th Generation Intel® Core™ i5 Processors	Q3'16	2	4	1.20 GHz	4 GT/s OPI
424	\$281.00	7th Generation Intel® Core™ i5 Processors	Q3'16	2	4	2.50 GHz	4 GT/s OPI

**Figure III.2:** *First instances of selected CPUs data*

## 2 Data Pre-computation

Some features in our data have values that need to be reformatted for easily later sorting and analyses. Therefore, we need to gain a good understand on the features.

Listing 3.5: A processing for selected features

```
1 cpu_columns <- colnames(cpu_data)
2 gpu_columns <- colnames(gpu_data)
3 intersect(cpu_columns, gpu_columns)
4 # Output: character(0)
```

Since the data files have no common features, let us take a look at them independently.