

Relationship of clinical and sociodemographic variables between normotensive and prehypertensive groups

ABSTRACT

Hypertension or high blood pressure is a cardiovascular health condition characterized by elevated blood pressure in patients. Nevertheless, its diagnosis is problematic since this health condition generally does not produce symptoms until it reaches a higher stage. On the other hand, the adoption of emerging technologies has opened the ability to collect data in a more agile way and in different industries such as health services to improve the welfare of society. Therefore, efficient and faster decision-making has become necessary. Taking advantage of this data collection, it is possible to analyze health conditions that afflict people, such as hypertension, whose inadequate treatment and late diagnosis could worsen the health of people chronically. This work presents a statistical analysis of the relationship between sociodemographic and clinical variables with hypertension in its early stages. The data used for this study was collected at the Guilin People's Hospital in Guangxi, China. The data contains the weight, gender, body mass index, heart rate, blood pressure, and hypertension stage. The variables were studied through four statistical hypotheses, in which the difference between normotension and prehypertensive subjects was studied. The Bayesian and Fisherian statistical inference techniques were used to analyze the association between the variables with early hypertension stages by considering the data distribution, which was log-normal instead of making normal distribution assumptions like the classical Frequentist inference. For the decision of the Bayesian approach, the overlapping of the posterior distribution curves of the compared groups was considered to not reject the null hypothesis. For the Fisherian approach, the maximum likelihood estimator and variance of the parameters were used to draw the normal curves of the mean of each group, the not rejection of the statistical hypothesis was based on their overlapping. The results obtained by applying both

techniques demonstrate that gender, body mass index, heart rate, and weight do not provide information related to detecting hypertension in its early stages. These results highlight a gap in data collection from clinical variables that could alert the presence of hypertension in its early stages, which can be addressed in future studies.

Keywords

Prehypertension

Sociodemographic information

Clinical information

Bayesian inference

Fisherian approach

1. Introduction

Hypertension or high blood pressure is a cardiovascular health condition characterized by high arterial blood pressure levels in the patient. The diagnosis of hypertension is performed by measuring a subject's systolic and diastolic blood pressure through a sphygmomanometer, the gold standard for blood pressure measuring [1]. Systolic blood pressure refers to the force generated by the blood when the heartbeat squeezes and pushes the blood against the arteries. On the other hand, diastolic blood pressure refers to the force exerts into the arteries when the heart is filling. Based on these readings, a subject is classified into four stages of hypertension, as shown in Table 1 [2].

According to the World Health Organization (WHO), this health condition affects 1.1 billion individuals globally. Its inadequate management could produce long-term diseases such as heart attack, heart failure, kidney disease, coronary heart disease, diabetes, or strokes [3]. Besides, in 2013, based on the WHO reports, high blood pressure was responsible for at least 45% of heart disease deaths and 51% of stroke deaths. The above factors have led to consider this health condition a global public health issue [3]. Early detection of hypertension, lifestyle changes, and strict control, according to medical research, could limit its progression and effects. Nevertheless,

hypertension rarely generates noticeable symptoms that could alert from its presence in an early stage which difficult its diagnose and treat [4]. Furthermore, managing this medical condition is difficult and expensive.

Table 1 Blood pressure classification.

Blood Pressure Classification	Diastolic Blood Pressure (mmHg)	Systolic Blood Pressure (mmHg)
Normal	<80	<120
Prehypertension	80-89	120-139
Stage I hypertension	90-99	140-159
Stage II hypertension	Equal or greater than 100	Equal or greater than 160

On the other hand, the Industry 4.0 concept has been accepted as a goal within multiple industry sectors [5], such as healthcare services, to provide a better quality of life [6]. The concept of industry 4.0 used in the healthcare industry lies in technologies that allow the collection of information from patients more efficiently, data processing in increasingly shorter times, and decision-making in real-time [7]. This data is gathered by health centers or hospitals into what is known as electronic health records (EHRs), where physiological signals and clinical and sociodemographic information from patients are available [8]. However, the problem with extensive data collection is how to process it and interpret it for decision-making.

Therefore, considering the data for the decision-making available through EHRs in healthcare services and the health concerns that high-blood pressure generates, this study presents a statistical inference analysis of the relationship between sociodemographic and clinical variables with the stages of hypertension. The above is done to identify risk factors associated with this health condition in its early stages, concretely studying the difference between normotensive and prehypertensive subjects. The data used to carry out this analysis was collected from patients admitted at the Guilin People's Hospital in Guanxi, China [9].

Two statistical inference methods were proposed to decide between four statistical hypotheses related to gender, body mass index (BMI), heart rate, weight, and their association with normal and prehypertension subjects, and blood pressure types. The first method was based on the Bayesian approach and the second one employed Fisherian inference; the purpose of testing both techniques is to provide a solid conclusion about the relationship of the mentioned clinical and sociodemographic variables with hypertension early developments as opposed to the literature in which Frequentist inference or regression techniques are commonly employed to analyze the relationship between sociodemographic and clinical variables with high blood pressure. Furthermore, few studies have made a comparison of the two proposed inference methods that consider the distribution of data without making normality assumptions.

This work has been structured as follows. Section 2 presents a literature review of how hypertension or high blood pressure has been studied from a statistical perspective by analyzing clinical and sociodemographic variables. Section 3 defines the problem and motivation of the present study and the research questions and statistical hypothesis proposed to tackle the problem and aims of the current research. Next, Section 4 presents the proposed statistical inference analysis case study. The Frequentist approach using the Fisherian and Bayesian approach was employed to analyze which variables could be associated with hypertension in its early stages and decide the statistical hypotheses. Finally, based on the results shown in Section 5 and its respective analysis, the conclusions about the study are presented, and future work is also discussed in Section 6.

2. Literature Review

The analysis of diseases allows us to relate their variables and create predictive models that help decision-making before complications occur and, in the best of cases, completely prevent the development of the disease. In high blood pressure, it is necessary to understand that it is made up of clinical and sociodemographic variables. The analyses that are usually carried out are both Frequentist and Bayesian. The articles from 2016 to 2021 show the relationships that statistical analyses looked for. The Scopus database search match is shown in Figure 1.

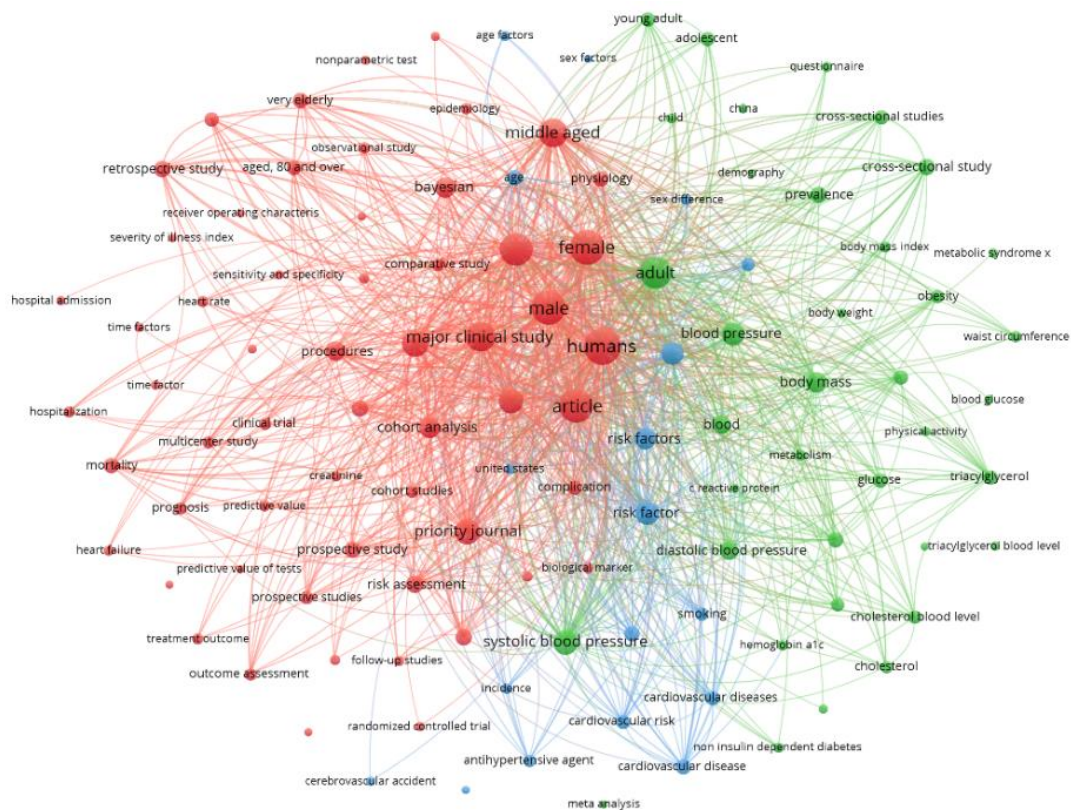


Figure 1. Scopus keyword coincidence about high blood pressure and statistical analysis. Vosviewer, own elaboration.

When analyzing the coincidences between the type of approach practiced on these data, the users' preference to use a Frequentist approach over Bayesian Inference (see Figure 2). Moreover, the relationships between sociodemographic and clinical variables such as gender, BMI, heart rate, or weight with hypertension by employing statistical inference have been studied. For gender, the

study of Everett et al.[10] demonstrate that there are excellent disparities between hypertension levels between men and women in their twenties. In this work, it was shown that women are less likely to have developed hypertension than men.

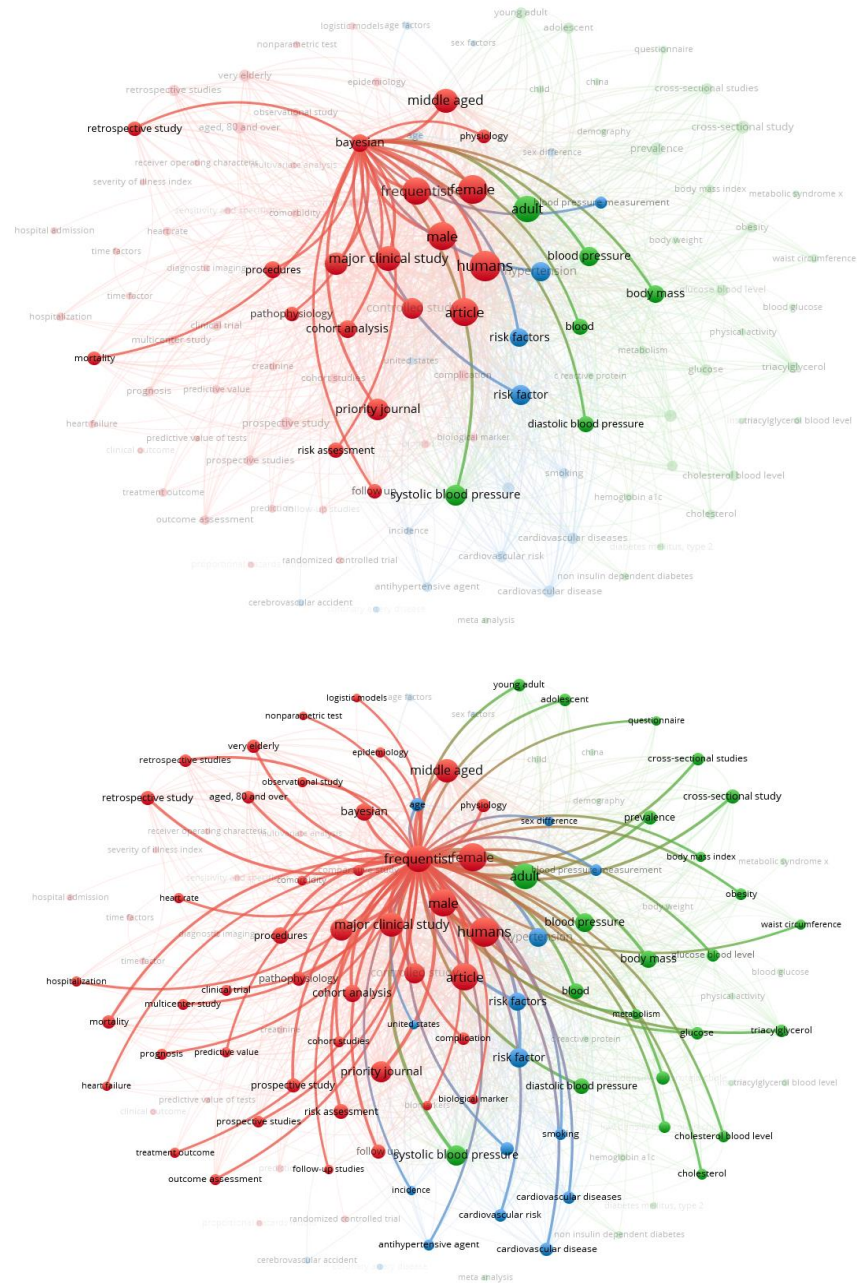


Figure 2. Top, Bayesian approach coincidence. Bottom, Frequentist approach coincidence.

These conclusions were based on the analysis of a multivariate logistic regression analysis in which the level of awareness of hypertension was measured in each group. The results showed that

women are more significantly aware of their hypertension condition than men, with an odds ratio of 0.37 ($p < .05$). On the other hand, Choi et al. [11] demonstrated a higher prevalence of hypertension in men than in women for Korean Adults. These results were based on the use of logistic regression models. According to the authors, the hypertension prevalence was higher for men (34.6%) than for women (30.8%).

In the case of BMI, some studies have also pointed out an influence between BMI and hypertension. According to Landi et al. [12], average systolic and diastolic blood pressure increases significantly and linearly through different BMI levels. To reach this conclusion, the authors employed a logistic regression analysis between different levels of BMI and hypertension. The results of this study showed that the prevalence of hypertension in subjects with average weight (BMI: 18.5–24.9) was 45%, 67% for overweight (BMI: 25.0–29.9) subjects, of 79% for obesity classes I and II (BMI: 30.0–39.9) subjects, and up to 87% for those who were obese class III (BMI greater than 40.0). Besides, Fariha et al. [13] explored the association between BMI and hypertension in South Asia according to their socio-economic group. For this case, a logistic regression model was proposed to estimate the odds ratio with a 95% confidence interval (CI) of high blood pressure for each five-unit increase in BMI. The results showed that for each increment of five-unit increase in the BMI, the odds ratio for hypertension increase 1.79 (95% CI: 1.65–1.93) in Bangladesh, 1.59 (95% CI: 1.58–1.61) in India, and 2.03 (95% CI: 1.90–2.16) in Nepal [13].

Besides, Yang et al. [14] studied the relationship between blood pressure with weight change with ages above 20 years. According to this report, systolic and diastolic blood pressure increased linearly with weight gain based on the results of a linear regression model, which led to obtaining an adjusted mean difference between the highest and the lowest quintiles of weight gain of 6.0 mmHg (95% CI: 5.6–6.5) for systolic blood pressure and 3.9 mmHg (95% CI: 3.6–4.2) for diastolic

blood pressure. On the other hand, Colangelo et al. [15] studied the association of resting heart rate (RHR) in young adults over 30 years with an incident of hypertension in a cohort of Black and White men and women. The hazard ratio for a ten bpm higher RHR was determined using a joint longitudinal time-to-event model composed of a mixed random effect sub-model, quadratic follow-up time, and a survival sub-model modified for confounders. This analysis led to hazard ratios of 1.47 (95% CI: 1.23–1.75), 1.51 (95% CI: 1.28–1.78), 1.48 (95% CI: 1.26–1.73), and 1.02 (95% CI: 0.89–1.17) for Black men, White men, White women, and Black women, respectively. Based on these results, the authors concluded that a higher heart rate is associated with a higher risk of hypertension incidents.

From this literature review, it is possible to notice that most of the studies generally employed a Frequentist approach to determine the associations between sociodemographic and clinical variables with high blood pressure levels. Nevertheless, an inference method that considers the data distribution through the goodness of fit was not shown. Besides, some cited works provide a decision through logistic regression or linear regression models, and other techniques were not tested. Moreover, an analysis of the different stages of hypertension has not been conducted primarily when comparing groups in the early stages of hypertension, such as normal and prehypertension groups. Taking this background as a reference, the following section presents a set of scientific and statistical hypotheses that will be tested by applying Fisherian and Bayesian inference techniques.

3. Problem Definition and Formulation

3.1 Motivation

In general, hypertension does not generate noticeable symptoms that could alert from elevated blood pressure. Due to the above, this health condition is known as the “silent killer” [4]. Hypertension could produce dizziness, chest pain, headaches, breathing difficulties, nosebleeds, or heart palpitations. However, these symptoms are not specific to this condition and usually occur until hypertension has reached a higher stage. On the other hand, medical evidence suggests that early detection of hypertension could reduce its development and consequences [16]. Therefore, high blood pressure detection and monitoring is still an open research topic, especially in lower-income countries with less pervasive and proactive healthcare services. Taking advantage of the increase in the available electronic health records collected in health centers, it is possible to use the data to apply statistical inference techniques that could demonstrate risk factors related to high blood pressure. Moreover, the literature presented in the previous section showed that the Frequentist inference tends to be more used than Bayesian inference, this shows a gap that can be fulfilled by applying and comparing techniques based on both approaches that can produce more reliable results by considering the proper distribution of the data without making normal distribution assumptions that are often considered in Frequentist inference. According to the reasons given, this study aims to determine what is the association between sociodemographic and clinical variables with high blood pressure and its stages, especially in its early developments, by comparing two statistical inference techniques without normal distribution assumptions.

3.2 Dataset

The selected dataset consists of clinical trial data for the noninvasive diagnosis of cardiovascular diseases. The dataset contains information about 219 patients with ages ranging

between 21 to 86 years old. The data collection was performed at the Guilin People's Hospital in China. The data collection was approved by the ethics committee of the Guilin People Hospital and the Guilin University of Electronic Technology in China [22]. The data was subduced to two screening stages. The first one was a screening of missing and abnormal values. If one or more items were missing or an abnormal value was present, the participant record was removed. The second screening was related to removing records that were not associated with cardiovascular diseases since the dataset was designed to focus on patients' clinical information with hypertension or hypertension risk. Although the database has different sociodemographic and clinical variables, for this study, the chosen variables are related to gender, weight, systolic and diastolic blood pressure, heart rate, and BMI of the subjects.

3.3 Research questions and statistical hypotheses

Considering the aim of this study, the motivation, and the variables contained in the dataset, the following research questions were established.

1. The sociodemographic information of the patient is related to a blood pressure type.

To answer this research question, gender was considered the sociodemographic variable. It is used to study its association with the systolic blood pressure of normal and prehypertension subjects. The initial assumption is that there is no difference between men's and women's mean systolic blood pressure. The analysis of systolic blood pressure is crucial since it describes the heart's force to pump the blood through a subject's body. Usually, medical experts based their diagnosis on this blood pressure type. Hence the importance of evaluating the difference in systolic blood pressure of these two groups.

2. The clinical information of the patient is related to a hypertension stage.

The clinical variables considered for this research question are the subjects' BMI, heart rate, and weight. In this case, the first two stages of hypertension (normal and prehypertension) are studied to find if the mentioned clinical variables could be related to the presence of high blood pressure in its early stages. The initial assumption is that there is no difference in the mean values of BMI, heart rate, and weight between normal and prehypertensive subjects. Since blood pressure increases are not so much different from a prehypertensive subject to a normal subject, it is crucial to investigate which variables can describe this difference better.

Table 2 shows the statistical hypotheses that arose from the research questions with their respective mathematical representation.

Table 2 Statistical hypotheses.

SH	Assumption	Statistical Hypothesis	
		Null Hypothesis	Alternative Hypothesis
1	There is no difference between systolic blood pressure means for men and women .	$H_o: \mu_{sf} = \mu_{sm}$	$H_a: \mu_{sf} \neq \mu_{sm}$
		μ_{sf} : Mean Systolic Blood Pressure of females μ_{sm} : Mean Systolic Blood Pressure of males	
2	There is no difference between the body mass index means of normal and prehypertensive subjects.	$H_o: \mu_{bmiN} = \mu_{bmiP}$	$H_a: \mu_{bmiN} \neq \mu_{bmiP}$
		μ_{bmiN} : Mean body mass index of normal subjects μ_{bmiP} : Mean body mass index of prehypertensive subjects	
	There is no difference between the heart rate means of normal and prehypertensive subjects.	$H_o: \mu_{HN} = \mu_{HP}$	$H_a: \mu_{HN} \neq \mu_{HP}$
		μ_{HN} : Mean heart rate of normal subjects μ_{HP} : Mean heart rate of prehypertensive subjects	
	There is no difference between the weight means of normal and prehypertensive subjects.	$H_o: \mu_{Wn} = \mu_{Wp}$	$H_a: \mu_{Wn} \neq \mu_{Wp}$
		μ_{Wn} : Mean weight of normal subjects μ_{Wp} : Mean weight of prehypertensive subjects	

4. Analytical Solution

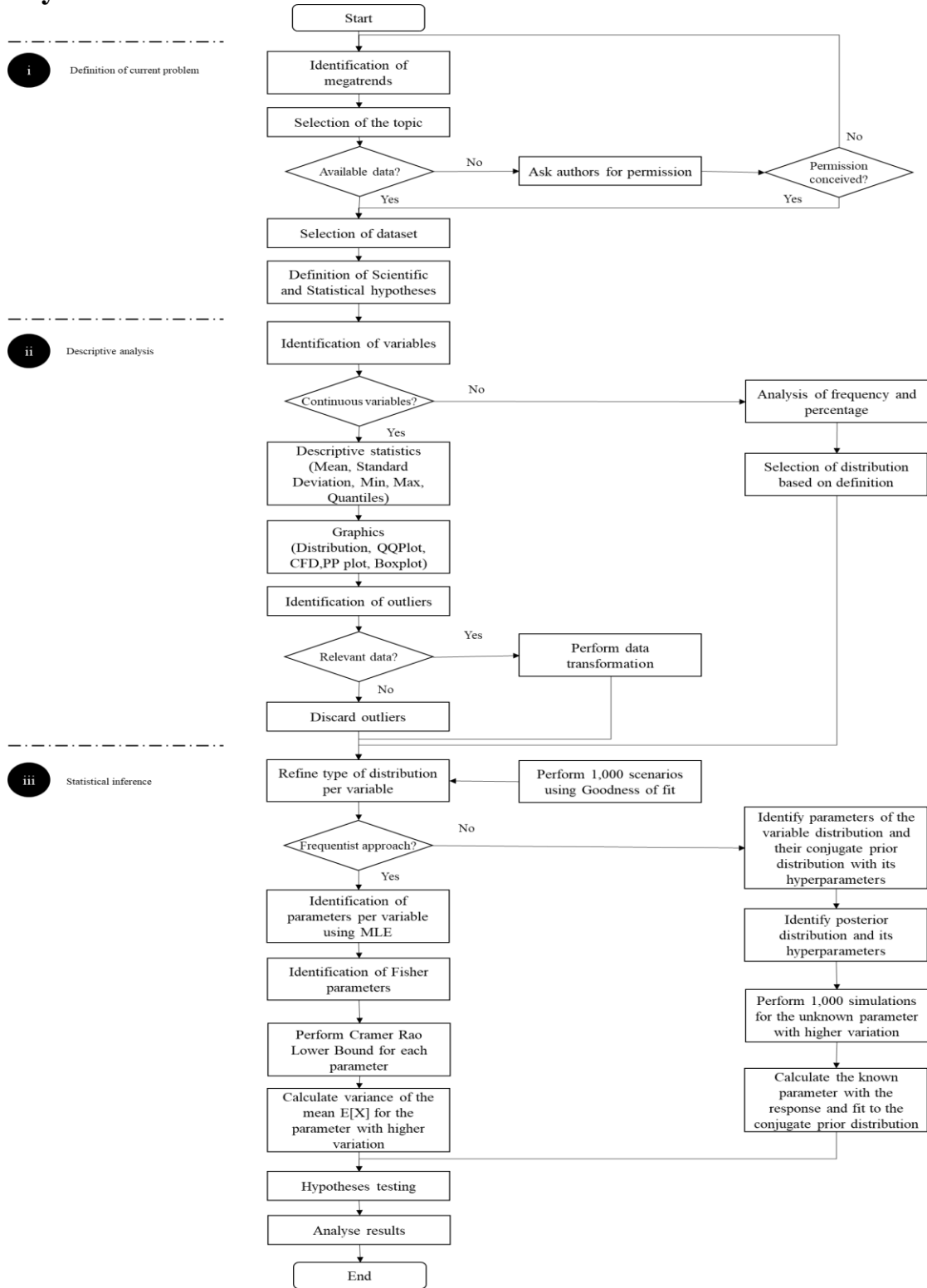


Figure 3. Data analysis flowchart.

During the development of this work, the steps suggested in Figure 3 were deployed. Three general stages were developed: i) Definition of the problem, ii) Descriptive analysis, and iii) Statistical inference techniques. The first stage has been delved previously. The second and third stages are delved into 4.1 and 4.2 sub-sections, respectively. A general guideline to data analysis and the stages are basic analysis activities. In this manner, a structure is provided to perform a statistical inference analysis.

4.1 Descriptive statistics and analysis of the distribution of the variables

The comprehensive clinical data created by Liang et al. [9] was selected to test the proposed statistical hypotheses. The description of each of the variables can be seen in Table 3.

Table 3 Variables description

Type	Variable	Description	Unit
Discrete	Gender	Female or Male participant	F/M
	Hypertension Stage	Current blood pressure classification	N/P/SI/SII
Continuous	Age	Number of years of participant	years
	Height	Height measurement in centimeters	cm
	Weight	Weight measurement in kilograms	Kg
	Blood Pressure (BP)	Blood pressure, systolic and diastolic, measured with an Omron HEM-7201 upper arm blood pressure monitor	mmHg
	Heart rate	Heart rate measure in beats per minute	bpm
	BMI	Body Mass Index	Kg/m ²

The distributions of the discrete variables are related to the frequency of occurrence of the event. In contrast, the continuous ones describe the behavior of the data, and for this, it is necessary to differentiate them. In addition, a brief description of the variables and units has been included in each of the cases. For the categorical variables present in the dataset, the gender proportion can

be appreciated in Figure 4a, while the ratio of blood pressure classification stages can be seen in Figure 4b.

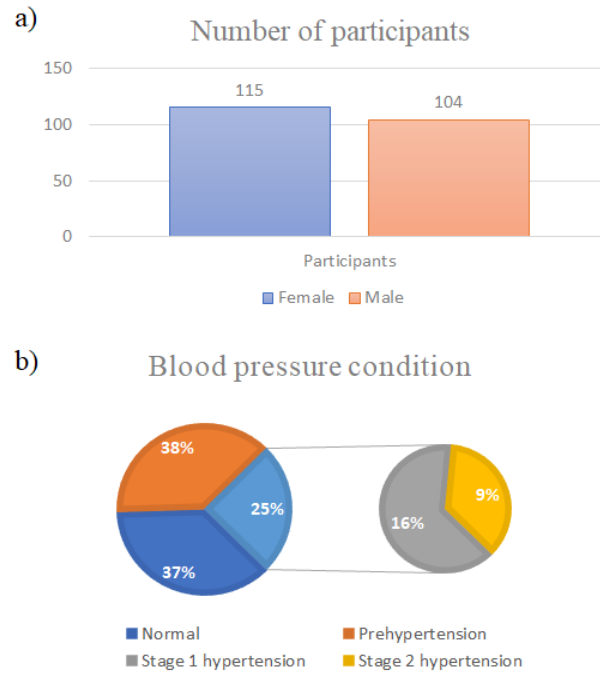


Figure 4. a) Number of participants according to their gender. b) Hypertension stages percentage.

On the other hand, Table 4 presents the summary of the descriptive statistics of the continuous variables of the dataset in which information related to the mean value of each variable, their standard deviation, the minimum and maximum values, and quantiles are also presented.

Table 4 Descriptive statistics of the continuous variables

Variable	Mean	Standard deviation	Min	25%	50%	75%	Max
Age (years)	57.1789	15.8743	21	48	58	67.5	86
Height (meters)	161.2283	8.2028	145	155	160	167	196
Weight (kilograms)	60.1917	11.8863	36	52.5	60	66.5	103
Systolic Blood Pressure (mmHg)	127.9452	20.3777	80	113.5	126	139	182
Diastolic Blood Pressure (mmHg)	71.8493	11.1112	42	64	79	78	107
Heart rate (beats/minute)	73.6392	10.7388	52	66	73	80	106
BMI kg/m^2	23.1072	4.0043	14.69	20.55	22.6	25	37.46

Besides, Figure 5 presents the boxplot of each of the variables involved in the statistical hypotheses where outliers in some of the variables can be appreciated.

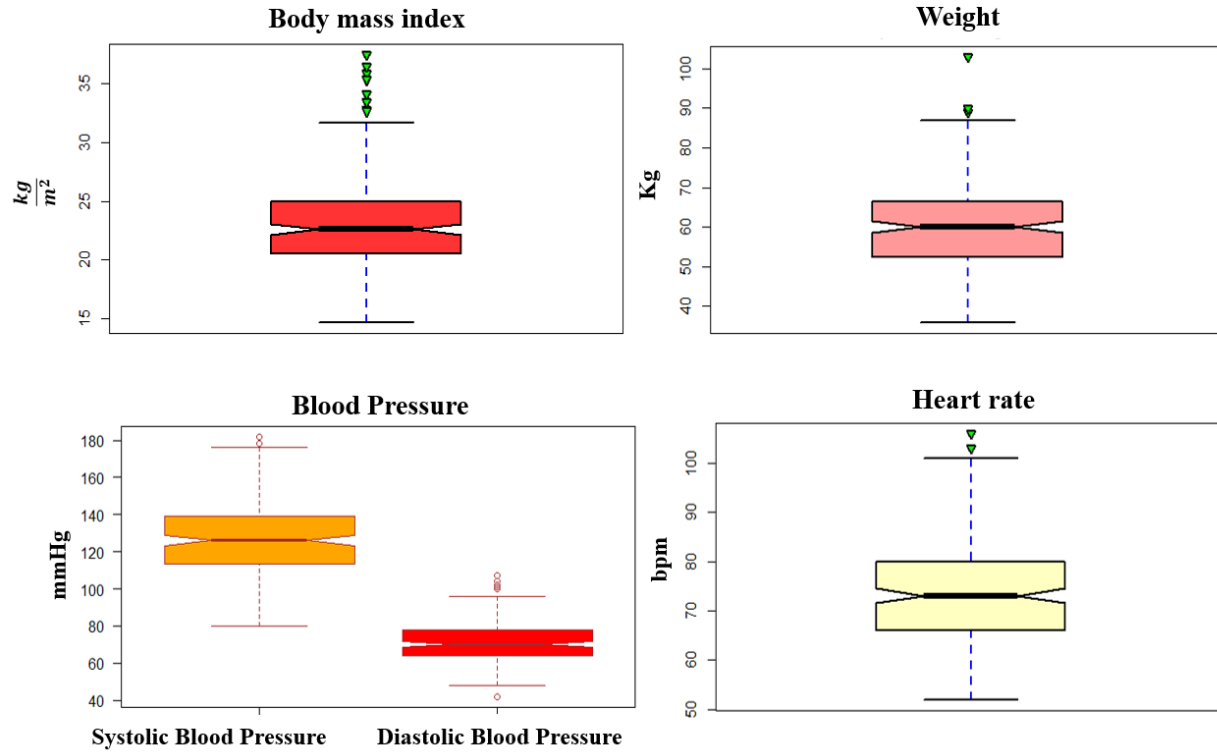


Figure 5. Boxplot of continuous variables.

Although each of the variables has several outliers, it will not be significant in any of the cases, and this can be seen in Table 5, where it is possible to notice that the number of outliers for each of the variables is below 3% of the data except for the BMI whose number of outliers reach 4% of the whole sample. Despite having more than 3% of outliers, the BMI variable was not transformed, and the outliers were not removed. This decision was taken since these values represent that the relationship between height and weight is higher in some patients. Therefore, they are above the average values of the weight in Guanxi.

Table 5 Number and percentage of outliers in each of the variables

Variable	Outliers	% of Outliers
Age	0	0
Height	2	0.91
Weight	5	2.28
Systolic Blood Pressure	3	1.37
Diastolic Blood Pressure	6	2.74
Heart rate	3	1.37
BMI	9	4.11

To select the best distribution of the variables involved in the statistical hypotheses, the typical type of univariate distributions used to model the variables is presented in Table 6 based on a brief literature review.

Table 6. Distributions are used in the literature to model the variables of the dataset.

Variable	Author	Distribution
Weight (kg)	Nouri et al. [17]	Normal distribution
	Burmaster et al. [18]	Log-normal distribution
Blood Pressure (mmHg)	Pater et al. [19]	Log-normal distribution
	Makutch et al. [20]	Normal Distribution
Heart Rate (b/m)	Ostchega et al. [21]	Normal Distribution
	Poole et al. [22]	Log-Normal Distribution
Body Mass Index $\frac{kg}{m^2}$	Penman et al. [23]	Normal Distribution/Log-normal distribution

As can be seen, the literature has commonly used the normal distribution to model variables such as weight, BMI, blood pressure, and heart rate. Nevertheless, in some cases, the authors have pointed out that if the data presents some skewness, the log-normal distribution could be a better fit. Furthermore, the goodness of fit was employed to determine the best distribution of the variables. Different univariate distributions were tested, including the Normal, Weibull, Gamma, and lognormal. Using the riskdistribution package of R, and selecting the best fit was based on the loglikelihood values; the closer its value to zero, the better the fit to the distribution.

For all the variables, the first iteration of the goodness of fit shown that the best distribution was the log-normal distribution based on the loglikelihood value closest to zero. Nevertheless, a simulation was performed for the goodness of fit for each of the variables involved in the statistical hypotheses. The simulation process was applied 1,000 times using bootstrapping. The results are shown in Table 7.

Table 7. The goodness of fit simulation using bootstrapping.

Variable	Frequency of the loglikelihood value closest to zero						Total
	Chi-Square	Normal	Gamma	Weibull	Logistic	Log-normal	
Weight	1	0	42	0	2	955	1,000
Systolic Blood Pressure	0	17	354	0	7	622	1,000
Diastolic Blood Pressure	4	2	133	0	32	829	1,000
Heart Rate	0	5	91	0	0	904	1,000
Body Mass Index	0	1	9	0	12	978	1,000

By making this procedure, the best distribution for the variables involved in the statistical hypotheses was the log-normal distribution. Since the statistical hypotheses referred to comparing the mean of two groups, it is essential to consider how the mean is calculated for the log-normal distribution. For a log-normal random variable (x) with parameters μ and σ , the probability density function is shown in Equation 1 [24].

$$f(x) = \frac{1}{x\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{\ln x - \mu}{\sigma}\right)^2} \quad (1)$$

In equation 1, μ is the mean log, and σ is the standard log of the log-normal distribution. The population means $E[X]$, the variance $V[X]$, and the standard deviation $SD[X]$ of the log-normal distribution are calculated as follows.

$$E[X] = e^{\mu + \frac{\sigma^2}{2}} \quad (2)$$

$$V[X] = e^{2\mu + \sigma^2} (e^{\sigma^2} - 1) \quad (3)$$

$$SD[X] = e^{\mu + \frac{\sigma^2}{2}} (\sqrt{e^{\sigma^2} - 1}) \quad (4)$$

Taking the previous analysis as a reference, it is possible to notice that the data does not follow a normal distribution as Frequentist inference methods assume to test statistical hypotheses. Therefore, due to the complexity that the data has, it is essential to apply techniques that can handle it and, in this way, provide a more reliable result and decision for the testing of the proposed

statistical hypotheses. To do the above, Bayesian and Fisherian inference techniques are proposed to handle this complexity. In the next section, the procedures used to apply them, and the results obtained for each of them are presented.

4.2 Statistical inference techniques

4.2.1 Bayesian Approach

In the Frequentist approach, the parameter θ of a population is assumed to be a fixed quantity estimated through the observed data. On the other hand, in Bayesian Inference, the unknown parameter θ , is treated as a random variable. In other words, it is assumed that there exists an initial guess about the distribution of the parameter. This distribution is known as the prior distribution. After observing new data, the distribution of θ is updated based on this new data. This updating process is performed through Bayes' Rule [25]. This updating produces a further distribution for θ known as posterior distribution. This relation can be express as shown in Equation 5.

$$P(\theta|D) = \frac{P(D|\theta)P(\theta)}{P(D)} \quad (5)$$

Where D refers to the collected data, $P(D|\theta)$ is the likelihood function of θ , $P(\theta)$ is the prior distribution of θ , $P(\theta|D)$ is the posterior distribution of θ , and $P(D)$ is known as the evidence. The evidence term refers to how probable it is to observe the collected data.

The computation of this conditional probability is complex. Therefore, to apply the Bayesian approach, it is necessary to determine the conjugate prior to the data distribution and calculate the prior and posterior hyperparameters. This conjugate prior can be consulted if the data distribution belongs to the exponential family as it is the case of the log-normal distribution. Table 8 shows the conjugate prior of the log-normal distribution for each of its parameters and how the prior and posterior hyperparameters should be computed assuming one of the parameters to be known.

Table 8. Conjugate prior and posterior hyperparameters of the log-normal distribution.

Likelihood	Model Parameter s	Conjugate Prior	Prior hyperpara- meters [26]	Posterior hyperparameters [26]
Log-normal (With known σ^2)	μ	Normal	μ_0, σ_0^2	$\mu \sim \text{normal}(\frac{\frac{\mu_0}{\sigma_0^2} + \frac{\sum_{i=1}^n \ln(x_i)}{\sigma^2}}{\frac{1}{\sigma_0^2} + \frac{n}{\sigma^2}}, (\frac{1}{\sigma_0^2} + \frac{n}{\sigma^2})^{-1})$
Log-normal (With known μ)	σ^2	Inverse Gamma	α, β	$\sigma^2 \sim \text{inverse gamma}(\alpha + \frac{n}{2}, \beta + \frac{\sum_{i=1}^n (\ln(x_i) - \mu)^2}{2})$
μ : mean log – lognormal distribution σ^2 : variance log – lognormal distribution μ_0 : Mean normal distribution σ_0^2 : Variance normal distribution				α : Shape parameter inverse gamma distribution β : Scale parameter inverse gamma distribution n : the size of the sample

In this case, each of the hypotheses' prior and posterior distribution curves was calculated based on the conjugate prior for the μ parameter. This selection was made since the μ parameter had a higher variation compared to σ according to the inverse of the Fisher Information of each parameter. The known parameter value was calculated by using the maximum likelihood estimation of σ before splitting the variables into the respective groups. Moreover, to estimate the prior hyperparameters, 1000 simulations using bootstrapping were run to get the value of the unknown parameter and then used the resultant data to fit it to the conjugate prior distribution.

4.2.1.1 Results in the Bayesian approach

The results obtained by applying the Bayesian Approach can be appreciated in Tables 9 to 12. The values of the known parameter (σ) and the values of the prior and posterior hyperparameters for μ are presented from which the prior and posterior distribution curves were drawn. Moreover, Figures 6 and 7 show the results and comparison of each group's prior and posterior distribution curves with their respective statistical hypothesis. Since the posterior distribution of each of the groups involved in the statistical hypotheses overlap, it is possible to state that the null hypotheses are not rejected.

Table 9. Prior and posterior hyperparameters for the first statistical hypothesis.

Hypothesis	Model Parameters	Group	Prior hyperparameters	Posterior hyperparameters
There is no difference between systolic blood pressure means for men and women.	μ (<i>Unknown</i>) $\sigma^2 = 0.025131$	Female	$\mu_{0,f} = 4.82331675$ $\sigma_{0,f}^2 = 0.00023402$	$\mu_f \sim \text{normal}(4.823429, 1.1207^{-4})$
		Male	$\mu_{0,m} = 4.856181$ $\sigma_{0,m}^2 = 0.0002133$	$\mu_m \sim \text{normal}(4.85617, 1.1329^{-4})$

f: Parameters of female group *m*: Parameters of male group

Table 10. Prior and posterior hyperparameters for the second statistical hypothesis.

Hypothesis	Model Parameters	Group	Prior hyperparameters	Posterior hyperparameters
There is no difference between the body mass index means of normal and prehypertensive subjects.	μ (<i>Unjkown</i>) $\sigma^2 = 0.0277$	Normal	$\mu_{0,n} = 3.0841953$ $\sigma_{0,n}^2 = 0.0003432$	$\mu_n \sim \text{normal}(3.08412, 1.725534^{-4})$
		Pre-hypertension	$\mu_{0,p} = 3.124587$ $\sigma_{0,p}^2 = 0.0002504$	$\mu_p \sim \text{normal}(3.12479, 1.417267^{-4})$

n: Parameters of Normal group *p*: Parameters of Prehypertensive group

Table 11. Prior and posterior hyperparameters for the third statistical hypothesis.

Hypothesis	Model Parameters	Group	Prior hyperparameters	Posterior hyperparameters
There is no difference between the heart rate means of normal and prehypertensive subjects.	μ (<i>Unknown</i>) $\sigma^2 = 0.02073$	Normal	$\mu_{0,n} = 4.2904850$ $\sigma_{0,n}^2 = 0.0002466$	$\mu_n \sim \text{normal}(4.290349, 1.263757^{-4})$
		Pre-hypertension	$\mu_{0,p} = 4.2788912$ $\sigma_{0,p}^2 = 0.0002324$	$\mu_p \sim \text{normal}(4.279013, 1.190213^{-4})$

n: Parameters of Normal group *p*: Parameters of Prehypertensive group

Table 12. Prior and posterior hyperparameters for the Fourth statistical hypothesis.

Hypothesis	Model Parameters	Group	Prior hyperparameters	Posterior hyperparameters
There is no difference between the weight means of normal and prehypertensive subjects.	μ (<i>unknown</i>) $\sigma^2 = 0.037167$	Normal	$\mu_{0,n} = 4.04155$ $\sigma_{0,n}^2 = 0.0004021$	$\mu_n \sim \text{normal}(4.041194, 2.152472^{-4})$
		Pre-hypertension	$\mu_{0,p} = 4.07582$ $\sigma_{0,p}^2 = 0.0004387$	$\mu_p \sim \text{normal}(4.075827, 2.190153^{-4})$

n : Parameters of Normal group p : Parameters of Prehypertensive group

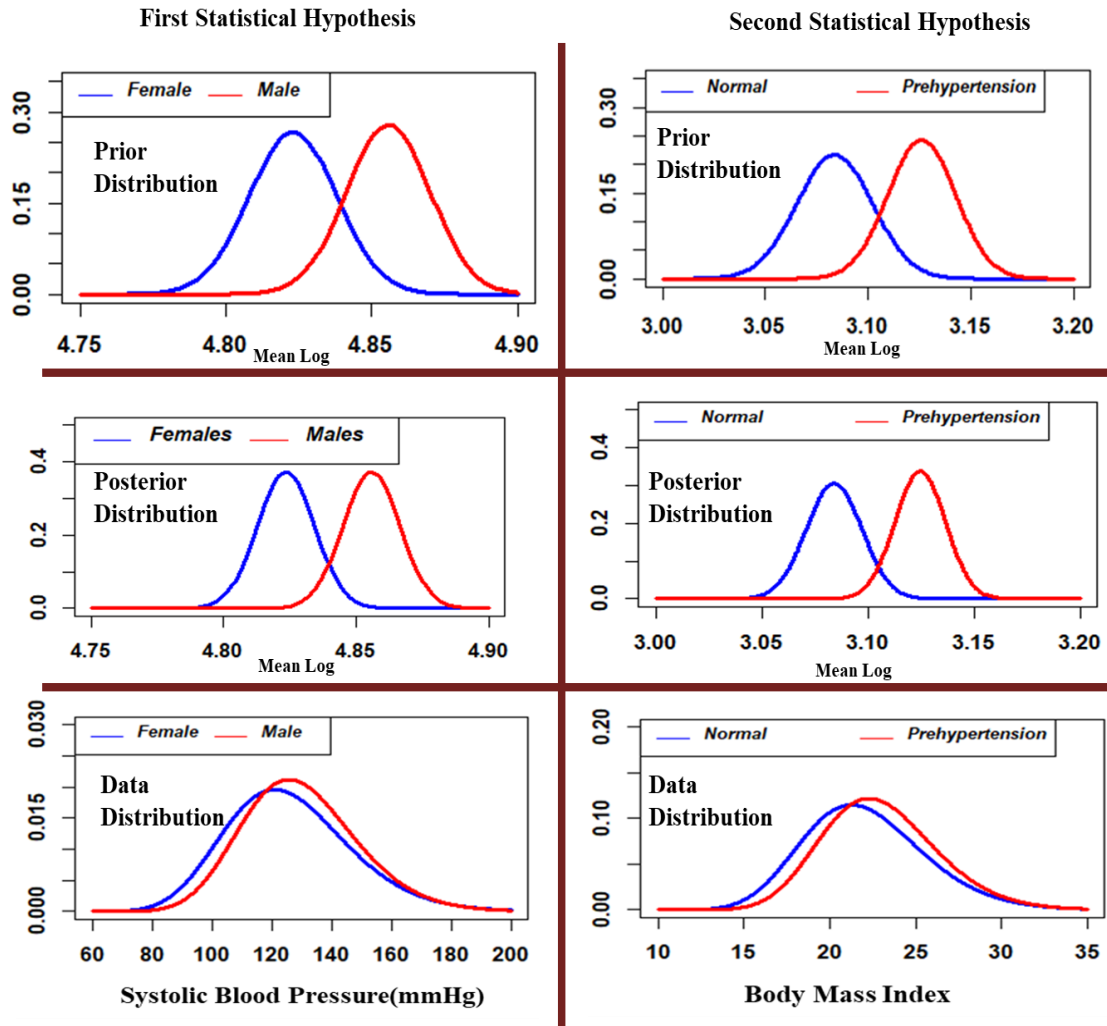


Figure 6. Prior, Posterior, and Data Distribution of the First (left) and Second (right) Statistical Hypotheses.

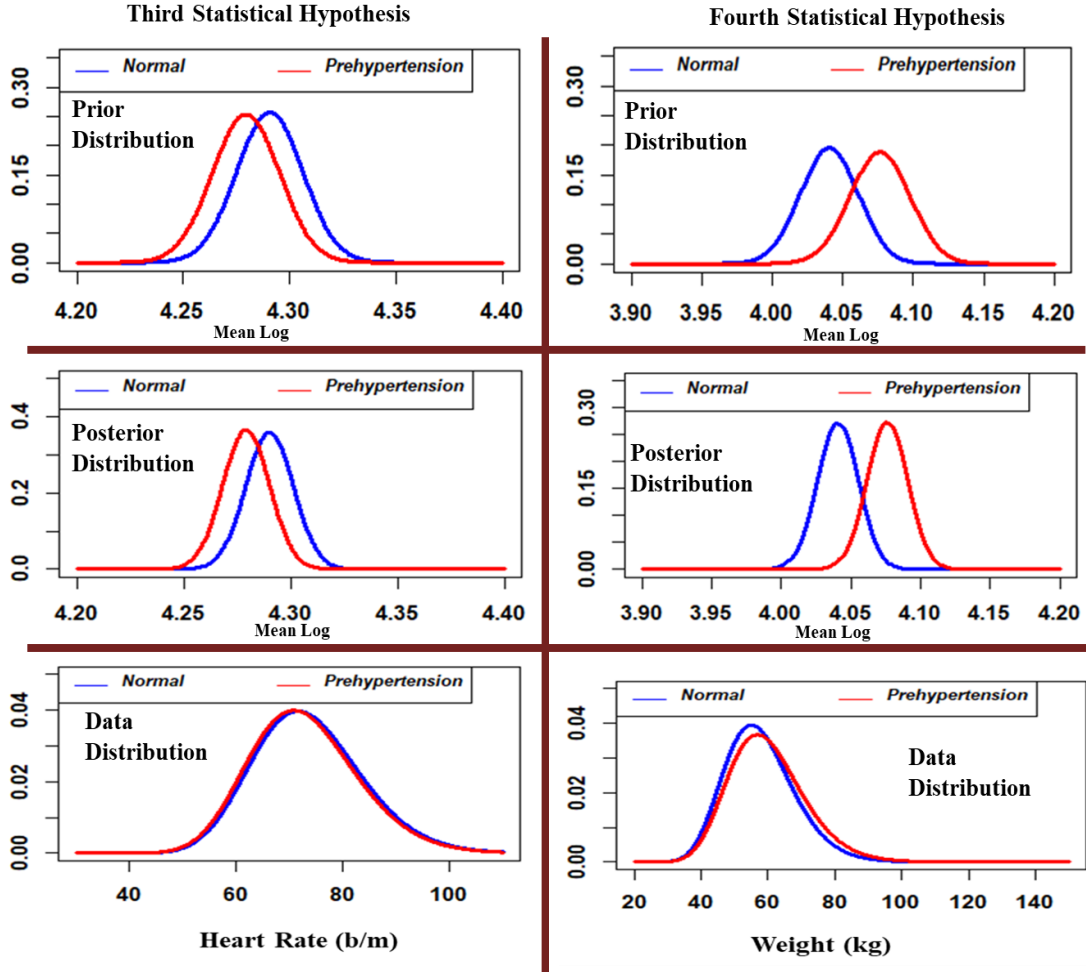


Figure 7. Prior, Posterior, and Data Distribution of the Third (left) and Fourth (right) Statistical Hypotheses.

4.2.2 Fisherian Approach

To apply the Fisherian approach, three concepts need to be considered. First, the maximum likelihood estimation (MLE), a technique used for estimating the parameters of a probability distribution by maximizing its likelihood function so that the observed data is most likely under the assumed statistical model. Fisher Information, a method used to calculate the amount of information that an observable random variable has about an unknown parameter of a distribution that models the random variable, and the Cramer's Roa Lower Bound (CRLB), which expresses a lower bound for the variance of an unbiased estimator, based on Fisher's information. The MLE

for a log-normal distribution was computed for each variable and group involved in the statistical hypotheses. Later the Fisher information was used to estimate the variance and CRLB of the distribution parameters, and later the normal curves of the mean of the groups involved in the hypothesis were plotted. The expressions used to apply the Fisherian approach and test the proposed statistical hypothesis can be appreciated in Table 13 [27].

The mean for a log-normal distribution is given by equation 2; therefore, it is necessary to calculate the mean-variance according to the Delta Method. Since the mean value depends on two parameters, μ , and σ , one is assumed constant depending on its variance. The parameter with the lower variation was considered fixed. Moreover, the CRLB of the mean for each group was compared against the mean square error by simulation using bootstring 1000 times. This procedure was done to ensure that the correct distribution of the data was selected.

Table 13. Likelihood function, Fisher Information, and CRLB of the parameters of a log-normal distribution.

Likelihood Function of the log-normal distribution	$\prod_{i=1}^n f(x_i) = \prod_{i=1}^n \frac{1}{x_i \sigma \sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{\ln x_i - \mu}{\sigma} \right)^2}$
The loglikelihood function of the log-normal distribution	$L = -\frac{n}{2} \ln(\sigma^2) - \frac{n}{2} \ln(2\pi) - \frac{1}{2\sigma^2} \sum_{i=1}^n (\ln x_i - \mu)^2 - \sum_{i=1}^n \ln(x_i)$
MLE of the μ parameter.	$\hat{\mu} = \frac{\sum_{i=1}^n \ln x_i}{n}$
MLE of σ^2 parameter	$\hat{\sigma}^2 = \frac{\sum_{i=1}^n (\ln x_i - \mu)^2}{n}$
Fisher Information of μ	$I(\hat{\mu}) = \frac{n}{\sigma^2}$
Fisher Information of σ	$I(\hat{\sigma}) = \frac{2n}{\sigma^2}$
Variance and CRLB of μ	$\text{var}(\hat{\mu}) = \frac{\sigma^2}{n}$
Variance and CRLB of σ	$\text{var}(\hat{\sigma}) = \frac{\sigma^2}{2n}$
Variance and CRLB of the mean $E[X]$ if μ have a higher variation.	$\text{VAR}[E[X]] = \text{CRLB}_{E[X]} = \left(e^{\mu + \frac{\sigma^2}{2}} \right)^2 \frac{\sigma^2}{n}$
Variance and CRLB mean $E[X]$ if σ has a higher variation.	$\text{VAR}[E[X]] = \text{CRLB}_{E[X]} = \left(\sigma e^{\mu + \frac{\sigma^2}{2}} \right)^2 \frac{\sigma^2}{2n}$

4.2.2.1 Results in the Frequentist approach

In Tables 14 to 17, the mean of each group, the variance, and the CRLB for the statistical hypotheses can be appreciated. In this case, the parameter with the higher variation was μ based on the inverse of the Fisher information; therefore, the variance of the mean $E[X]$ was calculated as follows.

$$VAR[E[X]] = \left(e^{\mu + \frac{\sigma^2}{2}} \right)^2 \frac{\sigma^2}{n} \quad (6)$$

Table 14. Mean, the variance of the mean, and CRLB for the First statistical hypothesis.

Hypothesis	Group	$E(X)$, $Var(E(X))$ and $CRLB(E(X))$
There is no difference between systolic blood pressure means for men and women.	Female	$E(X)_{sf} = 126.134$
		$Var(E(X)_{sf}) = 3.80798$
		$CRLB(E(X)_{sf}) = 3.80798$
	Male	$E(X)_{sm} = 129.9486$
		$Var(E(X)_{sm}) = 3.560454$
		$CRLB(E(X)_{sm}) = 3.560454$
Sf :Parameters of Females Group Sm : Parameters of Males Group		

Table 15. Mean, the variance of the mean, and CRLB for the Second statistical hypothesis.

Hypothesis	Group	E(X), Var (E(X)) and CRLB(E(X))
<p>There is no difference between the</p> <p>body mass index means of normal and prehypertensive</p> <p>subjects.</p>	Normal	$E(X)_{BMIn} = 22.1321$
		$Var(E(X)_{BMIn}) = 0.1588$
		$CRLB(E(X)_{BMIn}) = 0.15886$
	Prehype rtension	$E(X)_{BMIp} = 23.00325$
		$Var(E(X)_{BMIp}) = 0.13149$
		$CRLB(E(X)_{BMIp}) = 0.13149$
<p>BMIn: Parameters of Normal group</p> <p>BMIp: Parameters of Prehypertensive group</p>		

Table 16. Mean, the variance of the mean, and CRLB for the Third statistical hypothesis.

Hypothesis	Group	E(X), Var (E(X)) and CRLB(E(X))
There is no difference between the heart rate means of normal and prehypertensive subjects.	Normal	$E(X)_{Hn} = 73.69074$
		$Var(E(X)_{Hn}) = 1.313017$
		$CRLB(E(X)_{Hn}) = 1.313017$
	Pre-hyper-tension	$E(X)_{Hp} = 72.88665$
		$Var(E(X)_{Hp}) = 1.220483$
		$CRLB(E(X)_{Hp}) = 1.220483$
Hn: Parameters of Normal group		
Hp: Parameters of Prehypertensive group		

Table 17. Mean, the variance of the mean, and CRLB for the Fourth statistical hypothesis.

Hypothesis	Group	E(X), Var (E(X)) and CRLB(E(X)).
There is no difference between the weight means of normal and prehypertensive subjects.	Normal	$E(X)_{wn} = 57.80737$
		$Var(E(X)_{wn}) = 1.365242$
		$CRLB(E(X)_{wn}) = 1.365242$
	Pre-hyper-tension	$E(X)_{wp} = 59.94701$
		$Var(E(X)_{wp}) = 1.490524$
		$CRLB(E(X)_{wp}) = 1.490524$
Wn: Parameters of Normal group		
Wp: Parameters of Prehypertensive group		

The first column of Tables 14 to 17 collects the information on the null hypothesis to be tested, the second column describes the groups that are being tested (e.g., female, and male patients, patients with normal or prehypertensive stage hypertension), in the third column the data collected for the method. First, the population mean $E[x]$, the variance of the mean, and finally, the CRLB of each group. Figure 8 shows the Normal curves for the mean $E[x]$ for each group for the first and second statistical hypotheses, as well as the comparison of the mean square error (MSE) of each parameter and the CRLB of each group involved in the statistical hypothesis using 1000 simulations with bootstrapping. The same information can be found in Figure 9 for the third and fourth statistical hypotheses.

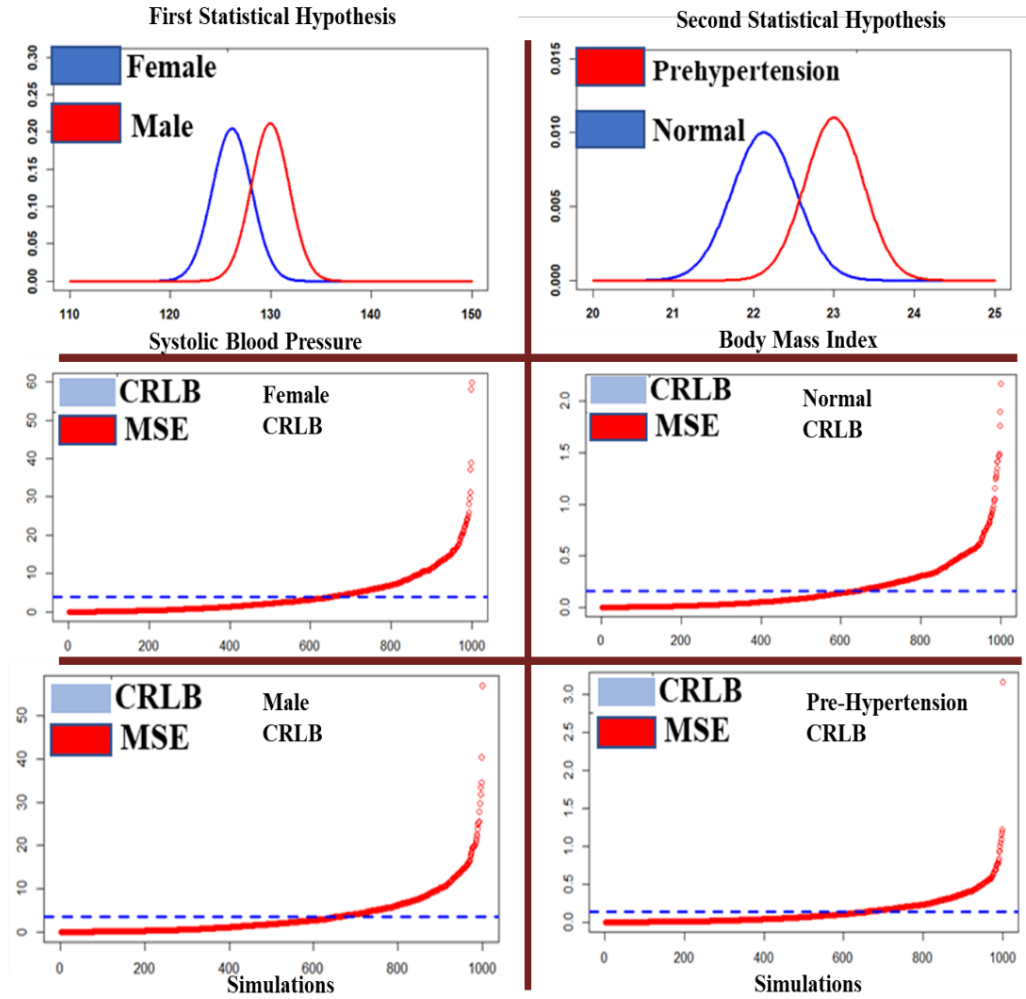


Figure 8. Results of the First (left) and Second (right) statistical hypotheses applying the Fisherian Approach.

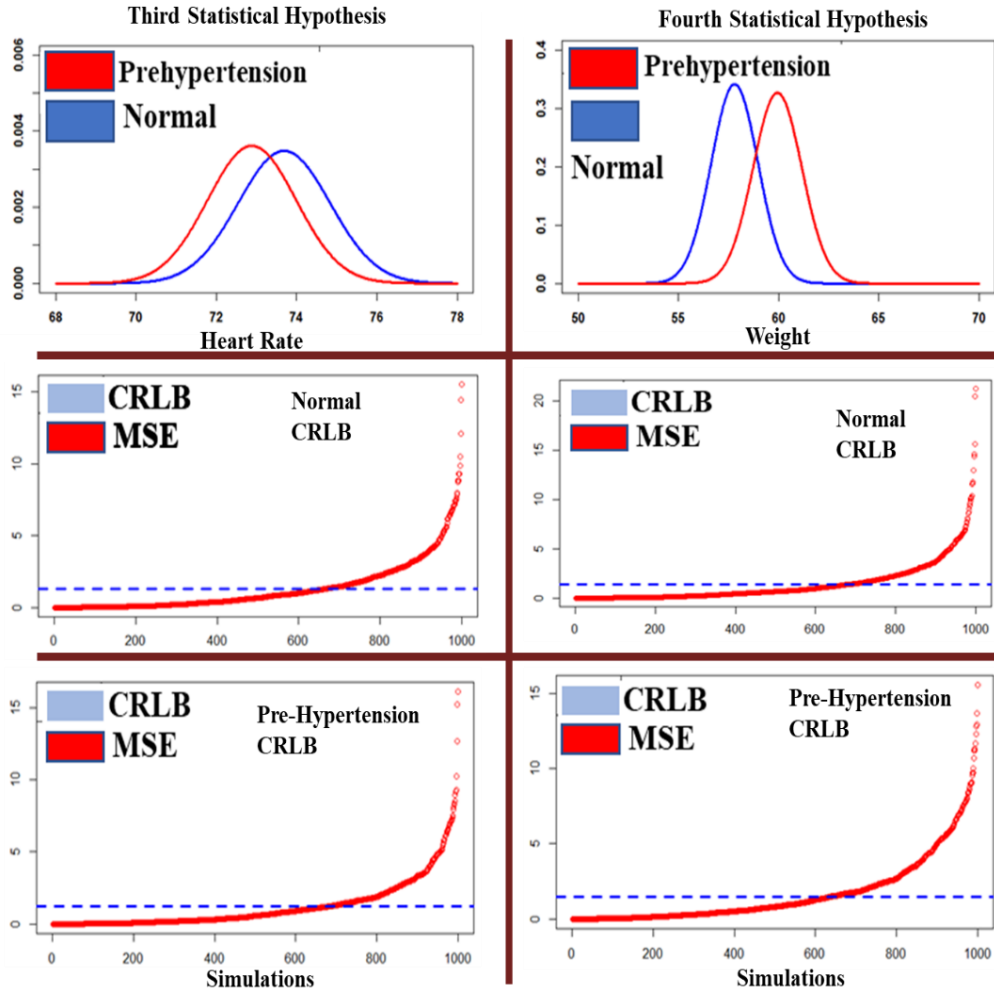


Figure 9. Results of the Third (left) and Fourth (right) statistical hypothesis applying the Fisherian Approach.

In each of the normal curves, both groups share an affinity, which is demonstrated by the overlapping of the blue and red curves in Figures 8 and 9. Due to this fact, it is asserted that the groups belong to the same space and can be compared. Thus, none of the null hypotheses are rejected.

5. Results Analysis

Table 18 presents the results of applying the Fisherian and Bayesian approaches. It is possible to observe that using both methods led to the same results in not rejecting the proposed null hypotheses. Since the Bayesian approached treats the parameters as random variables. At the same

time, the data is considered fixed, and the decision is based on the overlapping of the posterior distributions of each group. The posterior distribution will not vary too much in its location compared to the prior distribution. The above allowed that the posterior distribution of the compared groups overlapped.

On the other hand, in the Fisherian approach, the parameters compared are the mean of both groups, which requires computing the mean as expressed in Equation 2 and the variance as shown in Equation 6 to draw the normal curve of the parameter. Nevertheless, since the calculated mean depends on the data, it is expected that by changing the data, the value and variance of the mean will change and produce a different result or overlapping between the normal curves of the mean, in other words, the location of the normal curve of the mean will change. The above could produce different results that could lead to rejecting the null hypothesis [28]. Although either of the two approaches can be treated independently, in this work, both methods have been applied to handle the complexity of the data, which follows a log-normal distribution, whose results of not rejecting the null hypothesis were the same.

By comparing the results with the literature, it is possible to notice some discrepancies. Nevertheless, this can be attributed to how the relationship between the variables and hypertension was studied. For instance, in the works of Everett et al. [10] and Choi et al. [11], the presence of hypertension between females and males was concluded based on the level of awareness of the subjects, but not in the comparison of systolic blood pressure as in the present study. Besides the level of awareness does not differ too much between groups in the study of Choi et al. [11].

On the other hand, the studies of Landi et al. [12] and Fariha et al. [13] showed an increase in blood pressure related to the rise in BMI. Nonetheless, the presented analysis concluded that there is no difference in BMI between normal and prehypertension subjects. This may be due to the

criterion in which the null hypothesis is rejected because a non-overlap is required between the curves of the compared parameters to reject the null hypotheses. However, in Figure 6, the results of the Bayesian approach in the posterior distribution curves overlap in a smaller region. If the decision criterion were to be relaxed to a smaller interval or percentage of overlapping, it would be possible to interpret the results differently. Besides, the values of that posterior distribution curve of the prehypertension group are higher than the normal group. This suggests a slight increase in BMI between hypertension groups, but this increase is not sufficient to reject the null hypothesis according to the proposed decision criteria.

Similar behavior can be appreciated for the weight (Figure 7) in which a complete overlapping of the posterior distribution curves is not present and the ranges of values that the posterior curve of the weight of prehypertension subjects is slightly higher compared to the curve of normal subjects which is expected as suggested by Yang et al. [14] but again based on the criteria used to reject the null hypothesis the difference is not sufficient to reject the hypothesis. Finally, heart rate showed a major overlapping in the Bayesian and Fisherian approaches. Compared to the literature, the relationship between the heart rate and blood pressure has been less studied; nonetheless, the work of Mancia et al. [29] compared the heart rate variably of normotensive and prehypertensive subjects and showed that it was similar for both groups.

Table 18. Comparison between the Frequentist and Bayesian approaches.

Number	Null Hypothesis	Fisherian approach	Bayesian approach
1	There is no difference between systolic blood pressure means for men and women	Not Rejected	Not Rejected
2	There is no difference between the body mass index means of normal and prehypertensive subjects .	Not Rejected	Not Rejected
3	There is no difference between the heart rate means of normal and prehypertensive subjects .	Not Rejected	Not Rejected
4	There is no difference between the weight means of normal and prehypertensive subjects .	Not Rejected	Not Rejected

6. Conclusions and Future Work

Hypertension is a global public issue whose late diagnosis and inadequate treatment could develop other diseases such as heart attack, kidney diseases, diabetes, or strokes. Besides, this health condition does not produce noticeable symptoms that could alert from its presence in an early stage. Nevertheless, the data availability in EHRs has enabled the identification of possible risk factors related to high blood pressure. Considering the above, this work presented an analysis of the association between sociodemographic and clinical variables with the stages of hypertension in its early development by applying and comparing the Bayesian and Fisherman approaches in the proposed statistical hypotheses presented in Table 2 based on the clinical trial data provided by Guilin's People Hospital of Guanxi, China. The results of applying both methods revealed that gender, body mass index, weight, and heart rate are not different between normal and prehypertensive subjects. The not rejection of the proposed null hypotheses concluded by applying the Bayesian inference was further confirmed by the Fisherian approach. These results can be attributed to the groups of hypertensions analyzed, in this case, normal and prehypertension groups since their blood pressure characteristics do not differ too much. Nevertheless, it is essential to

mention that the results obtained are only applicable to the Guanxi region of China and cannot be generalized to other populations.

Moreover, the results obtained differ from the ones in the literature in the sense that there is no apparent association between BMI, weight, heart rate, or gender with the stages of hypertension and that two inference techniques were employed and compared to provide a decision about the statistical hypothesis in which the proper distribution of the data was considered for the analysis instead of making normal distribution assumptions. Since both approaches are relatively easy to compute and interpret, it could be a good practice to test both methods and compare the results obtained from the two perspectives to generate a reliable decision while analyzing other types of variables and their association with hypertension. Instead of only support the final decision in the Frequentist inference or regression models as it is commonly shown in the literature. Furthermore, the tested strategies have the advantage of considering the intricacy of the data without relying on normal distribution assumptions.

Furthermore, the obtained results demonstrate that the analyzed variables are not recommended to model the behavior of hypertension in its early stages for the case of the people living in Guanxi, China. Even though the variables in the literature were explored, in the present study, it was not possible to reject the hypothesis based only on the selected variables. The above highlights a gap that can be further explored, like studying other types of variables such as age, race, smoking status, oxygen levels, or cholesterol levels, to name a few, and its association with high blood pressure in its early stages. The above also requires the development of a clinical protocol that serves to collect the necessary data. These factors point out a future direction of the present work in which the comparison between normotension and prehypertension can be studied in more detail, considering the suggested variables.

References

- [1] B. C. Taylor, T. J. Wilt, and H. G. Welch, "Impact of diastolic and systolic blood pressure on mortality: implications for the definition of 'normal,'" *J. Gen. Intern. Med.*, vol. 26, no. 7, pp. 685–690, 2011.
- [2] A. V Chobanian *et al.*, "The Seventh Report of the Joint National Committee on Prevention, Detection, Evaluation, and Treatment of High Blood Pressure: the JNC 7 report.," *JAMA*, vol. 289, no. 19, pp. 2560–2572, May 2003, doi: 10.1001/jama.289.19.2560.
- [3] W. H. Organization, *Prevention of cardiovascular disease: guidelines for assessment and management of total cardiovascular risk*. World Health Organization, 2007.
- [4] W. H. Organization, "A global brief on hypertension: silent killer, global public health crisis: World Health Day 2013," World Health Organization, 2013.
- [5] A. Turkyilmaz, D. Dikhanbayeva, Z. Suleiman, S. Shaikholla, and E. Shehab, "Industry 4.0: Challenges and opportunities for Kazakhstan SMEs," *Procedia CIRP*, vol. 96, pp. 213–218, 2021, doi: <https://doi.org/10.1016/j.procir.2021.01.077>.
- [6] G. Aceto, V. Persico, and A. Pescapé, "Industry 4.0 and Health: Internet of Things, Big Data, and Cloud Computing for Healthcare 4.0," *J. Ind. Inf. Integr.*, vol. 18, no. February 2019, p. 100129, 2020, doi: 10.1016/j.jii.2020.100129.
- [7] M. Javaid and A. Haleem, "Industry 4.0 applications in medical field: A brief review," *Curr. Med. Res. Pract.*, vol. 9, no. 3, pp. 102–109, 2019, doi: 10.1016/j.cmrp.2019.04.001.
- [8] D. E. Adkins, "Machine learning and electronic health records: A paradigm shift," *Am. J. Psychiatry*, vol. 174, no. 2, pp. 93–94, 2017, doi: 10.1176/appi.ajp.2016.16101169.
- [9] Y. Liang, Z. Chen, G. Liu, and M. Elgendi, "A new, short-recorded photoplethysmogram dataset for blood pressure monitoring in China," *Sci. Data*, vol. 5, no. 1, p. 180020, 2018, doi: 10.1038/sdata.2018.20.
- [10] B. Everett and A. Zajacova, "Gender differences in hypertension and hypertension awareness among young adults," *Biodemography Soc. Biol.*, vol. 61, no. 1, pp. 1–17, 2015, doi: 10.1080/19485565.2014.929488.
- [11] H. M. Choi, H. C. Kim, and D. R. Kang, "Sex differences in hypertension prevalence and control: Analysis of the 2010-2014 Korea national health and nutrition examination survey," *PLoS One*, vol. 12, no. 5, pp. 1–12, 2017, doi: 10.1371/journal.pone.0178334.
- [12] F. Landi *et al.*, "Body mass index is strongly associated with hypertension: Results from the Longevity Check-up 7+ Study," *Nutrients*, vol. 10, no. 12, p. 1976, 2018.
- [13] F. B. Hossain, S. R. Shawon, G. Adhikary, and A. Chowdhury, "Association between body mass index (BMI) and hypertension in South Asian population: Evidence from Demographic and Health Survey," *bioRxiv*, pp. 1–9, 2019, doi: 10.1101/605469.
- [14] G. Yang *et al.*, "Body weight and weight change in relation to blood pressure in normotensive men," *J. Hum. Hypertens.*, vol. 21, no. 1, pp. 45–52, 2007, doi: 10.1038/sj.jhh.1002099.
- [15] L. A. Colangelo, Y. Yano, D. R. Jacobs, and D. M. Lloyd-Jones, "Association of Resting Heart Rate with Blood Pressure and Incident Hypertension over 30 Years in Black and White Adults: The CARDIA Study," *Hypertension*, pp. 692–698, 2020, doi: 10.1161/HYPERTENSIONAHA.120.15233.
- [16] L. A. Matura *et al.*, "Physical activity and symptoms in pulmonary arterial hypertension," *Chest*, vol. 150, no. 1, pp. 46–56, 2016.
- [17] S. Nouri Saeidlou, F. Babaei, and P. Ayremlou, "Malnutrition, overweight, and obesity

- among urban and rural children in north of west azerbaijan, Iran,” *J. Obes.*, vol. 2014, no. July, 2014, doi: 10.1155/2014/541213.
- [18] D. E. Burmaster and E. A. C. Crouch, “Log-normal distributions for body weight as a function of age for males and females in the United States, 1976-1980,” *Risk Anal.*, vol. 17, no. 4, pp. 499–505, 1997, doi: 10.1111/j.1539-6924.1997.tb00890.x.
 - [19] C. Pater, “The blood pressure ‘uncertainty range’ - A pragmatic approach to overcome current diagnostic uncertainties (II),” *Curr. Control. Trials Cardiovasc. Med.*, vol. 6, pp. 1–16, 2005, doi: 10.1186/1468-6708-6-5.
 - [20] R. W. Makuch, D. H. Freeman, and M. F. Johnson, “Justification for the log-normal distribution as a model for blood pressure,” *J. Chronic Dis.*, vol. 32, no. 3, pp. 245–250, 1979, doi: 10.1016/0021-9681(79)90070-5.
 - [21] Y. Ostchega, K. S. Porter, J. Hughes, C. F. Dillon, and T. Nwankwo, “Resting pulse rate reference data for children, adolescents, and adults: United States, 1999-2008,” *Natl. Health Stat. Report.*, no. 41, pp. 1999–2008, 2011.
 - [22] S. Poole and N. Shah, “Incorporating Observed Physiological Data to Personalise Pediatric Vital Sign Alarm Thresholds,” *Biomed. Inform. Insights*, vol. 11, p. 117822261881847, 2019, doi: 10.1177/1178222618818478.
 - [23] A. D. Penman and W. D. Johnson, “The changing shape of the body mass index distribution curve in the population: Implications for public health policy to reduce the prevalence of adult obesity,” *Prev. Chronic Dis.*, vol. 3, no. 3, pp. 3–6, 2006.
 - [24] R. Kissell and J. Poserina, *Advanced Math and Statistics*. 2017.
 - [25] R. S. Society, “Thomas Bayes’ s Bayesian Inference Author (s): Stephen M . Stigler Source : Journal of the Royal Statistical Society . Series A (General) , Vol . 145 , No . 2 (1982) , pp . ,” vol. 145, no. 2, pp. 250–258, 2013.
 - [26] D. Fink, “A Compendium of Conjugate Priors,” *Mont. Mag. West. Hist.*, no. 1994, pp. 1–47, 1997, [Online]. Available: <http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.157.5540&rep=rep1&type=pdf>.
 - [27] S. Wang and W. Gui, “Corrected maximum likelihood estimations of the log-normal distribution parameters,” *Symmetry (Basel)*, vol. 12, no. 6, 2020, doi: 10.3390/SYM12060968.
 - [28] P. R. Kileen, “Beyond statistical inference: A decision theory for science,” *Psychon. Bull. Rev.*, vol. 13, no. 4, pp. 549–562, 2006, doi: 10.3758/BF03193962.
 - [29] G. Mancia *et al.*, “Blood pressure and heart rate variabilities in normotensive and hypertensive human beings,” *Circ. Res.*, vol. 53, no. 1, pp. 96–104, 1983, doi: 10.1161/01.RES.53.1.96.