

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



PROBABILITY AND STATISTICS

THE EVOLUTION OF COMPUTER PROCESSORS: A STATISTICS OF COMMON PROPERTIES

Instructor: Dr. Nguyen Thi Mong Ngoc
Students: Chau Dang Minh - 2013748
Ha Khoi Nguyen - 2212287
Nguyen Thi Mai Anh - 2210103
Võ Ninh Giang - 2210834
Trinh Viet Cuong - 2210447

Ho Chi Minh City, May 2024

Evaluation

N.O.	Student	ID	Works	Completed
1	Chau Dang Minh	2212287	Dataset overview Preprocessing Final Review	100%
2	Ha Khoi Nguyen	2212287	Descriptive statistics	100%
3	Nguyen Thi Mai Anh	2210103	Theories Slides and Presentation	
4	Võ Ninh Giang	2210834	Inferential statistics	
5	Trinh Viet Cuong	2210447	Inferential statistics	

Introduction and Acknowledgement

Phenomena that are meaningful to humans appear not to be completely 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.

Contents

I	Dataset Overview	4
II	Foundations	6
1	Theory of Sampling	6
1.1	Terminology	6
1.2	Statistical characteristics	6
2	Statistical Hypothesis Testing	8
3	Linear Regression	8
4	Analysis of Variance	9
4.1	ANOVA Usecases	10
4.2	Types of ANOVA	10
4.2.1	One-way ANOVA	10
4.2.2	Two-way ANOVA	10
III	Data Preprocessing	11
1	Data Cleaning	11
2	Data Pre-computation	13
IV	Descriptive Statistics	16
1	CPU Data	16
2	GPU Data	18
V	Inferential Statistics	26
1	Linear Regression	26
2	Two-way ANOVA	27
2.1	Problem Statement	27
2.2	Tests of Conditional Hypotheses	27
2.3	ANOVA Test	27
VI	Conclusion and Discussion	31

Chapter I

Dataset Overview

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`. Let us get familiar to some technical features, given in the metadata or acquired over the internet.

Table I.1: Some technical features of CPUs and GPUs

N.O.	Feature name	Relevant to	Details
1	Lithography	CPU, GPU	The semiconductor technology used to manufacture an integrated circuit, and is reported in nanometer
2	Number of Cores	CPU	A hardware term that describes the number of independent central processing units
3	Number of Threads	CPU	A Thread, or thread of execution, is a software term for the basic ordered sequence of instructions that can be passed
4	Base Frequency	CPU	Describes the rate at which the processor's transistors open and close.
5	Cache	CPU	An area of fast memory located on the processor.
6	Thermal Design Power	CPU	Represents the average power, in watts, the processor dissipates when operating at Base Frequency with all cores active under an Intel-defined, high-complexity workload.
7	Embedded Availability	CPU	In essence, an embedded processor is a CPU chip used in a system which is not a general-purpose workstation, laptop or desktop computer.
8	Embedded Availability	CPU	In essence, an embedded processor is a CPU chip used in a system which is not a general-purpose workstation, laptop or desktop computer.
9	Memory Types	CPU	Single Channel, Dual Channel, Triple Channel, and Flex Mode.

Continued on next page

Table I.1: Some technical features of CPUs and GPUs (Continued)

10	Instruction Set	CPU	
11	Maximal Temperature	CPU	
12	Architecture	GPU	
13	Dedicated and Integrated	GPU	Whether the GPU is solely used or shares memory with a CPU
14	(Front-side) Bus Speed	CPU, GPU	The speed at which data is transferred between the processors and other components such as the memory, chipset, and peripherals.
15	No-Execute Bit	CPU	Hardware-based security feature that can reduce exposure to viruses and malicious-code attacks.
16	Thermal Monitoring Technologies	CPU	Protects the processor package and the system from thermal failure through several thermal management features.

Chapter II

Foundations

A statistical Hypothesis is a statement or assumption about the parameters of one or more populations, which may be true or false. Each such truth value can be expressed in the form of a hypothesis, usually denoted as H_0 , called the null hypothesis, and an alternative hypothesis, usually denoted as H_1 or H_a , which is its complement.

1 Theory of Sampling

1.1 Terminology

1. A population is a set $\{x_i\}_{i=1}^N$ of elements belonging to the object under study, which need to be observed, collected, and analyzed based on certain characteristics. A parameter is a numerical description of a population.
2. A sample is a subset $\{x_{i_j}\}_{j=1}^n$ of units selected from the population according to a certain sampling method. A statistic is a numerical description of a sample.
3. Let the population belong to a distribution represented by a random variable X . A random sample is a collection $\{X_i\}_{i=1}^n$ such that X, X_1, \dots, X_n are independent and identically distributed.

1.2 Statistical characteristics

Statistical characteristics are important properties on a population, a sample or a random sample directly related to the content of the study and survey, which need to be collected from the units of the population, generally divided into two groups

1. Central tendency includes mean, median, mode, giving us an idea of the trend of the data points clustering around certain values.
2. Dispersion includes variance, standard deviation, which describe the spread or variability of the data points around the central tendency.

In details, these quantities are formulated as followings.

1. Population mean

$$\mu = \frac{1}{N} \sum_{i=1}^N x_i. \quad (2.1.1)$$

2. Sample mean

$$\bar{x} = \frac{1}{n} \sum_{j=1}^n x_{i_j}. \quad (2.1.2)$$

3. Random sample mean

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i. \quad (2.1.3)$$

4. Population variance

$$\sigma^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \mu)^2. \quad (2.1.4)$$

5. Sample variance

$$s^2 = \frac{1}{n-1} \sum_{j=1}^n (x_{i_j} - \bar{x})^2. \quad (2.1.5)$$

The quantity s^2 is also called corrected sample variance using Bessel's correction, because the denominator is $n-1$, instead of n . This is briefly explained as that since

$$\mu = \arg \min_y \sum_{i=1}^N (x_i - y)^2 \text{ and } \bar{x} = \arg \min_y \sum_{i=1}^n (x_{i_j} - y)^2,$$

we have

$$\begin{aligned} \frac{1}{n} \sum_{j=1}^n (x_{i_j} - \bar{x})^2 &\leq \frac{1}{n} \sum_{j=1}^n (x_{i_j} - \mu)^2 \\ &= \frac{1}{N} \sum_{j=1}^n (x_i - \mu)^2. \end{aligned} \quad (\text{assume that the sample is unbiased})$$

Hence, dividing by $n-1$ slightly increases the sample variance to better estimate the population variance.

6. Random sample variance

$$S^2 = \frac{1}{N-1} \sum_{n=1}^N (X_i - \bar{X})^2. \quad (2.1.6)$$

7. Population, sample and random sample derivations are the square roots σ , s and S , respectively.

8. Sample ratio

$$f = \frac{n}{N} \quad (2.1.7)$$

9. Coefficient of variation

$$CV = \left(\frac{s}{\bar{x}} \right) \cdot 100\% \quad (2.1.8)$$

10. Suppose that the sample $\{x_{i_j}\}_{j=1}^n$ is in increasing order. The sample median is calculated as

$$\text{med} = \begin{cases} x_{i_{k+1}}, & \text{if } n = 2k+1 \\ \frac{1}{2}(x_{i_k} + x_{i_{k+1}}), & \text{if } n = 2k. \end{cases} \quad (2.1.9)$$

11. Quartiles include three values Q_1, Q_2 and Q_3 . Where Q_2 is the median, Q_1 and Q_3 is the median of the set of samples less than and more than Q_2 , respectively. The quantity $\text{IQR} = Q_3 - Q_1$ is called the interquartile range.

12. An outlier is a point outside the range $(Q_1 - 1.5 \times \text{IQR}, Q_3 + 1.5 \times \text{IQR})$.

2 Statistical Hypothesis Testing

The process of testing a statistical hypothesis is a standardized procedure for making decisions to reject or not reject a hypothesis based on sample data. This process is called hypothesis testing and typically involves four steps:

1. **Setting Hypotheses** states the null hypothesis (H_0) and the alternative hypothesis (H_1). The hypotheses are formulated in a mutually exclusive manner, meaning if one is true, the other must be false.
2. **Constructing an Analysis Plan** outlines how sample data will be used to evaluate the null hypothesis. Evaluation typically centers around a single test statistic.
3. **Analyzing Sample Data** determines the values of the sample statistic (mean, proportion, t -statistic, z -score, etc.) as described in the analysis plan.
4. **Interpreting Results** applies decision rules outlined in the analysis plan. If the observed result is inconsistent with the null hypothesis, reject it.

Similar to estimation, statistical hypothesis testing does not provide results with 100% certainty but rather with a certain level of confidence, and errors can occur. Errors can be classified into two types:

- **Type I Error:** Rejecting the null hypothesis H_0 and accepting the alternative hypothesis H_1 while H_0 is true.
- **Type II Error:** Failing to reject the null hypothesis H_0 and not accepting the alternative hypothesis H_1 while H_1 is true.

Both types of errors can have adverse consequences. Depending on the situation, it's assessed which type of error leads to more serious consequences and should be minimized.

To make decisions about rejecting the null hypothesis, statisticians rely on specific rules. These rules are listed in the analysis plan. Traditionally, statisticians describe these decision rules in two ways: referencing a p -value or referencing a region of acceptance.

- **p -value:** The strength of evidence in favor of a null hypothesis is measured by the p -value. Suppose the test statistic is denoted by S . The p -value is the probability of observing a test statistic as extreme as S , assuming the null hypothesis is true. If the p -value is smaller than the significance level, we reject the null hypothesis.
- **Region of Acceptance:** The region of acceptance is a range of values. If the sample statistic falls within the region of acceptance, the null hypothesis is not rejected. The region of acceptance is set so that the probability of committing a Type I error is equivalent to the significance level.

Values outside the region of acceptance are called the region of rejection. If the sample statistic falls within the region of rejection, the null hypothesis is rejected. In such cases, it's said that the null hypothesis has been rejected at the significance level.

3 Linear Regression

Linear regression constructs a linear model between independent variables $\mathbf{x} = (0, x_1, \dots, x_d) \in \mathbb{R}^{d+1}$ and a dependent variable $y \in \mathbb{R}$. Traditionally, the relation is expressed as

$$y = w_0 + \sum_{i=1}^n w_i x_i = \mathbf{w}^\top \mathbf{x}, \quad (2.3.1)$$

where $\{w_i\}_{i=0}^d$ are parameters to learn. Let $\{(\mathbf{x}_n, \hat{y}_n)\}_{n=1}^N$ be collected observations. Our objective is to

minimize the error

$$\begin{aligned}
\mathcal{L}(\mathbf{w}) &= \sum_{n=1}^N (y - \hat{y})^2 \\
&= \sum_{n=1}^N (\mathbf{w}^\top \mathbf{x}_n - \hat{y}_n)^2 \\
&= \|\mathbf{X}\mathbf{w} - \hat{\mathbf{y}}\|^2,
\end{aligned} \tag{2.3.2}$$

where $\hat{\mathbf{y}} = (\hat{y}_1, \dots, \hat{y}_N)^\top \in \mathbb{R}^N$ and $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_N) \in \mathbb{R}^{N \times (d+1)}$. Taking the derivative of \mathcal{L} , we have

$$\frac{d\mathcal{L}(\mathbf{w})}{d\mathbf{w}} = 2\mathbf{X}^\top (\mathbf{X}\mathbf{w} - \hat{\mathbf{y}}).$$

Equivalently, $\mathbf{X}^\top \mathbf{X}\mathbf{w} = \mathbf{X}^\top \hat{\mathbf{y}}$. If $\mathbf{X}^\top \mathbf{X}$ is invertible i.e. the independent variables are linearly independent, we have the unique solution

$$\mathbf{w} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \hat{\mathbf{y}}. \tag{2.3.3}$$

Linear regression can be formulated extensively in a probabilistic way. Let $\epsilon \sim \mathcal{N}(0, \sigma^2)$ be a normal-distributed noise, we have the model

$$y = \mathbf{w}^\top \mathbf{x} + \epsilon \sim \mathcal{N}(\mathbf{w}^\top \mathbf{x}, \sigma^2). \tag{2.3.4}$$

Hence, for each data $\mathbf{x}_n, n = 1, \dots, N$, the probability density of \hat{y}_n given \mathbf{w} and σ^2 is

$$p(\hat{y}_n | \mathbf{x}_n, \mathbf{w}, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma}} \exp\left(-\frac{1}{2\sigma^2} (\hat{y}_n - \mathbf{w}^\top \mathbf{x}_n)^2\right). \tag{2.3.5}$$

Maximizing the likelihood $\prod_{i=1}^N p(\hat{y}_n | \mathbf{x}_n, \mathbf{w}, \sigma^2)$ yields the same solution \mathbf{w} as the original linear regression. Additionally,

$$\sigma^2 = \frac{1}{N} \sum_{i=1}^N (\hat{y}_n - \mathbf{w}^\top \mathbf{x}_n)^2. \tag{2.3.6}$$

In summary, some requirements need to be satisfied for an accurate and reliable solution.

1. The independent and dependent variables correlate. This means that the change in one or more independent variables induces the change in the dependent variable.
2. The observed data in the dataset are independent of each other. This means that the value of the dependent variable for one observed data point does not depend on the value of the dependent variable for another observed data point.
3. Across all levels of the independent variables, the variance of the errors remains constant. This indicates that the number of independent variables does not affect the variability of the errors.
4. The errors in the model follow a normal distribution.

4 Analysis of Variance

Analysis of Variance (ANOVA) is a statistical analysis tool used to partition the total observed variability found in a dataset into two parts: Systematic factors and Random factors. Systematic factors statistically influence a certain dataset, while random factors do not. We use ANOVA to determine the influence of independent variables on the dependent variable in regression studies.

The t -test and z -test methods, developed in the 20th century, were used for statistical analysis until 1918 when Ronald Fisher introduced the Analysis of Variance method. Both serve as hypothesis tests to assess whether there is any significant difference between means. The t -test is used when the population variance is unknown or the sample size is small ($n < 30$). Meanwhile, the z -test is applied when the population variance is known and the sample size is large ($n > 30$). The t -test uses the Student's t -distribution, while the z -test uses the standard normal distribution. As the sample size increases, the T -distribution converges to the standard normal distribution. The t -test can be understood as a statistical test used to compare and analyze whether the mean values of two populations are different when the standard deviation is unknown. Conversely, the z -test is a parametric test applied when the standard deviation is known to determine whether the means of two datasets are different. ANOVA is also known as Fisher's analysis of variance and is an extension of both t -tests and z -tests.

4.1 ANOVA Usecases

ANOVA test is the initial step in analyzing factors influencing a specific dataset. Upon completing a test, additional checks on measurable methodological factors contributing to dataset inconsistency are conducted. The results, using F -distribution, can be utilized to generate supplementary data suitable for proposed regression models. The test enables comparison of more than two groups simultaneously to determine whether a relationship exists among them. The outcome of ANOVA analysis, the F -statistic (also known as the F -ratio), allows for analysis across multiple datasets to ascertain variation between samples and within samples.

If no significant differences exist among the tested groups (known as the null hypothesis), the F -statistic result of ANOVA will be close to 1. The distribution of all possible values of the F -statistic follows the F -distribution. This is a family of probability distributions characterized by two parameters, known as the degrees of freedom of the numerator and the degrees of freedom of the denominator.

4.2 Types of ANOVA

There are two main types of ANOVA: one-way ANOVA and two-way ANOVA. Additionally, there are variations of ANOVA such as Multivariate ANOVA (MANOVA), which differs from ANOVA in that the former tests multiple dependent variables simultaneously while the latter evaluates only one dependent variable at a time. One-way or two-way refers to the number of independent variables in the analysis of variance. Furthermore, ANOVA relies on assumptions. ANOVA tests assume that the data are normally distributed, the variance is nearly equal across groups, and all observations are independent. If these assumptions are not met, ANOVA may not be useful for group comparisons.

4.2.1 One-way ANOVA

One-way ANOVA is used to assess the effect of a single factor on a single response variable. It determines whether all samples are similar i.e. determine whether there is a statistically significant difference between the means of three or more independent (unrelated) groups. One-way ANOVA is formulated as

$$Y_{i,j} = \mu_j + \varepsilon_{i,j}, \quad (2.4.1)$$

where $i = 1, \dots, n_j$ are indices of observations within the same group, $j = 1, \dots, k$ are indices of influenced groups, μ_j is the mean for each affected group and $\varepsilon_{i,j} \sim N(0, \sigma^2)$ is normal-distributed errors, independent for all i and j . Mean and variance of $Y_{i,j}$ are $E(Y_{i,j}) = \mu_j$ and $Var(Y_{i,j}) = \sigma^2$, respectively.

4.2.2 Two-way ANOVA

Two-way ANOVA is an extension of one-way ANOVA, where there are two independent factors. It is used to observe the interaction between two factors and test the effects of two factors simultaneously. It is formulated as

$$Y_{i,j} = \mu_j + \gamma_j + (\mu\gamma)_{i,j} + \varepsilon_{i,j}, \quad (2.4.2)$$

Assumptions about random variables are similar to one-way ANOVA.

Chapter III

Data Preprocessing

1 Data Cleaning

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

Listing 3.1.1: 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 3.1.2: 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 select only columns such that 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. For CPU data, we concern about the price, so let us choose a quota for which the column of prices is maintained.

Listing 3.1.3: Cleaning functions

```
1 # Check if a cell has a valid value
2 is_valid <- function(value) {
3   return(!is.na(value))
}
```

Product_Collection	Vertical_Segment	Processor_Number	Status	Launch_Date	Lithography	Recommended_Customer_Price	nb_of_Cores	nb_of_Threads	Pr
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 3.1.1: First instances of CPUs data

```

4      & !is.null(value)
5      & !value == ""
6      & !value == "N/A"
7      & !trimws(value) == "-"
8      & !value == "missing"
9      & !value == "unknown")
10     # Add your criteria
11   }
12
13
14   # Select columns with enough valid cells
15   filtered_data <- function(data, valid_percentage=0.5) {
16     selected_columns <- character(0)
17
18     for (col in colnames(data)) {
19       valid_count <- sum(is_valid(data[[col]]))
20       total_instances <- length(data[[col]])
21
22       if ((valid_count / total_instances) >= fill) {
23         selected_columns <- c(selected_columns, col)
24       }
25     }
26
27     return(data[selected_columns])
28   }

```

Listing 3.1.4: Cleaned CPU data

```

1   filtered_cpu_data <- filtered_data(cpu_data, valid_percentage=0.5)
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[] # Adjust selected columns for your later needs
8   selected_cpu_data <- unique(selected_cpu_data)
9   kable(head(selected_cpu_data), format = "html") %>%
10  kable_styling()

```

Listing 3.1.5: Cleaned GPU data

```

1 filtered_gpu_data <- filtered_data(gpu_data, valid_percentage=0.5)
2
3 processed_gpu_data <-
4   filtered_gpu_data[
5     apply(filtered_gpu_data, 1, function(row) all(sapply(row, is_valid))), ]
6
7 selected_gpu_data <- processed_gpu_data[] # Adjust selected columns for your later needs
8 selected_gpu_data <- unique(selected_gpu_data)
9 kable(head(selected_gpu_data), format = "html") %>%
10  kable_styling()

```

	Product_Collection	Vertical_Segment	Processor_Number	Status	Launch_Date	Lithography	Recommended_Customer_Price	nb_of_Cores	nb_of_Threads
1	7th Generation Intel® Core™ i7 Processors	Mobile	i7-7Y75	Launched	Q3'16	14 nm	\$393.00	2	4
2	8th Generation Intel® Core™ i5 Processors	Mobile	i5-8250U	Launched	Q3'17	14 nm	\$297.00	4	8
3	8th Generation Intel® Core™ i7 Processors	Mobile	i7-8550U	Launched	Q3'17	14 nm	\$409.00	4	8
5	7th Generation Intel® Core™ i5 Processors	Mobile	i5-7Y57	Launched	Q1'17	14 nm	\$281.00	2	4
11	Intel® Pentium® Processor 2000 Series	Mobile	2020M	Launched	Q3'12	22 nm	\$134.00	2	2
14	Intel® Pentium® Processor 4000 Series	Mobile	4405U	Launched	Q3'15	14 nm	\$161.00	2	4

Figure 3.1.2: First instances of selected CPUs data

It's worth finding the key of the data to know which feature uniquely determine an instance. We prioritize categorical features, specifically pure names.

Listing 3.1.6: Keys of the data

```

1 nrow(selected_cpu_data) - nrow(selected_cpu_data[, c("Product_Collection",
2   "Processor_Number")])
3 nrow(selected_gpu_data) - nrow(selected_gpu_data[, c("Name")])
4 # Outputs: 0

```

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 3.2.1: 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, we took a look at them independently and decided to pre-compute some features.

	Architecture	Best_Resolution	Core_Speed	DVI_Connection	Dedicated	Direct_X	HDMI_Connection	Integrated	L2_Cache	Manufacturer	Max_Power
2	R600 XT	1366 x 768	-	2	Yes	DX 10		0 No	0KB	AMD	215 Watts
3	R600 PRO	1366 x 768	-	2	Yes	DX 10		0 No	0KB	AMD	200 Watts
5	RV630	1024 x 768	-	2	Yes	DX 10		0 No	0KB	AMD	45 Watts
6	RV630	1024 x 768	-	2	Yes	DX 10		0 No	0KB	AMD	50 Watts
7	R700 RV790 XT	1920 x 1080	870 MHz	1	Yes	DX 10.1		1 No	0KB	AMD	190 Watts
8	R600 GT	1024 x 768	-	2	Yes	DX 10		0 No	0KB	AMD	150 Watts

Figure 3.1.3: First instances of selected GPUs data

1. Use the common column name `Release_Date` for both data. Extract `Release_Year` and `Release_Quarter` for each instance.
2. If `Recommended_Customer_Price` is a range, compute the average.

Listing 3.2.2: Extract Year and Quarter from Dates

```

1 month_to_quarter <- function(month) {
2   quarter <- switch(month,
3     "Jan" = "1",
4     "Feb" = "1",
5     "Mar" = "1",
6     "Apr" = "2",
7     "May" = "2",
8     "Jun" = "2",
9     "Jul" = "3",
10    "Aug" = "3",
11    "Sep" = "3",
12    "Oct" = "4",
13    "Nov" = "4",
14    "Dec" = "4",
15    "Unknown")
16   return(quarter)
17 }
18
19 names(selected_cpu_data)[names(selected_cpu_data) == "Launch_Date"] <- "Release_Date"
20 selected_cpu_data$Release_Year <- as.integer(sub("Q[1-4]'(\\d+)", "\\1",
21   gsub("\\s+", "", selected_cpu_data$Release_Date))) +
22   2000
23 selected_cpu_data$Release_Quarter <- as.integer(sub("Q([1-4])'.*", "\\1",
24   gsub("\\s+", "", selected_cpu_data$Release_Date)))
25
26 selected_gpu_data$Release_Year <-
27   as.integer(sub(".*-(\\d{4}) .*", "\\1", selected_gpu_data$Release_Date))
28 selected_gpu_data$Release_Quarter <-

```

```
28 as.character(sub(".*-(\\w+)-\\d{4}.*", "\\1", selected_gpu_data$Release_Date))
29 selected_gpu_data$Release_Quarter <- sapply(selected_gpu_data$Release_Quarter, month_to_quarter)
```


Chapter IV

Descriptive Statistics

1 CPU Data

```
1 litho_year <- CPUs_processed %>%
2   group_by(Release_Year) %>%
3   summarize(mean_litho = mean(Lithography),
4             median_litho = median(Lithography),
5             .groups = "drop"
6           )
7
8 ggplot(litho_year, aes(x = Release_Year)) +
9   geom_line(aes(y = mean_litho, color = "Mean")) +
10  geom_line(aes(y = median_litho, color = "Median")) +
11  scale_color_manual(values = c("Mean" = "blue", "Median" = "red")) +
12  labs(x = "Year", y = "Lithography", title = "Mean and Median Lithography by Year") +
13  scale_x_continuous(breaks = seq(min(litho_year$Release_Year), max(litho_year$Release_Year),
14                                by = 1)) +
15  theme_minimal()
```

```
1 options(repr.plot.width = 15, repr.plot.height = 8)
2 ggplot(data = CPUs_processed, aes(y = Status, x = Launch_Year, fill = Status)) +
3   geom_boxplot() +
4   labs(x = "Status", y = "Launch Date", title = "Boxplot of Status over Year") +
5   scale_x_continuous(breaks = seq(min(Litho_year$Launch_Year), max(Litho_year$Launch_Year), by
6                                 = 1))
```

```
1 options(repr.plot.width = 15, repr.plot.height = 8)
2 ggplot(data = CPUs_processed, aes(y = Vertical_Segment, x = Launch_Year, fill =
3   Vertical_Segment)) +
4   geom_boxplot() +
5   labs(x = "Year", y = "Vertical Segment", title = "Boxplot of Vertical Segment over Year") +
6   scale_x_continuous(breaks = seq(min(Litho_year$Launch_Year), max(Litho_year$Launch_Year), by
7                                 = 1))
```

```
1 options(repr.plot.width = 15, repr.plot.height = 8)
2 ggplot(data = CPUs_processed, aes(y = Instruction_Set, x = Launch_Year, fill =
3   Instruction_Set)) +
4   geom_violin() +
5   labs(x = "Launch Date", y = "Instruction set", title = "Violinplot of Instruction set over
6   Year") +
7   scale_x_continuous(breaks = seq(min(Litho_year$Launch_Year), max(Litho_year$Launch_Year), by
```

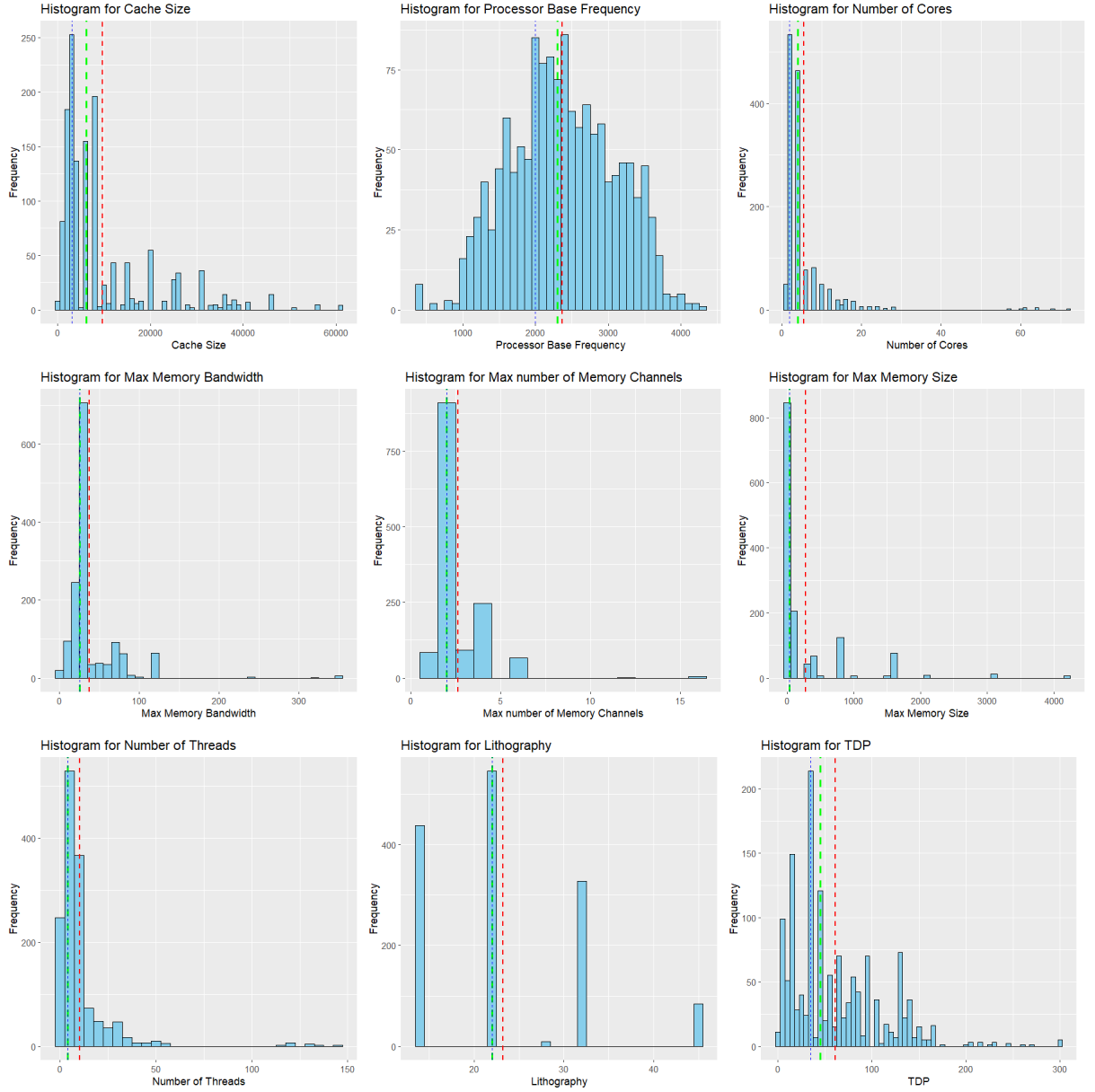


Figure 4.1.1: Histograms of several features in CPU data

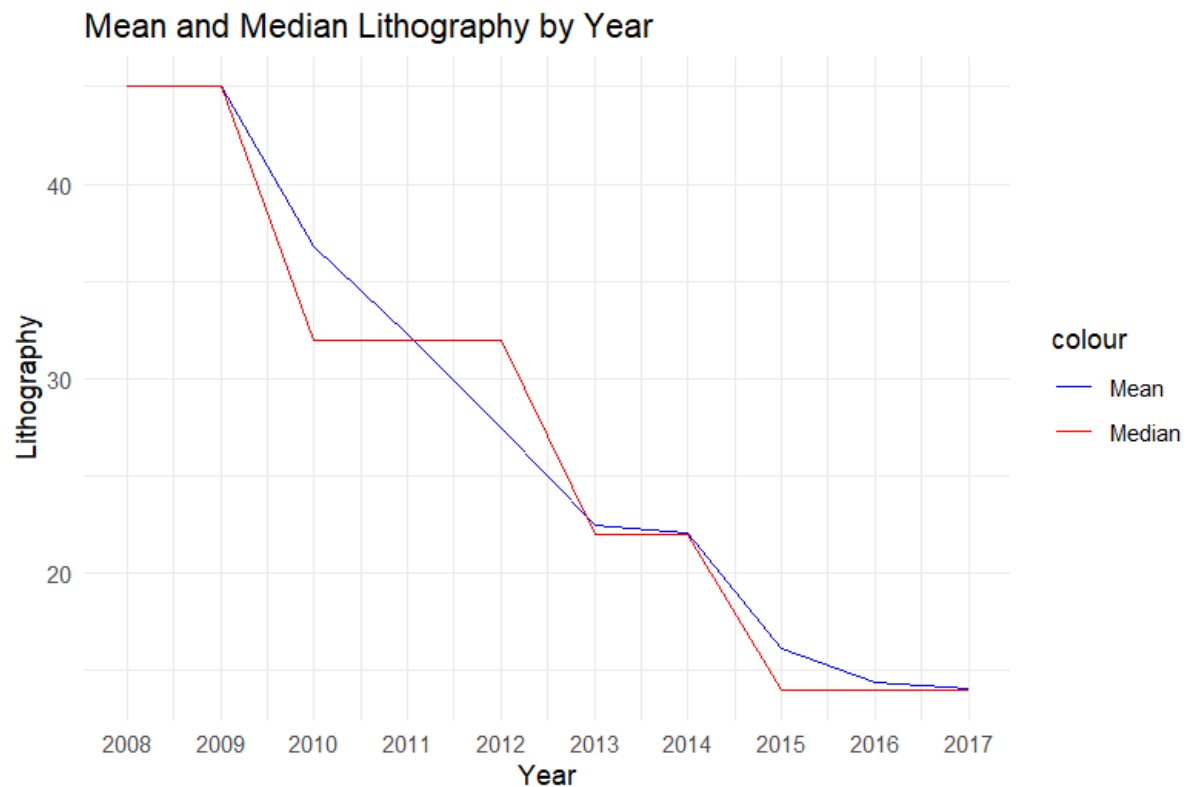


Figure 4.1.2: Mean and median of lithography over years

```
= 1))
```

2 GPU Data

```
1 freq <- table(GPUs_processed$Manufacturer, GPUs_processed$Release_Year)
2 total_count <- colSums(freq)
3 percentage <- prop.table(freq, margin = 2) * 100
4
5 barplot(freq,
6         legend.text = TRUE,
7         main = "Counts of Each Year",
8         xlab = "Year",
9         ylab = "Count",
10        col = c("skyblue", "salmon", "lightgreen", "yellow"),
11        border = "black")
```

```
1 barplot(percentage,
2         legend.text = TRUE,
3         main = "Market Share Percentage by Year",
4         xlab = "Year",
5         ylab = "Percentage",
6         col = c("skyblue", "salmon", "lightgreen", "yellow"),
7         border = "black")
```

```
1 scatter_plot <-
2   ggplot(GPUs_processed, aes(x = Release_Year + Release_Month/12, y = Memory, color =
3     Manufacturer)) +
4   geom_point() +
```

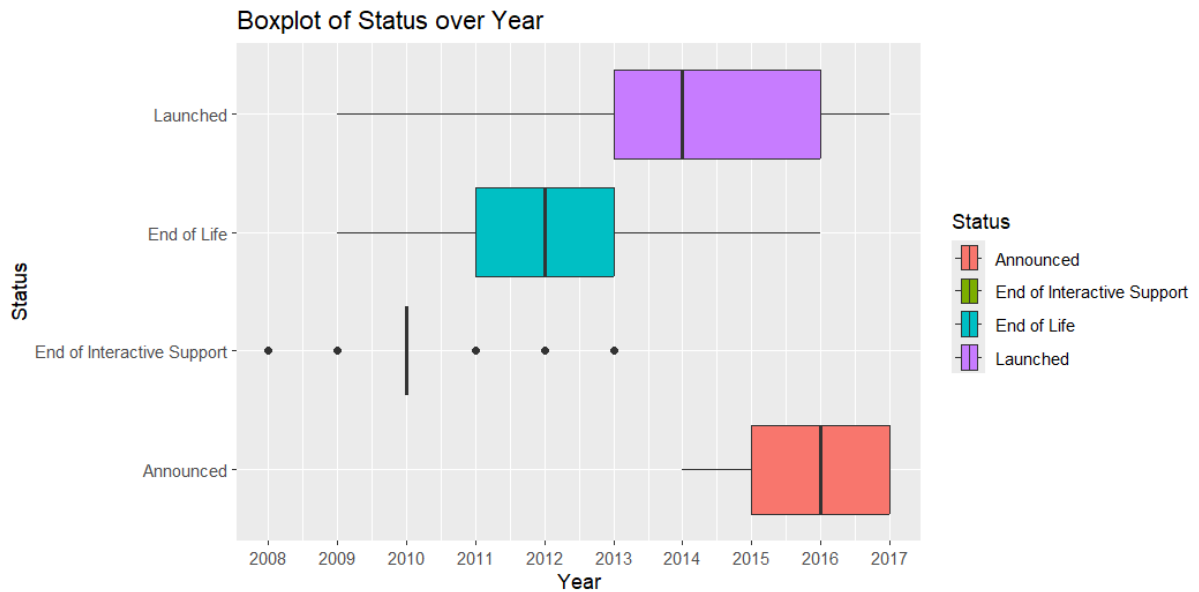


Figure 4.1.3: Status over years

```
4 scale_color_manual(values = c("skyblue", "salmon", "lightgreen", "yellow")) +
5 labs(x = "Year", y = "GPU Memory", title = "Scatter Plot of GPU Memory over Years") +
6 theme_minimal()
7 print(scatter_plot)
```

```
1 memory_summary <- GPUs_processed %>%
2   group_by(Release_Year) %>%
3   summarise(mean_memory = mean(Memory),
4             median_memory = median(Memory))
5
6 line_plot <- ggplot(memory_summary, aes(x = Release_Year)) +
7   geom_line(aes(y = mean_memory, color = "Mean")) +
8   geom_line(aes(y = median_memory, color = "Median")) +
9   scale_color_manual(values = c("Mean" = "blue", "Median" = "red")) +
10  labs(x = "Year", y = "Memory", title = "Mean and Median Memory by Year") +
11  theme_minimal()
12 print(line_plot)
```

```
1 moore_law <- data.frame(
2   Release_Year = seq(min(memory_summary$Release_Year), max(memory_summary$Release_Year), 1),
3   Memory = 2^(0.5 * (seq(min(memory_summary$Release_Year), max(memory_summary$Release_Year), 1)
4     - min(memory_summary$Release_Year - 8)))
5 )
6 memory_summary$log_mean_memory <- log(memory_summary$mean_memory)
7 memory_summary$log_median_memory <- log(memory_summary$median_memory)
8
9 line_plot <- ggplot(memory_summary, aes(x = Release_Year)) +
10  geom_line(aes(y = log_mean_memory, color = "Logarithm of Mean"), size = 1) +
11  geom_line(aes(y = log_median_memory, color = "Logarithm of Median"), size = 1) +
12  geom_line(data = moore_law, aes(y = log(Memory), color = "Moore's Law"), size = 1, linetype =
13    "dashed") +
14  scale_color_manual(values = c("Logarithm of Mean" = "blue", "Logarithm of Median" = "red",
15    "Moore's Law" = "green4")) +
16  labs(x = "Year", y = "Logarithm of Memory", title = "Logarithm of Mean and Median Memory by
17    Year") +
18  theme_minimal()
```

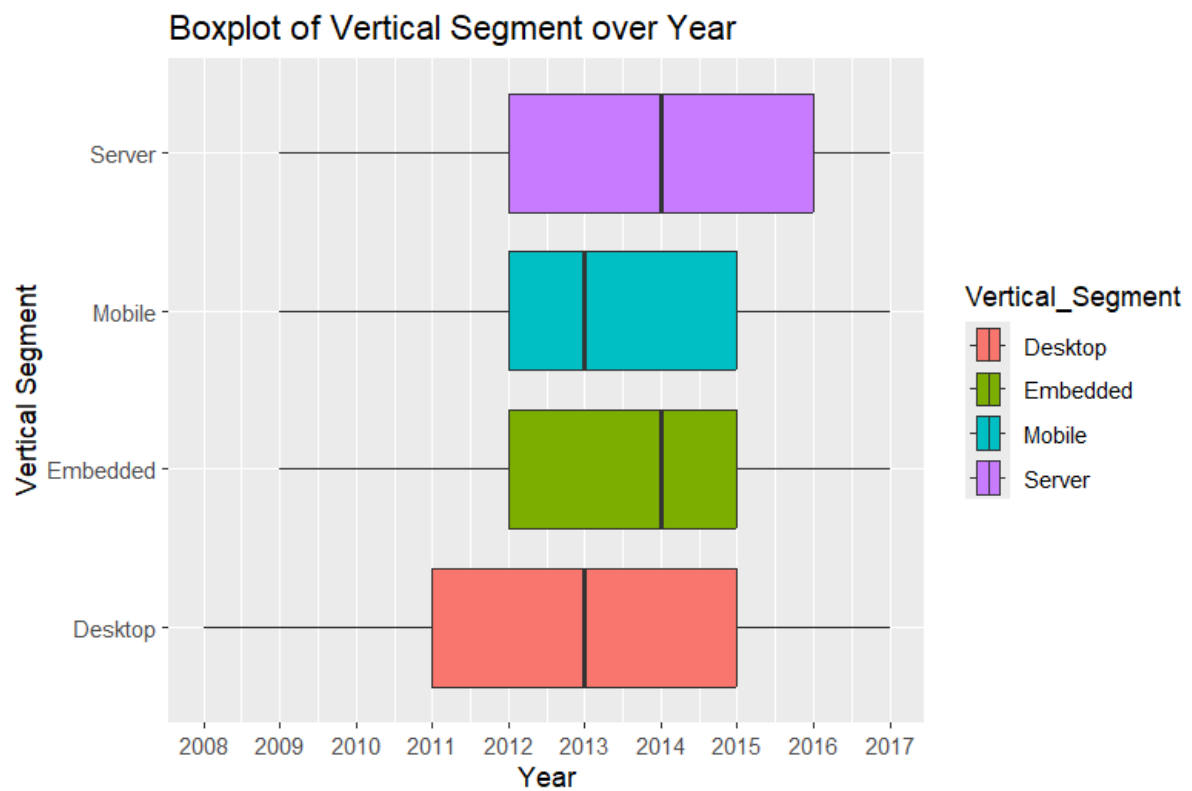


Figure 4.1.4: Segment over years

```
16  
17 print(line_plot)
```

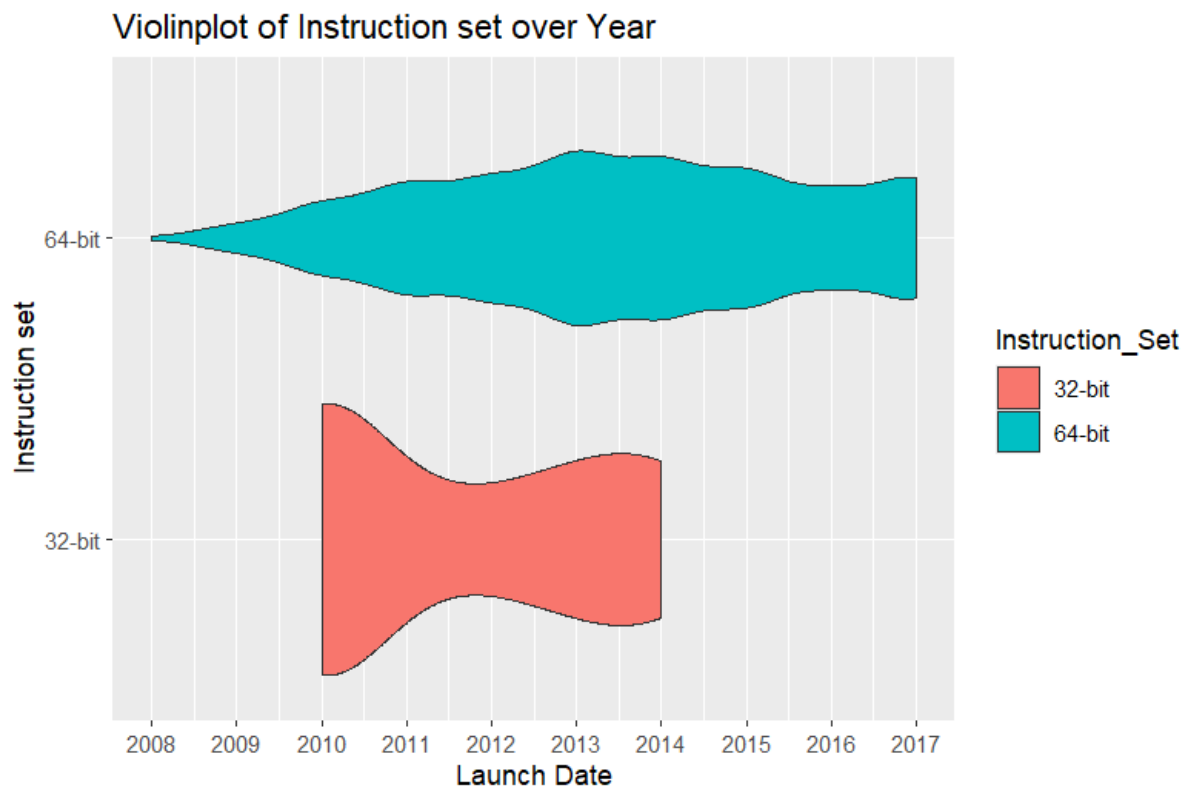


Figure 4.1.5: Instruction set over years

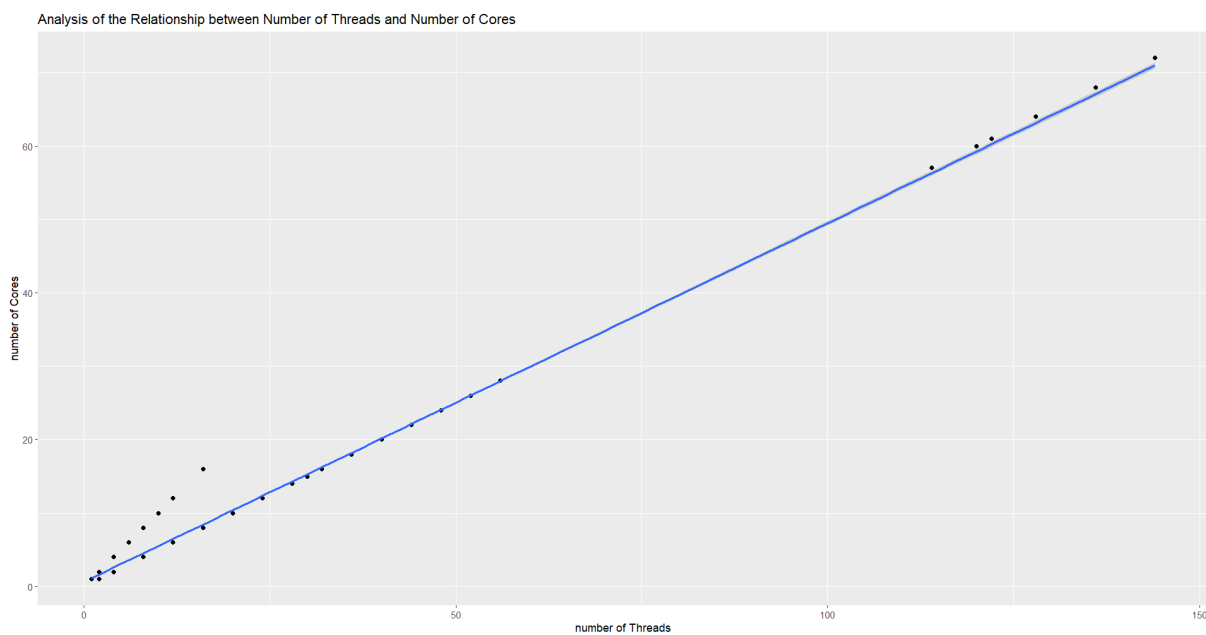


Figure 4.1.6: Relation between Number of Threads and Number of Cores

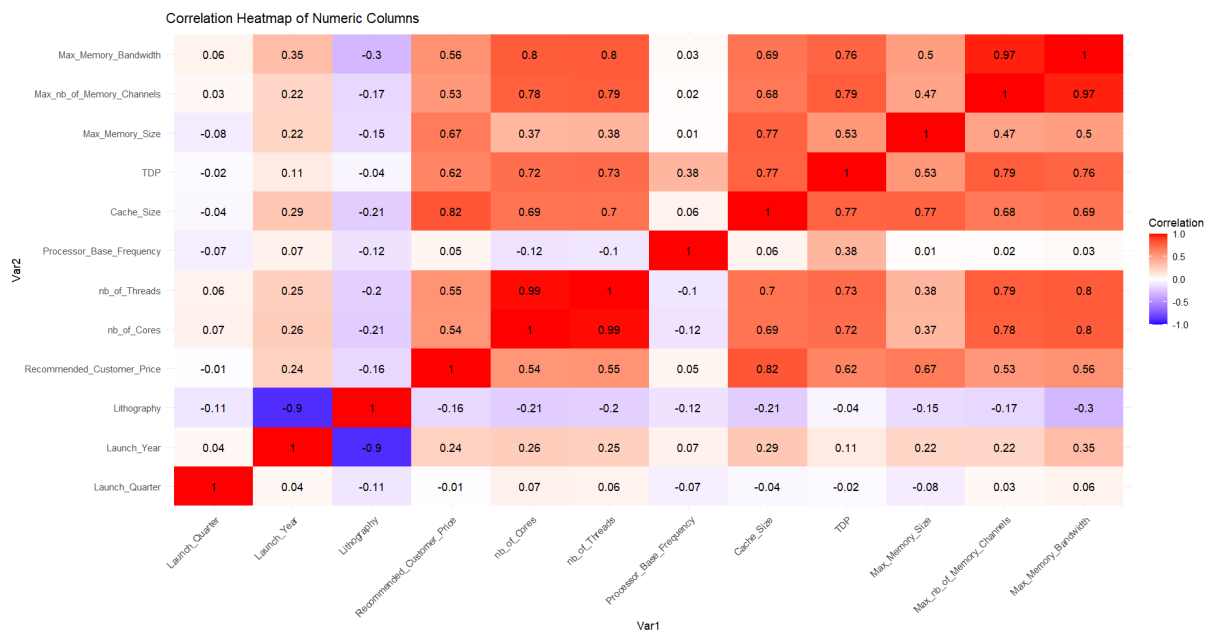


Figure 4.1.7: Overall Correlation

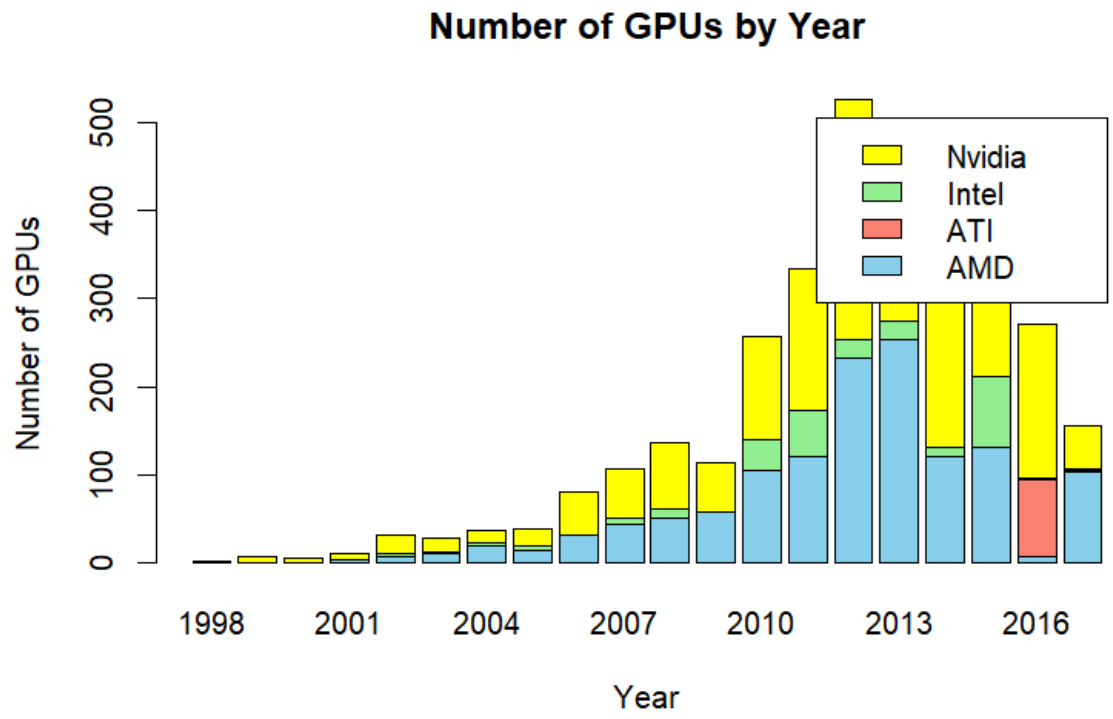


Figure 4.2.1: Sales of GPU manufacturer over years

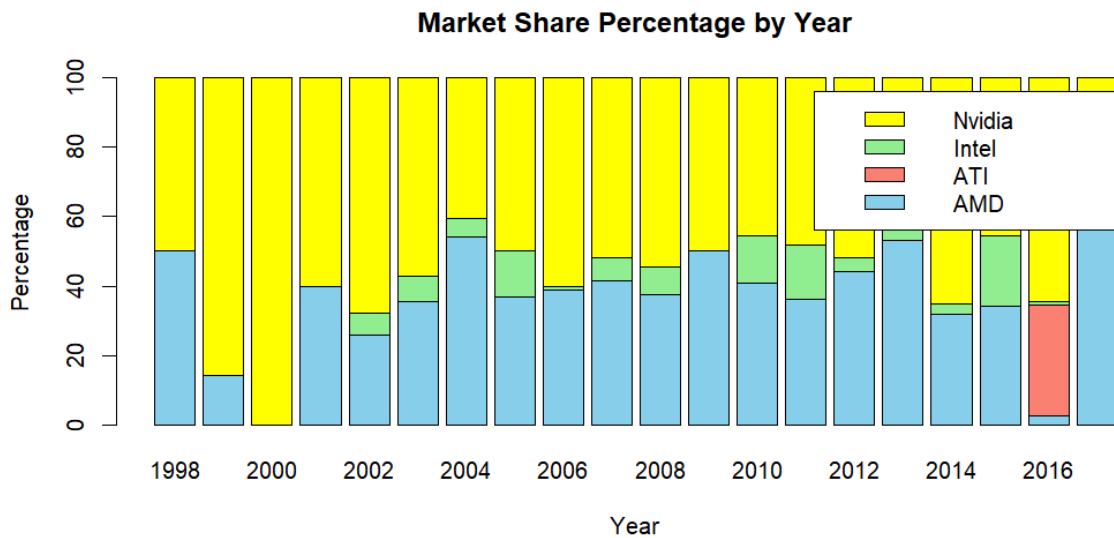


Figure 4.2.2: Market Share of GPU manufacturer over years

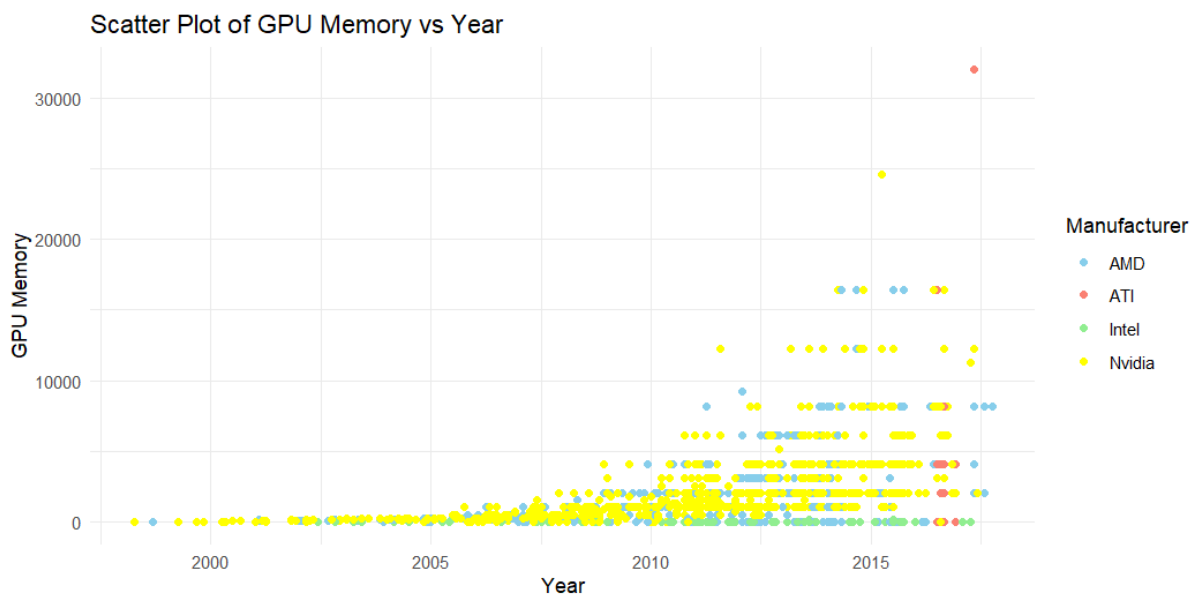


Figure 4.2.3: GPU Memory over Years

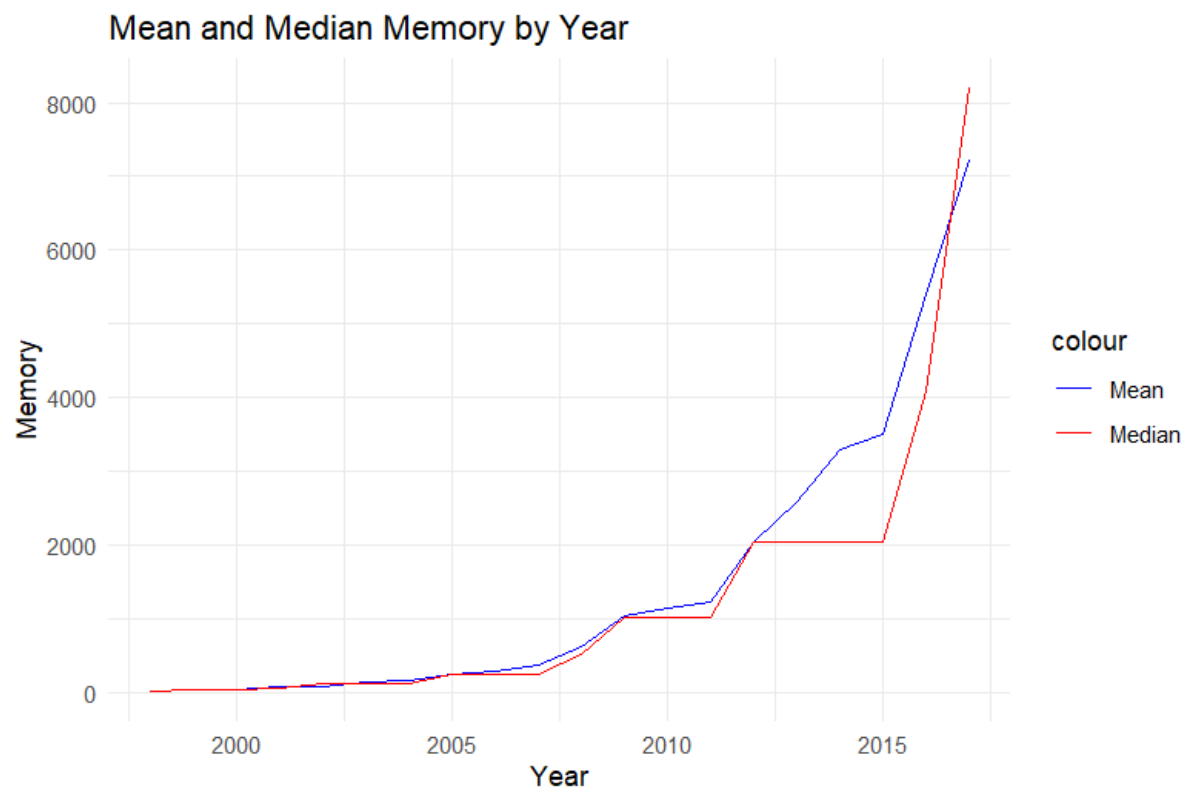


Figure 4.2.4: Mean and median of memory over years

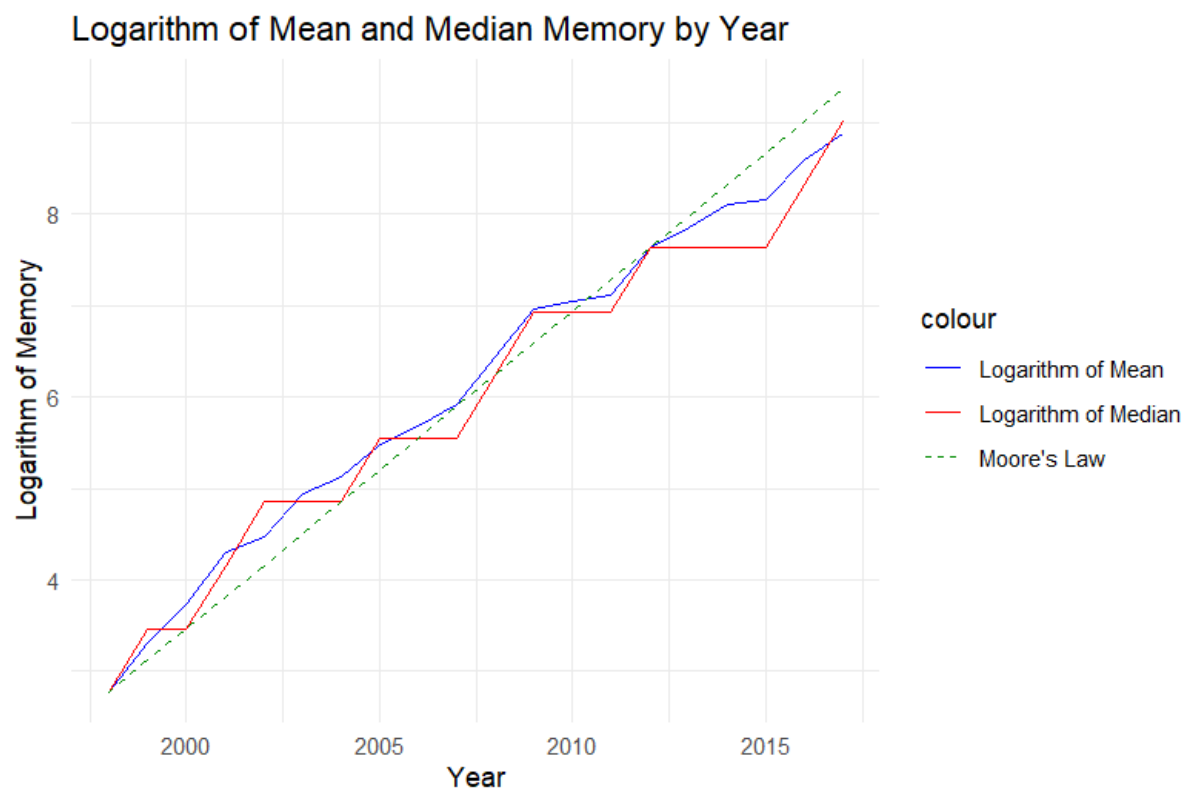


Figure 4.2.5: Logarithm of Mean and Median Memory by Year

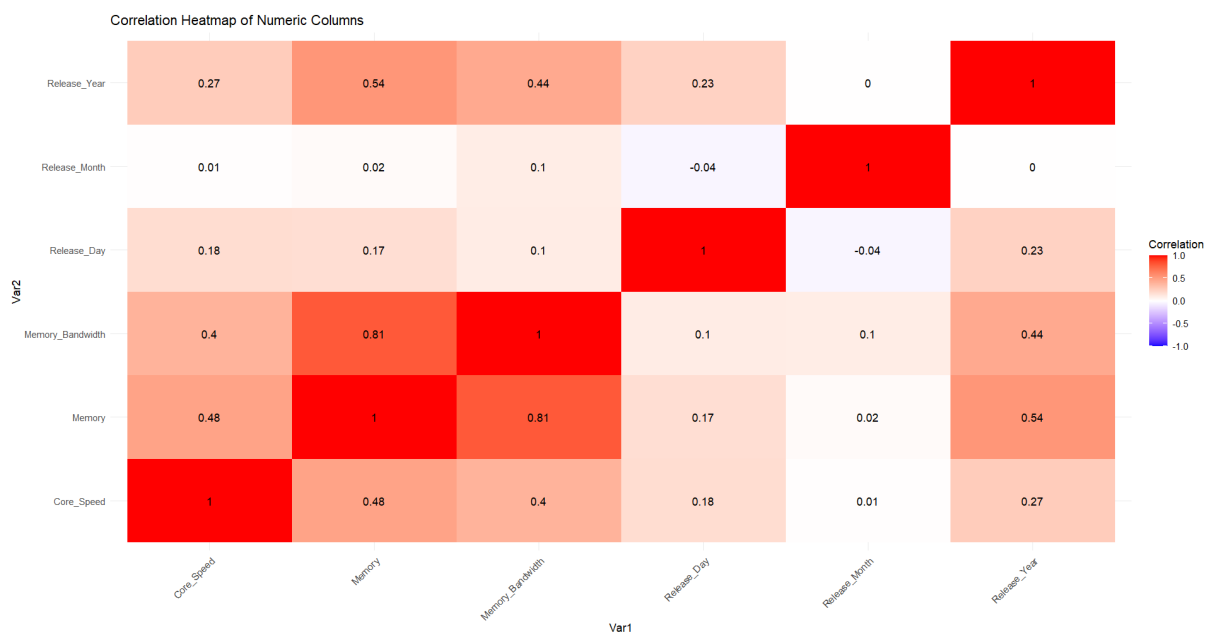


Figure 4.2.6: Enter Caption

Chapter V

Inferential Statistics

1 Linear Regression

We will try to infer if there is a linear relation between the price and other features. Since the percentage of GPUs having price is low (16.3%), we work on the CPU data. As illustrated in Figure 5.1.1, linear relations can be seen through Max_Memory_Bandwidth, Cache, Lithography and nb_of_Cores towards Recommended_Customer_Price. Let us see how well a linear model can fit the price based on these features.

```
1 CPUs_LR_input <- get_valid(CPUs_processed[, c(features, "Recommended_Customer_Price")])
2
3
4 # Convert Vertical_Segment to a factor with specified levels
5 CPUs_LR_input$Vertical_Segment <- factor(CPUs_LR_input$Vertical_Segment,
6     levels = c("Embedded", "Mobile", "Desktop", "Server"),
7     labels = c(1, 2, 3, 4))
8
9 CPUs_LR_input$Vertical_Segment <- as.numeric(as.character(CPUs_LR_input$Vertical_Segment))
10
11
12 lr_CPUs <- lm(Recommended_Customer_Price ~
13     Max_Memory_Bandwidth*
14     Cache*
15     Lithography*
16     nb_of_Cores, data = CPUs_LR_input)
17
18 summary(lr_CPUs)
19
20 # Compute RMSE
21 RMSE <- sqrt(mean(residuals(lr_CPUs)^2))
22 print(RMSE)
23
24 # Output: 435.6237
```

In Figure 5.1.2, the coefficients of Cache and Lithography are positive and negative, respectively, matching our intuition that larger cache and smaller lithographic technology cost more. The increase of nb_of_Cores does not lead to higher price can be interpreted as a consequence of Moore's law.

2 Two-way ANOVA

2.1 Problem Statement

Many customers believe that the product lines in each applied segment will have different numbers of cores, so perform ANOVA testing to check if there is a relationship between `Product_Collection` and `Vertical_Segment` to `nb_of_Cores`.

```

1 CPUs_ANOVA_input <-
2   get_valid(
3     CPUs_selected %>%
4     filter(Product_Collection %in% c("Core", "Legacy") &
5           Vertical_Segment %in% c("Desktop", "Embedded", "Mobile")) %>%
6     select(Product_Collection, Vertical_Segment, nb_of_Cores)
7   )

```

2.2 Tests of Conditional Hypotheses

Before ANOVA test, let us test whether initial conditions are satisfied.

1. Testing whether `nb_of_Cores` is normal-distributed (H_0) or not (H_1). Since $p_{\text{value}} < 0.05$, we reject H_0 i.e. `nb_of_Cores` is not normal-distributed.

```

1 library(nortest)
2 shapiro.test(CPUs_ANOVA_input$nb_of_Cores)
3
4 # Output: Shapiro-Wilk normality test
5
6 # data: CPUs_ANOVA_input$nb_of_Cores
7 # W = 0.62213, p-value < 2.2e-16

```

2. Testing whether the variances are equal (H_0) or there exist two groups of different variances (H_1). Since $p_{\text{value}} = 1.339e - 07 < 0.05$, we reject H_0 .

```

1 library(car)
2 leveneTest(CPUs_ANOVA_input$nb_of_Cores ~
3           CPUs_ANOVA_input$Product_Collection *
4           CPUs_ANOVA_input$Vertical_Segment)
5
6 # Output
7 # Levene's Test for Homogeneity of Variance (center = median)
8 # Df F value Pr(>F)
9 # group 5 8.2835 1.339e-07 ***
10 # 645

```

Thus, `nb_of_Cores` is not normally distributed, and the variances of the groups are not equal. However, the sample size is large (more than 30). Therefore, we can disregard these violations.

2.3 ANOVA Test

Our objective is to examine the impact of the two factors `Product_Collection` and `Vertical_Segment` on `nb_of_Cores`. Since this is a large sample, even though the sample violates the assumptions of normal distribution and homogeneity of variances, the two-factor ANOVA method can still be applied. However, the results are only for reference. For more accurate results, other methods can be used.

```

1 av <- aov(nb_of_Cores ~
2           Product_Collection *
3           Vertical_Segment,
4           data = CPUs_ANOVA_input)

```

```

5 summary(av)
6 # Output
7 #
8 #Product_Collection      1  94.6  94.64 46.895 1.75e-11 ***
9 #Vertical_Segment        2 198.2  99.08 49.093 < 2e-16 ***
10 #Product_Collection:Vertical_Segment 2 27.4  13.71  6.795  0.0012 **
11 #Residuals               645 1301.7  2.02

```

From the mean squares, we see that the effects of `Product_Collection` and `Vertical_Segment` appear to be similar, and both effects are statistically significant, as the p -values for both factors are very low. Specifically, for `Product_Collection`, the hypotheses are the average number of CPU cores in the two product types, Core and Legacy, are the same (H_0) or different (H_1). For `Vertical_Segment`, they are the average number of CPU cores in the application segments Desktop, Embedded, and Mobile are the same (H_0) or at least two segments have different average numbers of CPU cores (H_1). Both H_0 hypotheses are rejected. We conclude that the average number of CPU cores is different for different product lines (`Product_Collection`) and application segments (`Vertical_Segment`), so the average number of CPU cores depends on the product line and application segment.

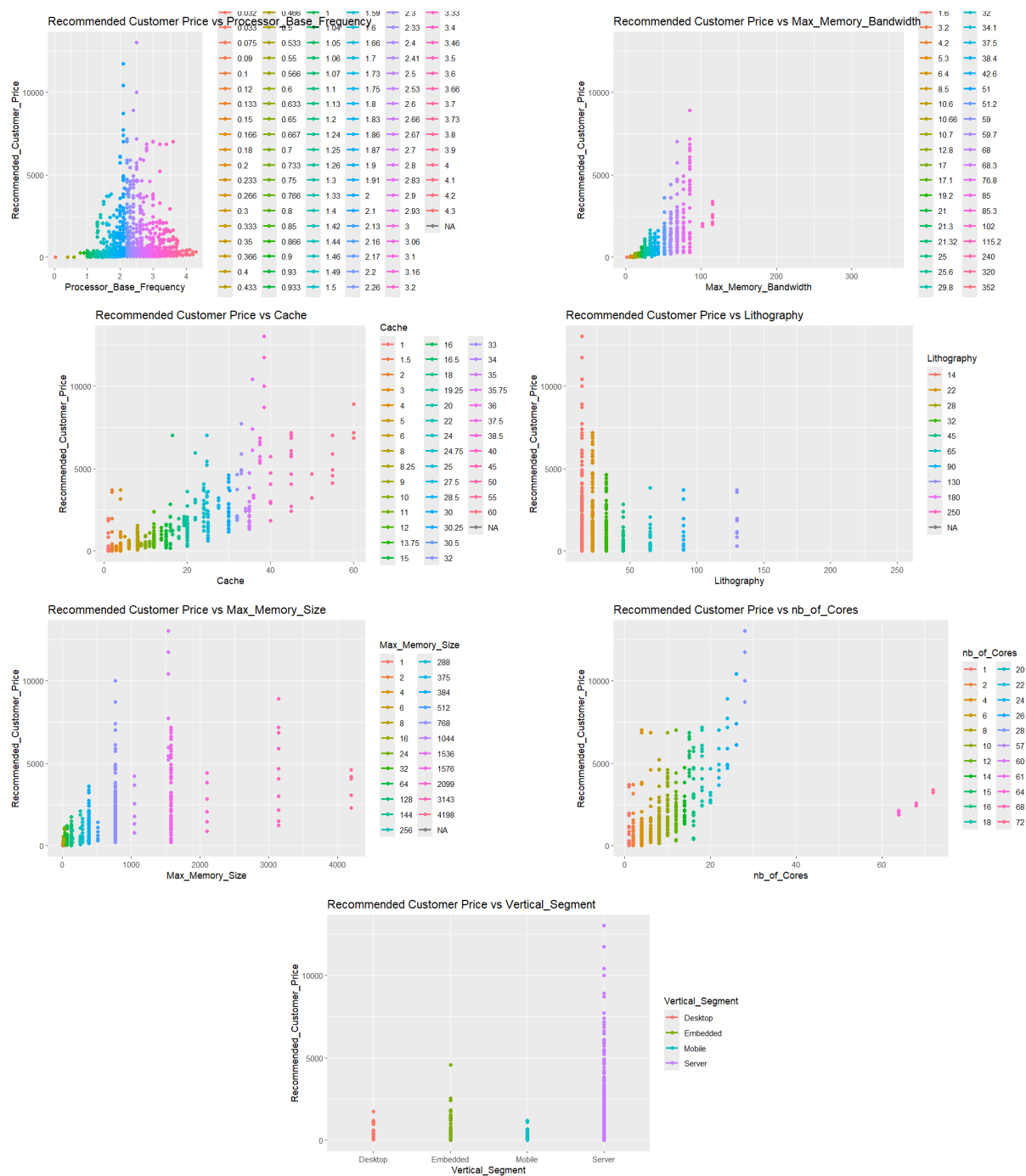


Figure 5.1.1: Plots of several features against the price

Coefficients:

(Intercept)

Max_Memory_Bandwidth

Cache

Lithography

nb_of_Cores

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	1.210e+03	3.164e+02	3.823	0.000
Max_Memory_Bandwidth	-1.115e+01	8.670e+00	-1.286	0.199
Cache	4.479e+01	4.121e+01	1.087	0.282
Lithography	-3.937e+01	1.386e+01	-2.841	0.004
nb_of_Cores	-3.640e+02	1.222e+02	-2.978	0.003

Figure 5.1.2: Linear regression summary

```
> table(Intel_clean$Product_Collection, Intel_clean$Vertical_Segment)
```

	Desktop	Embedded	Mobile	Server
Atom	3	22	48	26
Celeron	27	8	44	0
Core	118	32	165	0
Legacy	119	41	168	56
Pentium	42	3	27	6
Quark	0	8	0	0
Xeon	0	47	7	379

Figure 5.2.1: Data Partition based on Product_Collection and Vertical_Segment

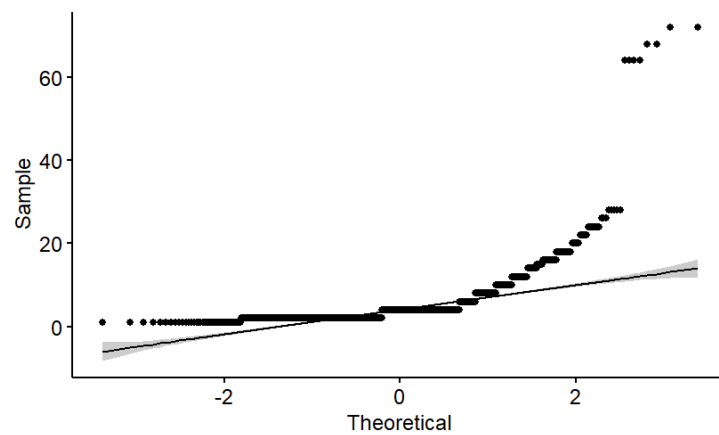


Figure 5.2.2: Q-Q plot returned by `ggqqplot(nb_of_Cores)` function. Observed values are far away from normal expectation.

Chapter VI

Conclusion and Discussion

Through the project, we have reviewed our knowledge on probability and statistics, particularly ANOVA and Linear Regression. Also, we have an opportunity to enhance our practice skills with R language. We did not work much on GPU data. The major reason is that many features are not familiar to our studying fields but relevant to computer graphics. There may be many insights can be dug out, which are our intention in the future.