

Sorting algorithms: how are they affected by memory faults?

Alexandro Vladno

Centre for Informatics and Systems

University of Coimbra

Coimbra, Portugal

alexandro.vladno@gmail.com

Fabiano Papaiz

Centre for Informatics and Systems

University of Coimbra

Coimbra, Portugal

fabianopapaiz@gmail.com

Leo Moreira Silva

Centre of Informatics and Systems

University of Coimbra

Coimbra, Portugal

leo.moreira@me.com

Abstract

Although most people in the world use technology devices for many tasks, they don't know how the devices work and how they deal with faults. When those faults occur in memory, software behavior could be affected. Together with the software-specific algorithms are the sorting algorithms used to solve problems like ordering a list of products by their price. This work presents a discussion about how quicksort, mergesort, insertionsort and bubblesort algorithms are affected by memory faults. The statistical tests showed that bubblesort, in general, was the less affected algorithm considering the variables involved in this study. On the other hand, insertionsort was the worse algorithm.

Index Terms

sorting algorithms, memory faults.

I. INTRODUCTION

Technology is deeply introduced in people's quotidian supporting a massive number of tasks, for example: searching for a shared car, surfing on the web, sending a message to someone, automating the company's production or using the company's software. Nevertheless, most people don't know that devices are continually dealing with memory failures, faults and errors. These devices were made with large and inexpensive memories, which are also error-prone [1].

Software behavior may be affected by the problems mentioned before, especially those from memory. We have a memory fault when the correct value that should be stored in a memory location gets altered because of a soft failure. In particular, the content of a location can change unexpectedly, i.e., faults may happen at any time: real memory faults are indeed highly dynamic and unpredictable [2].

In the beginning steps of software development, the designer has a general idea of the structure and functions. For each one of these, some algorithms will be produced or used. In the following stages, the outcome software (and its algorithms) will be tested and, then, delivered to the user. Different kinds of algorithms could be written or used in the software, and one of these is the sorting algorithms.

A good algorithm is that which gives satisfactory results for every range of data set. Sorting is a fundamental concept and important for solving other problems like is prerequisite for Binary Search. Sorting is often used in a large variety of critical applications and is a fundamental task that is used by most computers [3].

In this paper, we present a discussion about how these sorting algorithms, particularly Quicksort, Mergesort, Insertion Sort and Bubblesort, are affected by memory faults.

II. BACKGROUND

In this section we describe basic concepts about memory faults and sorting algorithms.

A. Memory Faults

Even the best digital system, with high-quality components and design techniques, may not be infallible to faults. Despite the title of this subsection, when the entire digital system (or software) is considered, there are three terms for computing fault and they have different meanings: failure, fault and error [4].

- *Error*: An error is a manifestation of a fault in a system, in which the logical state of an element differs from its intended value. An error occurs for a particular system state and input when an incorrect next state and/or output results.
- *Fault*: A fault is an anomalous physical condition. Causes include design errors, manufacturing problems, damage, fatigue, or other deterioration. Faults resulting from design errors and external factors are especially difficult to model and protect against because their occurrences and effects are hard to predict. A fault in a system does not necessarily result in an error;
- *Failure*: A failure denotes an element's inability to perform its functions because of error in the element itself or its environment, which in turn are caused by various faults;

B. Sorting Algorithms

Sorting algorithms are widely used in many aspects of data processing, information searches, business finance, computer encryption, etc. This work uses four sorting algorithms: quicksort, mergesort, insertionsort, and bubblesort. In the following subsections, we'll give an overview of them.

1) *Quicksort*: Quicksort algorithm, created by Hoare [5], is considered as one of the fastest and best sorting algorithms [6]. The algorithm is based on the paradigm of divide and conquer.

This algorithm has a execution time of $\theta(n^2)$ in the worst case over n numbers as input. Despite that execution time, quicksort is often the best option for sorting because of its remarkable average efficiency: $\theta(n \lg n)$ [7].

The basic steps of this algorithm are [6]:

- Pick an element, which is called a pivot, from the list waiting to be sorted;
- Perform partition operation to realize that all elements in the list with values smaller than the pivot came before the pivot. Otherwise, all elements in the list with values bigger than the pivot come after it (elements which are equal to pivot can go either way). After this partition, the pivot is in the final position of the list;
- Recursively sort the sub-list of smaller elements and the sub-list of the bigger elements.

2) *Mergesort*: Mergesort was invented by John Von Newman and is one of the most elegant algorithms to appear in the sorting literature. It is the first sorting algorithm to have $\theta(n \lg n)$ execution time bound. It is important to observe that this algorithm spends a lot of time on data transfer operations. In fact, standard Mergesort incurs about $2n$ data move operations [8].

Conceptually, Mergesort works as follows [8]:

- Divide the unsorted array into two sub arrays of about half the size;
- Sort each sub array recursively;
- Merge the two sub arrays back into one array.

3) *Insertionsort*: This algorithm sorts the array by shifting the elements one at time. It is efficient in sorting a small number of elements. The overall execution time of this algorithm is $\theta(n^2)$ [7]. The basic sorting steps are:

- If there are more than one element, pick the next element;
- Compare with all the elements in sorted sub-list;
- Shift all the elements in sorted sub-list that is greater than the value to be sorted;

- Insert the value;
- Repeat until list is sorted.

4) *Bubblesort*: The bubble sort is the oldest and simplest sorting method in use. It works by comparing each item in the list with the item next to it, and swapping them if required. The algorithm repeats this process until it makes a pass all the way through the list without swapping any items (in other words, all items are in the correct order) [9].

Table I below shows the time complexity comparison between the sorting algorithms presented. The n is the number of input elements.

TABLE I: Sorting algorithms complexity time comparison [10]

Algorithm	Time Complexity		
	<i>Best Case</i>	<i>Average Case</i>	<i>Worst Case</i>
Bubblesort	$O(n^2)$	$O(n^2)$	$O(n^2)$
Insertionsort	$O(n)$	$O(n^2)$	$O(n^2)$
Quicksort	$O(n \lg n)$	$O(n \lg n)$	$O(n^2)$
Mergesort	$O(n \lg n)$	$O(n \lg n)$	$O(n \lg n)$

III. EXPERIMENTAL SETUP

In this section, we present the setup of this study. We first state our problem, then define the independent and dependent variables. After that, we show the hypothesis and describe our dataset, showing all its characteristics.

A. Problem Statement

As introduced in the first section of this paper, sorting is a fundamental concept and essential for solving other problems. The content of memory location can change unexpectedly, i.e., faults may happen at any time. Considering this, the main objective of this work is to design experiments to answer the following question: *How are sorting algorithms affected by memory faults?*

B. Variables

For this experimental study, we assume that the independent and dependent variables are as shown in Table II and Table III below:

TABLE II: Independent variables.

Variable	Description
Probability of failure	Probability of a fault to occur
Array size	Size of the array of integers to be sorted
Sorting algorithm	Algorithm used to sort the array

TABLE III: Dependent variables.

Variable	Description
Largest subarray size	Size of the largest sorted subarray produced under the memory fault
Percentage of largest subarray size	Percentage of <i>largest subarray size</i> related to <i>array size</i> independent variable
Unordered elements quantity	Quantity of elements out of position after sorting algorithm execution. Adapted of the <i>k-unordered sequence</i> measure of disorder defined in [11]
Percentage of unordered elements quantity	Percentage of <i>unordered elements quantity</i> related to <i>array size</i> independent variable

C. Hypothesis

The set of hypothesis defined to test and draw some conclusions about this experiment are listed below. The confidence degree defined for hypothesis testing was 95% ($\alpha = 0.05$ and $\alpha - 1 = 0.95$).

- **Hypothesis 1:** For a given probability of failure and array size, tested algorithms will produce a different percentage of unordered elements quantity.
- **Hypothesis 2:** For a given probability of failure and array size, tested algorithms will produce a different percentage of the largest subarray size.
- **Hypothesis 3:** For each algorithm, the array size and probability of failure have a significative impact on the percentage of unordered elements quantity.
- **Hypothesis 4:** For each algorithm, the array size and probability of failure have a significative impact on the percentage of the largest subarray size.

D. Dataset

To conduct the proposed study, we define the values of the independent variables, as shown in Table IV:

TABLE IV: Values of the independent variables.

Variable	Values
Probability of failure	1%, 2% and 5%
Array size	100, 1000 and 10000
Sorting algorithm	Bubblesort, Quicksort, Mergesort and Insertionsort

Based on these variables, we ran an existing script *gen.py* to produce input files. We define that our sample was composed by 30 input files for a given combination of the probability of failure and array size. So, considering this, we ran 30 times for each combination of these independent variables, producing 30 inputs, totalizing 270 files. Figure 1 shows an example of produced input files:

0.01	100	9	48	37	6	26	7	24	44	17	50	48	30	49	33	22	13	42	29	39	13	19	13	9	28	34	1	33	27	14	45	48	40	11	17	6	50	9	44	20	16	37	45	23	14	38	29	10	49	44	46	35	45	15	2	22	1	46	40	8	48	23	23	32	35	3	15	8	36	17	24	27	48	28	5	28	50	44	4	25	6	9	1	11	44	26	50	44	12	7	20	30	20	37	20	6	8	13	15	20	49
------	-----	---	----	----	---	----	---	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	---	----	----	---	----	----	----	----	----	----	----	----	---	----	---	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	---	----	---	----	----	---	----	----	----	----	----	---	----	---	----	----	----	----	----	----	---	----	----	----	---	----	---	---	---	----	----	----	----	----	----	---	----	----	----	----	----	---	---	----	----	----	----

Fig. 1: Example of input file.

The input data shown in the Figure 1 is divided as follows:

- *Probability of Failure*: the first number of the sequence (0.01) is the probability of memory failure when sorting;
- *Sequence size*: the second number (100) means the size of the integers sequence used by sorting;
- *Sequence*: the rest of the numbers represents a list of n positive integers that will be sorted.

With this input data, we ran, for each one of these, all four algorithms considered in this study. The sorting algorithms used already existed. For example, using all 270 input files, we ran bubblesort, creating 270 output files, and so on for the other algorithms. At the end of executions, we get a total of 1080 output files. An output file looks like shown in Figure 2:

```
#1  9 48 37 6 26 7 24 44 17 50 48 30 49 33 22 13 42 29 39 13 19 13 9 28 34 1
33 27 14 45 48 40 11 17 6 50 9 44 20 16 37 45 23 14 38 29 10 49 44 46 35 45
15 2 22 1 46 40 8 48 23 23 32 35 3 15 8 36 17 24 27 48 28 5 28 50 44 4 25 6 9
1 11 44 26 50 44 12 7 20 30 20 37 20 6 8 13 15 20 49
#2  1 1 1 2 3 4 5 6 6 6 6 7 7 8 8 8 9 9 9 9 10 11 11 12 13 13 13 13 14 14 15
15 15 16 17 17 17 19 20 20 20 20 22 22 23 23 23 24 24 25 26 26 27 27 28 28 28
29 29 30 30 32 33 33 34 35 35 36 37 37 37 38 39 40 40 42 44 44 44 44 44 44 20
45 45 45 46 46 48 48 48 48 48 49 49 50 49 50 50 50
#3  1 1 1 2 3 4 5 6 6 6 6 7 7 8 8 8 9 9 9 9 10 11 11 12 13 13 13 13 14 14 15
15 15 16 17 17 17 19 20 20 20 20 20 22 22 23 23 23 24 24 25 26 26 27 27 28 28
28 29 29 30 30 32 33 33 34 35 35 36 37 37 37 38 39 40 40 42 44 44 44 44 44 44
45 45 45 46 46 48 48 48 48 48 49 49 49 50 50 50 50
#4  82
```

Fig. 2: Example of output file.

The output file gives four essential data, as enumerated below:

- *#1*: the original sequence of integers contained in the input file;
- *#2*: the sequence processed by the sorting algorithm under the memory fault model;
- *#3*: the sequence sorted correctly;
- *#4*: the size of the largest sorted subsequence in item 2. This number can be interpreted as a quality measure.

As higher, most successful was the sorting operation.

After generating the dataset, we developed a Python program that reads the 1080 output files and produces a single CSV file (first lines showed in Figure 3). This file contains the columns listed below and was used to perform the exploratory data analysis and to run the statistical tests.

- *sorting_algorithm*: the algorithm used to sort the array;
- *probability_of_failure*: the probability of failure used when sorting;
- *array_size*: the size of the sorted array;
- *largest_subarray_size*: the largest sorted subarray after sorting;
- *unordered_elements_quantity*: quantity of unordered sequence elements after sorting.
- *perc_unordered_elements_quantity*: percentage of unordered elements after sorting related to original array;
- *perc_largest_subarray_size*: percentage of largest sorted subarray after sorting related to original array.

```

sorting_algorithm;probability_of_failure;array_size;largest_subarray_size;
unordered_elements_quantity;perc_unordered_elements_quantity;
perc_largest_subarray_size
quick;0.01;100;35;4;4.00;35.00
quick;0.01;100;36;8;8.00;36.00
quick;0.01;100;31;5;5.00;31.00
quick;0.01;100;20;6;6.00;20.00
quick;0.01;100;31;5;5.00;31.00

```

Fig. 3: Example of output CSV file.

We use Python libraries to make data analysis and plot graphs. These libraries were:

- *Pandas*¹: open source library providing data structures and data analysis tools;
- *NumPy*²: library for scientific computing with Python;
- *SciPy*³: ecosystem of open-source software for mathematics, science, and engineering;
- *StatsModels*⁴: module that provides classes and functions for the estimation of many different statistical models, as well as for conducting statistical tests, and statistical data exploration;
- *Matplotlib*⁵: plotting library;
- *Seaborn*⁶: data visualization library based on matplotlib.

IV. DATA ANALYSIS

In this section, we present our results after the execution of sorting algorithms over the input files. We analyzed only two of the four dependent variables, which were *percentage of the largest subarray size (%LSS)* and *percentage of unordered elements quantity (%UEQ)*. These variables, because they are a percentage value, already were normalized (i.e., the same order of magnitude) related to dependent variable *array size*.

A. Exploratory Data Analysis (EDA)

Firstly, we performed an analysis of the distribution of the dependent variables %LSS and %UEQ. To help in this task, we produced histograms, boxplot graphs, tables containing data about mean, median, standard deviation, and the minimum and maximum values.

The following Figures 4, 5 and 6 illustrates examples of histograms and boxplot graphs for each combination of *Algorithm X Probability of Failure X Array Size*. In each of those figures, the graphs were exhibited over the dependent variables %LSS and %UEQ.

¹<https://pandas.pydata.org>

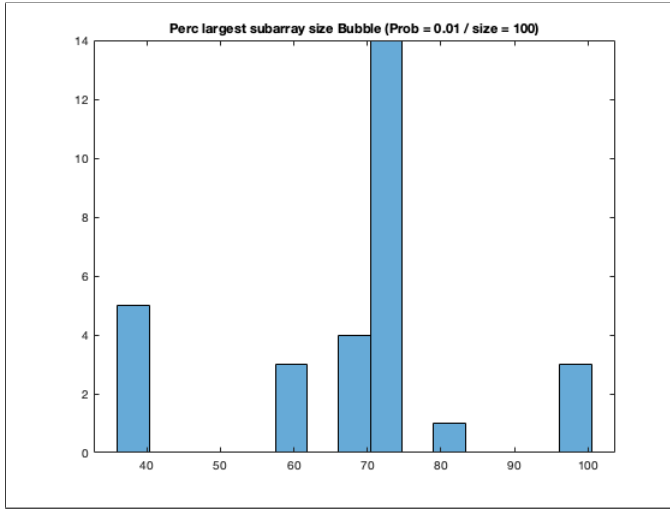
²<https://numpy.org>

³<https://www.scipy.org>

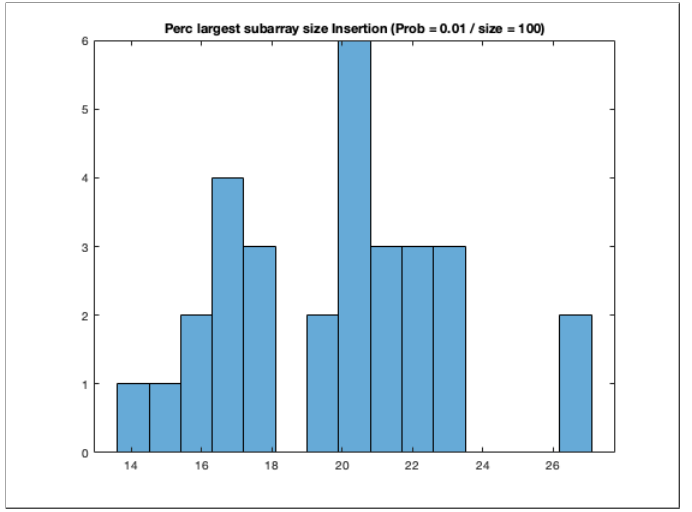
⁴<https://www.statsmodels.org>

⁵<https://matplotlib.org>

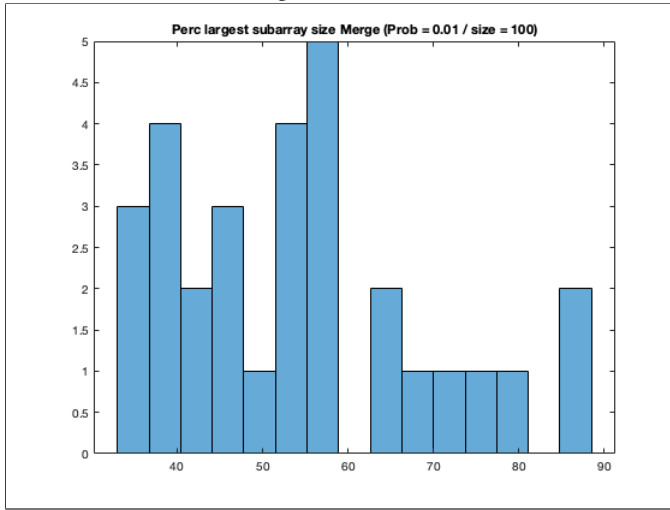
⁶<https://seaborn.pydata.org>



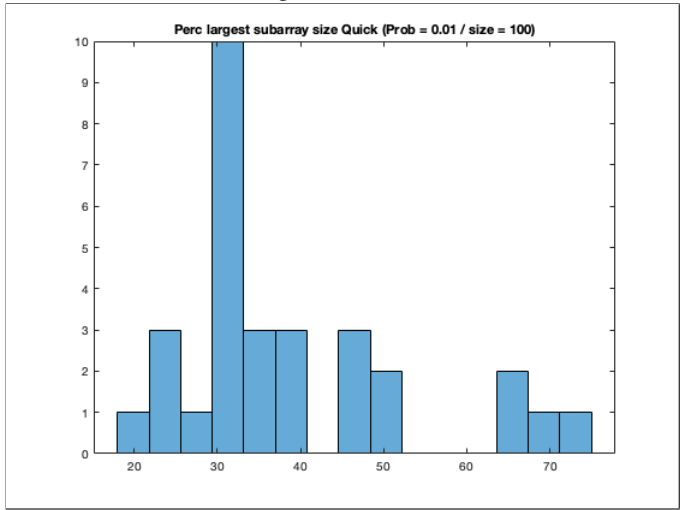
(a) Histogram for Bubblesort.



(b) Histogram for Bubblesort.



(c) Histogram for Mergesort.



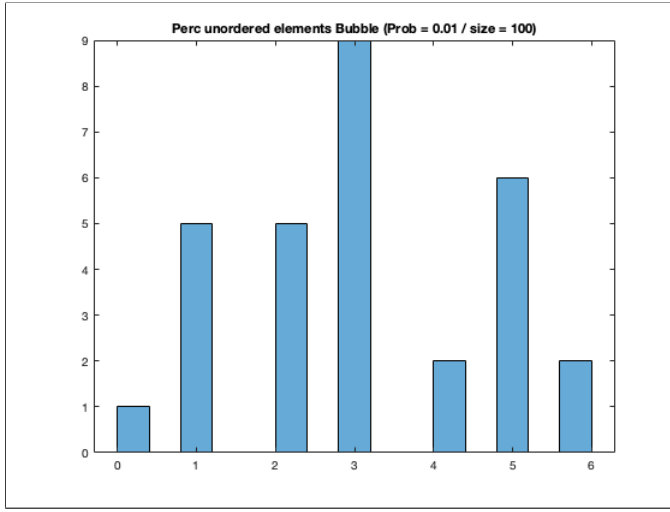
(d) Histogram for Quicksort.

Fig. 4: Histograms for %LSS, with *probability of failure* of 0.01 and *array size* of 100.

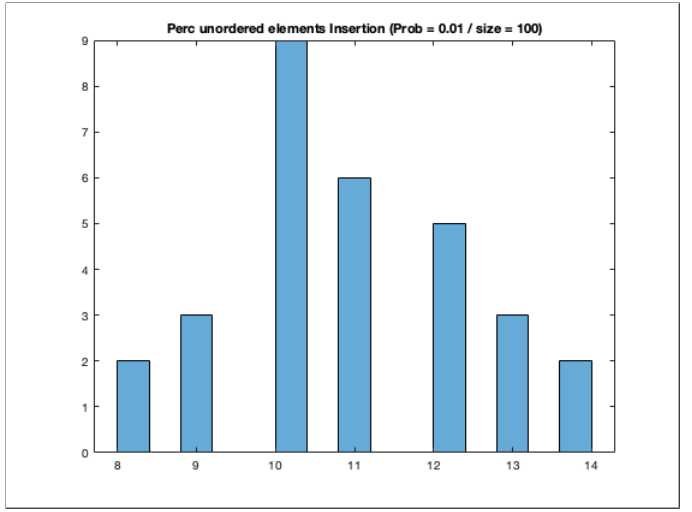
Based on Figures 4, 5 and 6, the following Tables V and VI illustrates the information about the dependent variables (%LSS and %UEQ) distribution.

TABLE V: Dependent variable %UEQ distribution.

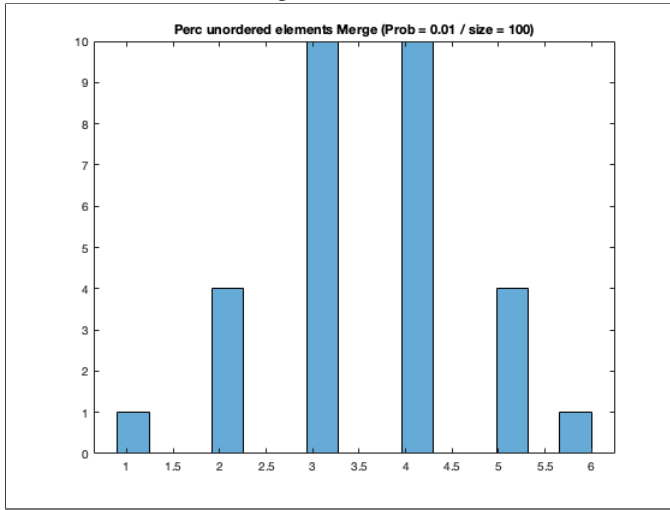
Prob. of Failure	Array Size	Algorithm	Percentage of unordered elements quantity (%UEQ)				
			Mean	Median	Std. Deviation	Minimum	Maximum
0.01	100	bubble	3.07	3.0	1.64	0.0	6.0
0.01	100	insertion	10.87	11.0	1.59	8.0	14.0
0.05	10000	merge	28.32	28.34	0.31	27.69	28.86
0.05	10000	quick	19.56	19.58	0.26	18.95	20.07



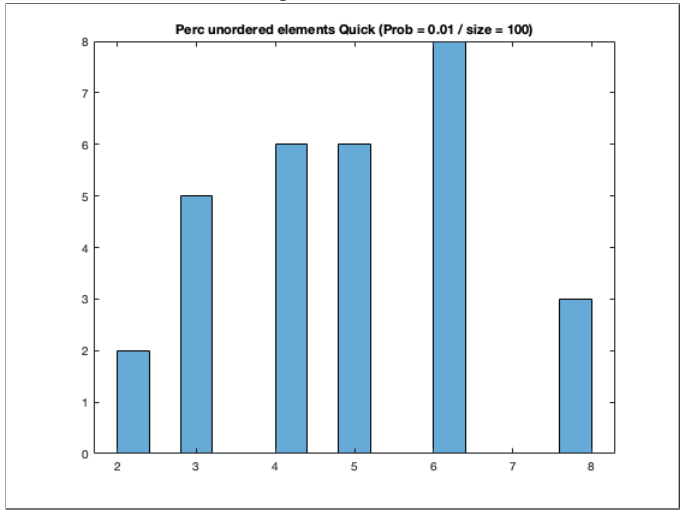
(a) Histogram for Bubblesort.



(b) Histogram for Bubblesort.



(c) Histogram for Mergesort.



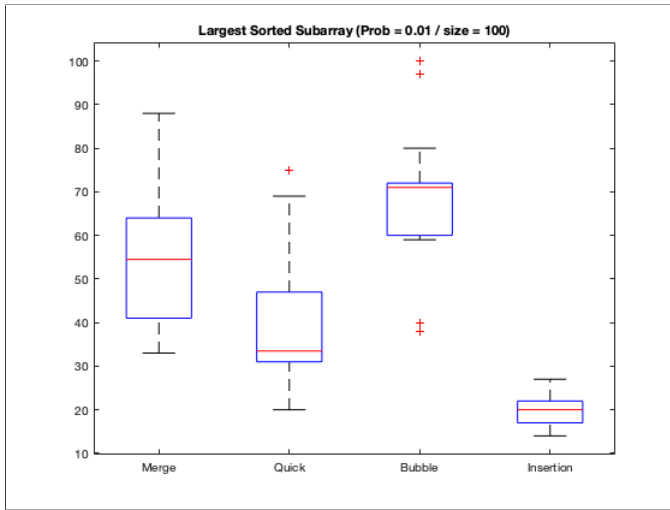
(d) Histogram for Quicksort.

Fig. 5: Histograms for %UEQ, with *probability of failure* of 0.01 and *array size* of 100.

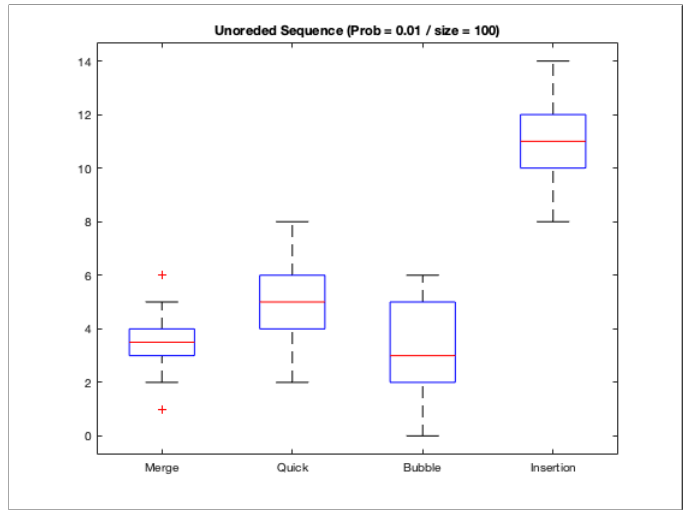
TABLE VI: Dependent variable %LSS distribution.

Prob. of Failure	Array Size	Algorithm	Percentage of largest subarray size (%LSS)				
			Mean	Median	Std. Deviation	Minimum	Maximum
0.01	100	bubble	67.4	71.0	15.84	38.0	100.0
0.01	100	insertion	19.77	20.0	3.11	14.0	27.0
0.05	10000	merge	0.19	0.2	0.03	0.15	28.86
0.05	10000	quick	0.3	0.3	0.03	0.23	0.34

We produced graphs with the same data showed in Tables V and VI. Figure 7 shows an example of these graphs. On it, we illustrate on the top-left graph that, considering the mean, bubblesort produces fewer unordered elements related to the total than other algorithms.



(a) Boxplot for %LSS.



(b) Boxplot for %UEQ.

Fig. 6: Boxplots for %LSS and %UEQ variables, with *probability of failure* of 0.01 and *array size* of 100.

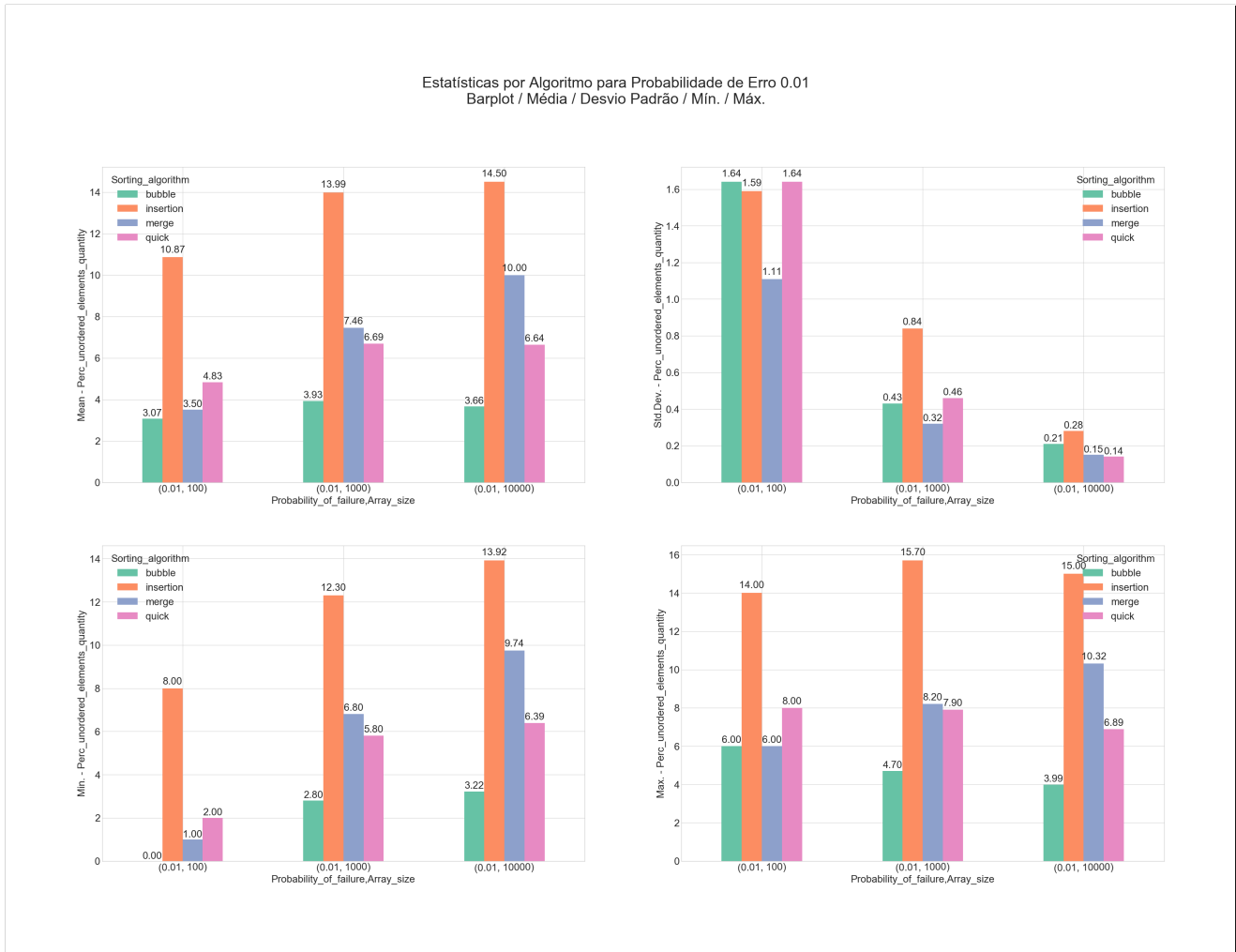


Fig. 7: Distribution data with *probability of failure* of 0.01, all algorithms and *array sizes*.

B. Dependent Variables Normality

An essential step to continue the analysis is to verify if dependent variables have a normal distribution. In this work, we produced the Q-Q plot to analyze the normality visually. We also ran the Shapiro-Wilk [12] test to confirm if they follow a normal distribution or not.

After the Shapiro-Wilk test, we determine that only the dependent variable %UEQ (*percentage of unordered elements quantity*) has a normal distribution related to mean for all algorithms. This situation occurred considering all possible combinations between independent variables: *sorting algorithm*, *probability of failure* and *array size*.

Based on these results, we choose to test just the hypothesis associated with the dependent variable %UEQ. We made use of the ANOVA method to test the hypothesis, and this method is premised the normal distribution of variables values.

Tables VII and VIII shows examples of results we got running Shapiro-Wilk test. The independent variables assume these values: *probability of failure* of 0.01 and *array size* of 100. The bold *p-values* confirms that the data is normally distributed.

TABLE VII: Shapiro-Wilk test for %LSS with *probability of failure* of 1% and *array size* of 100.

Algorithm	W	p-value
Bubblesort	0.854840	0.0007
Insertion Sort	0.961790	0.3439
Mergesort	0.937173	0.0763
Quicksort	0.869708	0.0016

TABLE VIII: Shapiro-Wilk test for %UEQ with *probability of failure* of 1% and *array size* of 100.

Algorithm	W	p-value
Bubblesort	0.934980	0.0666
Insertion Sort	0.949881	0.1678
Mergesort	0.936667	0.0739
Quicksort	0.936049	0.0712

Q-Q plot shows that how much more blue points close to the red line, most normal is the distribution. Figure 8 below presents this graph for the dependent variable *percentage of the largest subarray size* considering a *probability of failure* of 1% and a *array size* of 100 for all sorting algorithms.

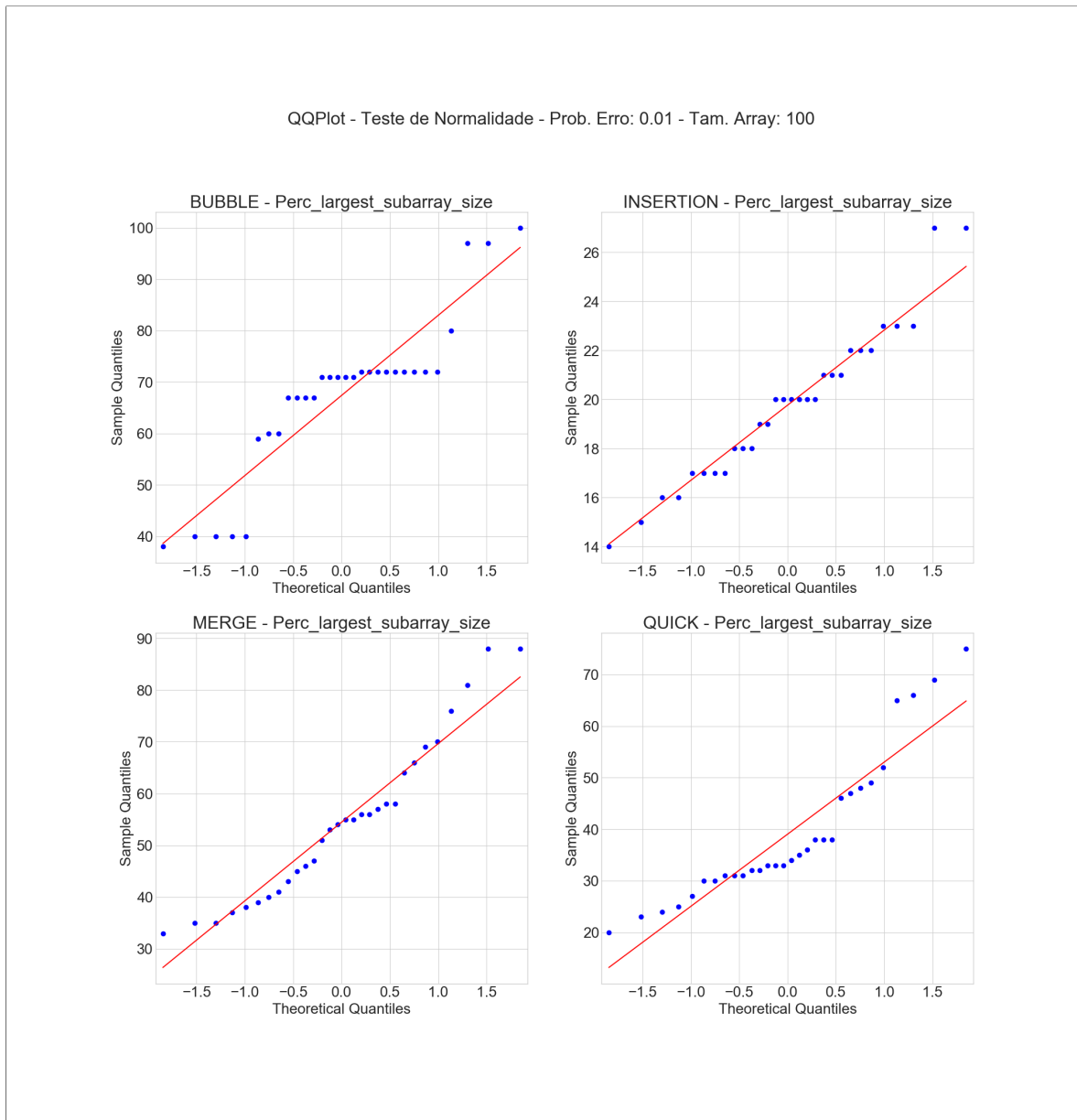


Fig. 8: Q-Q plot for %LSS with *probability of failure* of 1% and *array size* of 100.

Figure 9 shows the same graph for the dependent variable *percentage of unordered elements quantity* considering same values for independent variables.

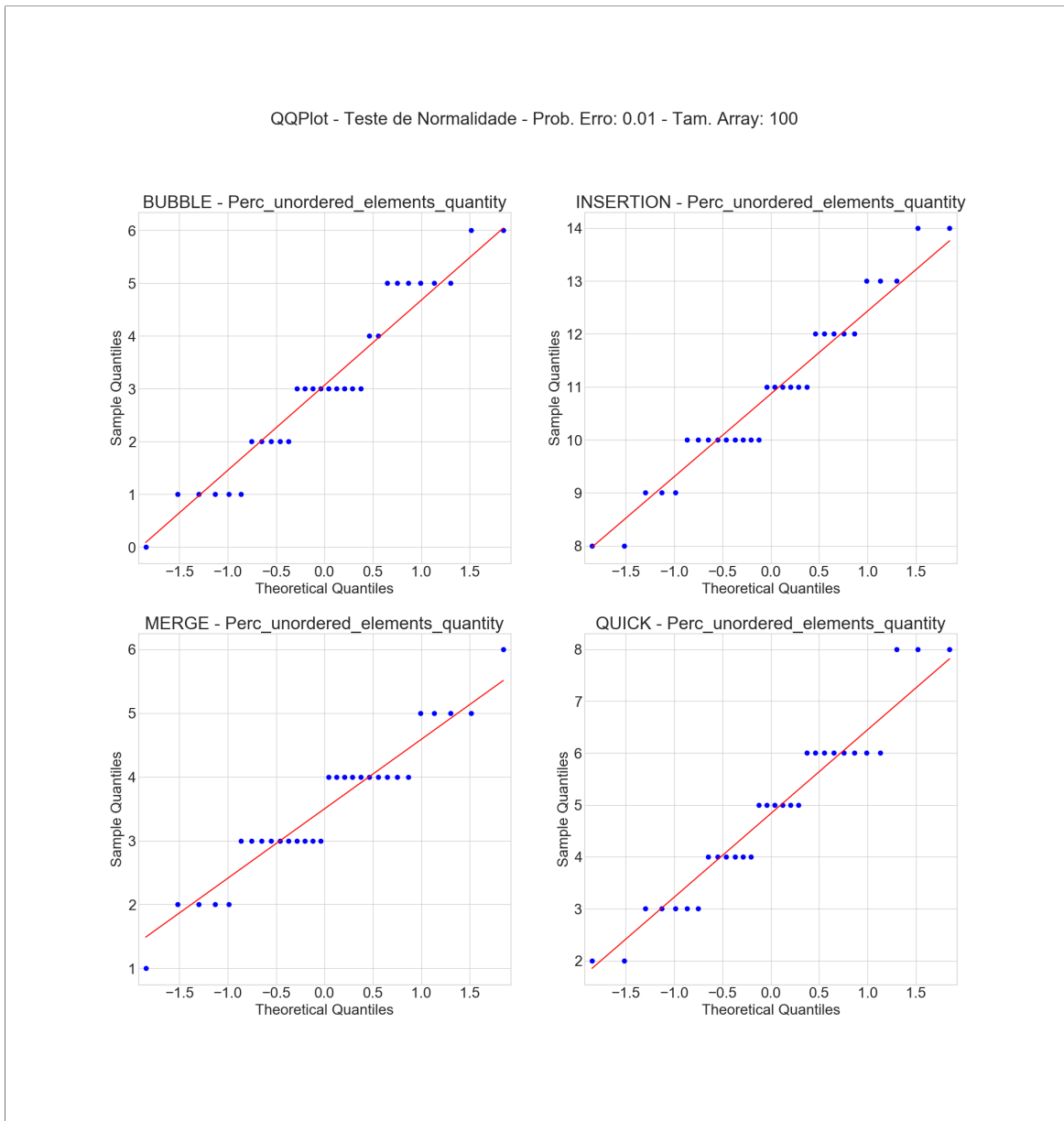


Fig. 9: Q-Q plot for %UEQ with *probability of failure* of 1% and *array size* of 100.

Tables IX and X shows another examples of results we got running Shapiro-Wilk test. The independent variables assume these values: *probability of failure* of 0.05 and *array size* of 1000. The bold *p-values* confirms that the data is normally distributed.

TABLE IX: Shapiro-Wilk test for %LSS with *probability of failure* of 5% and *array size* of 1000.

Algorithm	W	p-value
Bubblesort	0.891642	0.0052
Insertion Sort	0.921918	0.0300
Mergesort	0.878308	0.0025
Quicksort	0.917583	0.0232

TABLE X: Shapiro-Wilk test for %UEQ with *probability of failure* of 5% and *array size* of 1000.

Algorithm	W	p-value
Bubblesort	0.936459	0.0730
Insertion Sort	0.931681	0.0544
Mergesort	0.961920	0.3465
Quicksort	0.963954	0.3892

Other examples of distribution normality can be seen in Q-Q plot in Figures 10 and Y below. The considered dependent variable is *percentage of the largest subarray size* with a *probability of failure* value of 5% and *array size* value of 10000 for all sorting algorithms.

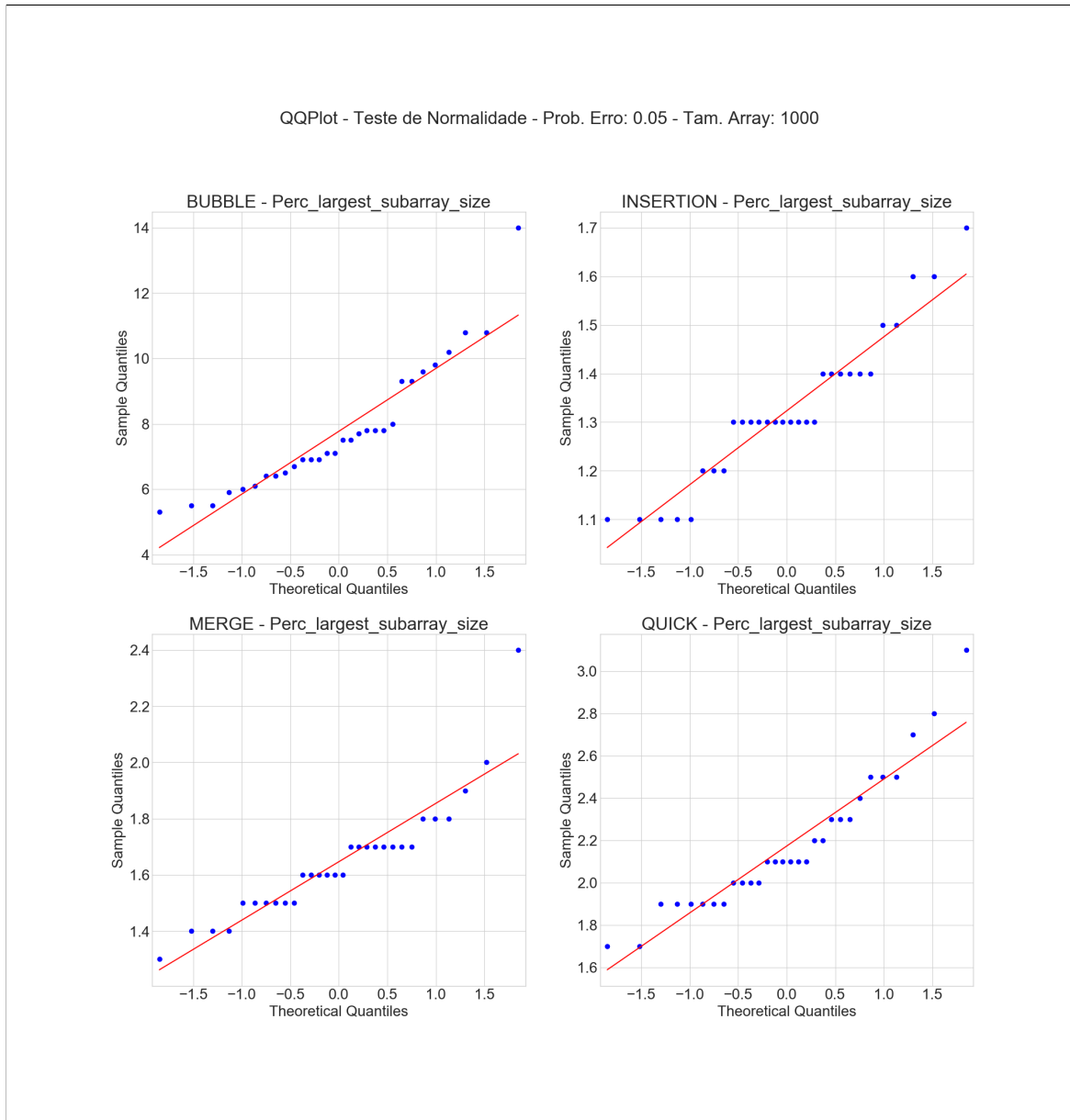


Fig. 10: Q-Q plot for %LSS with *probability of failure* of 5% and *array size* of 1000.

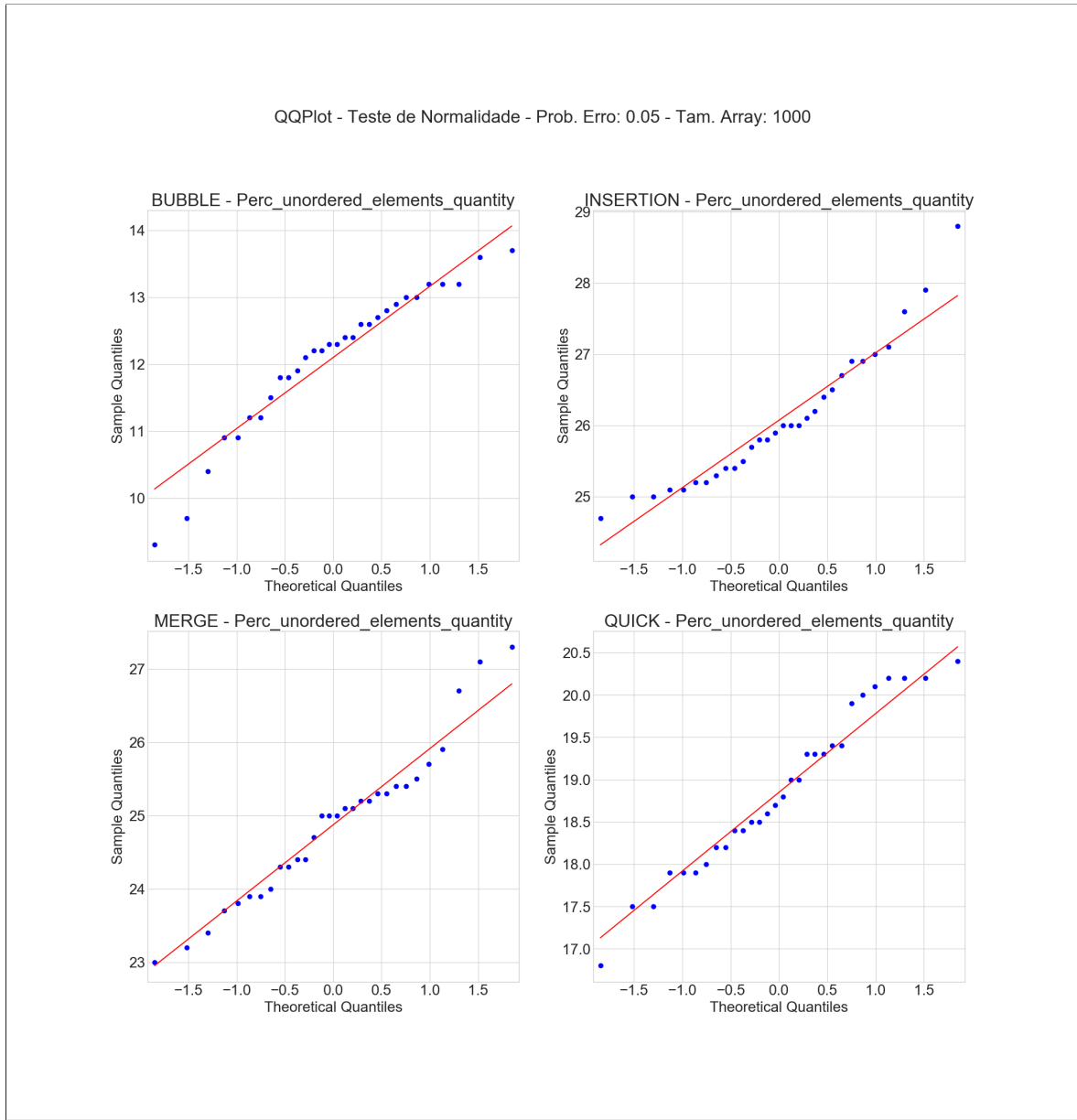


Fig. 11: Q-Q plot for %UEQ with *probability of failure* of 5% and *array size* of 1000.

C. Hypothesis Testing

For hypothesis testing, we define a confidence degree of 95%. As mentioned before, we only tested the hypothesis related to the dependent variable *%UEQ percentage of unordered elements quantity*. The hypothesis was:

- **Hypothesis 1:** For a given probability of failure and array size, tested algorithms will produce a different percentage of unordered elements quantity;
- **Hypothesis 3:** For each algorithm, the array size and probability of failure have a significative impact on the percentage of unordered elements quantity.

Because our dataset has 3 independent variables, we perform ANOVA method to simplify the hypothesis testing. This method allows comparing various groups of variables at the same time.

1) *Testing Hypothesis 1:* We formulated the following hypothesis for the combination of independent variables *probability of failure* and *array size*:

Hypothesis H_0 : There is no significant difference between algorithms related to variable %UEQ.

Hypothesis H_1 : There is significant difference between algorithms related to variable %UEQ.

After we run the testing method, we got the following results (Table XI):

TABLE XI: Hypothesis 1 ANOVA results.

	Probability of Failure = 1% and Array Size = 100			
	<i>sum_sq</i>	<i>df</i>	<i>F</i>	<i>PR(>F)</i>
algorithm	1174.466667	3.0	171.368721	1.845762e-42
residual	265.000000	116.0		
	Probability of Failure = 1% and Array Size = 1000			
	<i>sum_sq</i>	<i>df</i>	<i>F</i>	<i>PR(>F)</i>
algorithm	1635.180250	3.0	1817.615579	2.610034e-97
residual	34.785667	116.0		
	Probability of Failure = 1% and Array Size = 10000			
	<i>sum_sq</i>	<i>df</i>	<i>F</i>	<i>PR(>F)</i>
algorithm	1948.204989	3.0	15831.708335	2.332995e-151
residual	4.758210	116.0		
	Probability of Failure = 2% and Array Size = 100			
	<i>sum_sq</i>	<i>df</i>	<i>F</i>	<i>PR(>F)</i>
algorithm	1106.5	3.0	88.617785	7.054900e-30
residual	482.8	116.0		
	Probability of Failure = 2% and Array Size = 1000			
	<i>sum_sq</i>	<i>df</i>	<i>F</i>	<i>PR(>F)</i>
algorithm	2286.83025	3.0	1990.144336	1.507899e-99
residual	44.43100	116.0		
	Probability of Failure = 2% and Array Size = 10000			
	<i>sum_sq</i>	<i>df</i>	<i>F</i>	<i>PR(>F)</i>
algorithm	3081.306943	3.0	18442.105239	3.406264e-155
residual	6.460427	116.0		
	Probability of Failure = 5% and Array Size = 100			
	<i>sum_sq</i>	<i>df</i>	<i>F</i>	<i>PR(>F)</i>
algorithm	2274.291667	3.0	225.58173	3.104916e-48
residual	389.833333	116.0		
	Probability of Failure = 5% and Array Size = 1000			
	<i>sum_sq</i>	<i>df</i>	<i>F</i>	<i>PR(>F)</i>
algorithm	3704.037583	3.0	1202.21591	3.625400e-87
residual	119.132333	116.0		
	Probability of Failure = 5% and Array Size = 10000			
	<i>sum_sq</i>	<i>df</i>	<i>F</i>	<i>PR(>F)</i>
algorithm	4948.821397	3.0	19487.120298	1.402080e-156
residual	9.819533	116.0		

Based on these results, we reject the null hypothesis (H_0) for all testing combinations. Therefore, we can affirm with 95% of confidence that there is a significant difference between algorithms related do variable %UEQ.

2) *Testing Hypothesis 3*: We formulated the following hypothesis for all algorithms considered in this work:

Hypothesis H_0 : The values of variables *probability of failure* and *array size* don't have a significant impact on values of %UEQ variable.

Hypothesis H_1 : The values of variables *probability of failure* and *array size* have a significant impact on values of %UEQ variable.

After we run the testing method, we got the following results (Table XII):

TABLE XII: Hypothesis 3 ANOVA results

	Bubblesort			
	<i>sum_sq</i>	<i>df</i>	<i>F</i>	<i>PR(>F)</i>
size_of_array	0.019548	2.0	0.01350	0.90759
probability_of_failure	3590.644437	2.0	2479.71637	7.570027e-137
residual	385.169623	266.0		
	Mergesort			
	<i>sum_sq</i>	<i>df</i>	<i>F</i>	<i>PR(>F)</i>
size_of_array	1784.873256	2.0	383.585800	1.702008e-53
probability_of_failure	13009.963175	2.0	2795.961630	3.801570e-143
residual	1237.731651	266.0		
	Quicksort			
	<i>sum_sq</i>	<i>df</i>	<i>F</i>	<i>PR(>F)</i>
size_of_array	121.551362	2.0	32.964196	2.556033e-08
probability_of_failure	5626.955128	2.0	1526.005558	3.457996e-112
residual	980.841817	266.0		
	Insertionsort			
	<i>sum_sq</i>	<i>df</i>	<i>F</i>	<i>PR(>F)</i>
size_of_array	309.843991	2.0	75.905868	3.236586e-16
probability_of_failure	6806.075083	2.0	1667.358570	1.414758e-116
residual	1085.798822	116.0		

Based on these results, we reject the null hypothesis (H_0) for algorithms Mergesort, Quicksort, and Insertion sort. Therefore, we can affirm with 95% of confidence that the values of variables *probability of failure* and *array size* have a significant impact on values of %UEQ variable.

For Bubblesort algorithm, however, there is no confidence to reject the null hypothesis (H_0). Analyzing the result, we can affirm with 95% of confidence that the value of *probability of failure* variable had a significant impact on the value of the dependent variable %UEQ. But we can not affirm the same related to independent variable *array size*. These results demonstrate that, for the Bubblesort algorithm, the array sizes used in this experiment were not decisive to define values concerning the percentage of unordered elements quantity.

D. Performance analysis related to percentage of unordered elements quantity variable

In this section, we present an analysis of overall performance related to dependent variable *percentage of unordered elements quantity*. To help this analysis, we show below 3 bar graphs, each of these for a different value of *probability of failure* independent variable. These graphs show information about values of %UEQ variable grouped by *array size* and *sorting algorithm* variables. Figures 12, 13, and 14 below exhibit these graphs.

The bars in each graph represent the sample mean value for variable %UEQ inside each group. Above each bar, there is a vertical line that indicates the lower and upper limits of a confidence interval for a significance level of 5% ($\alpha = 0.05$).

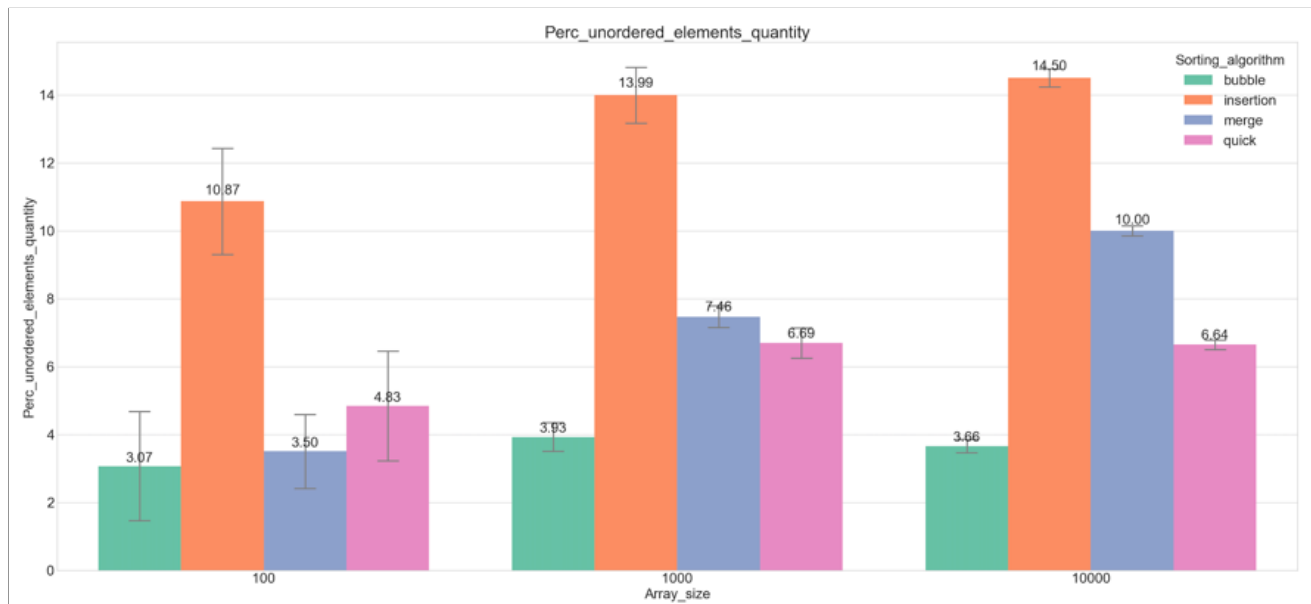


Fig. 12: Data for *probability of failure* of 1%.

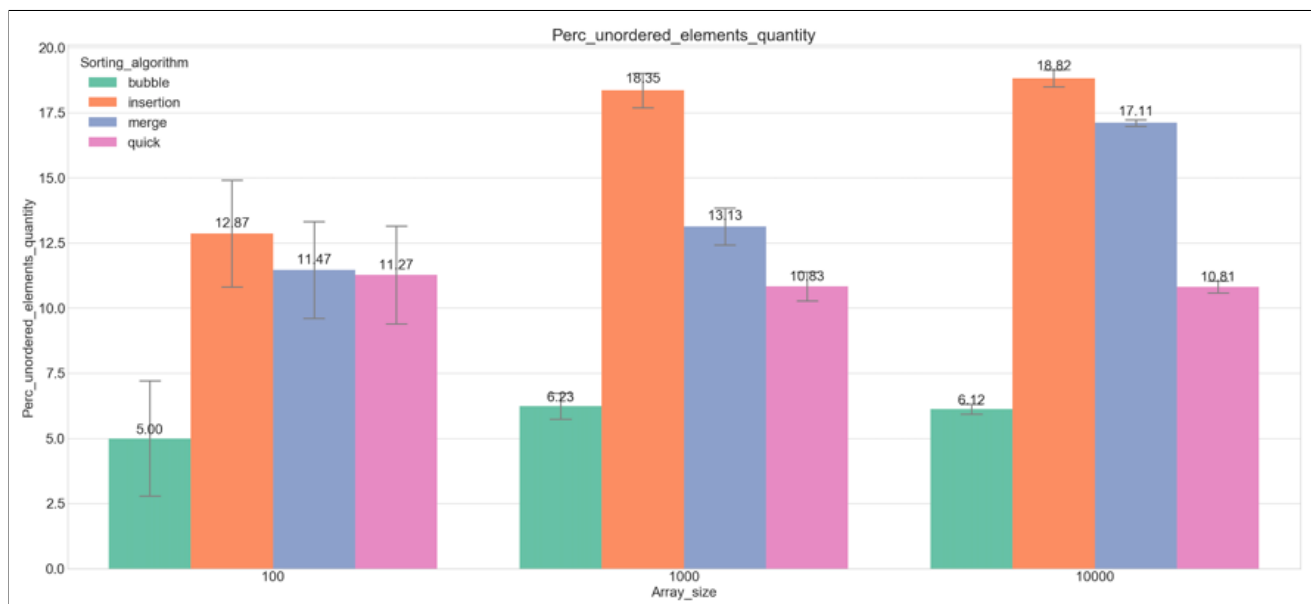


Fig. 13: Data for *probability of failure* of 2%.

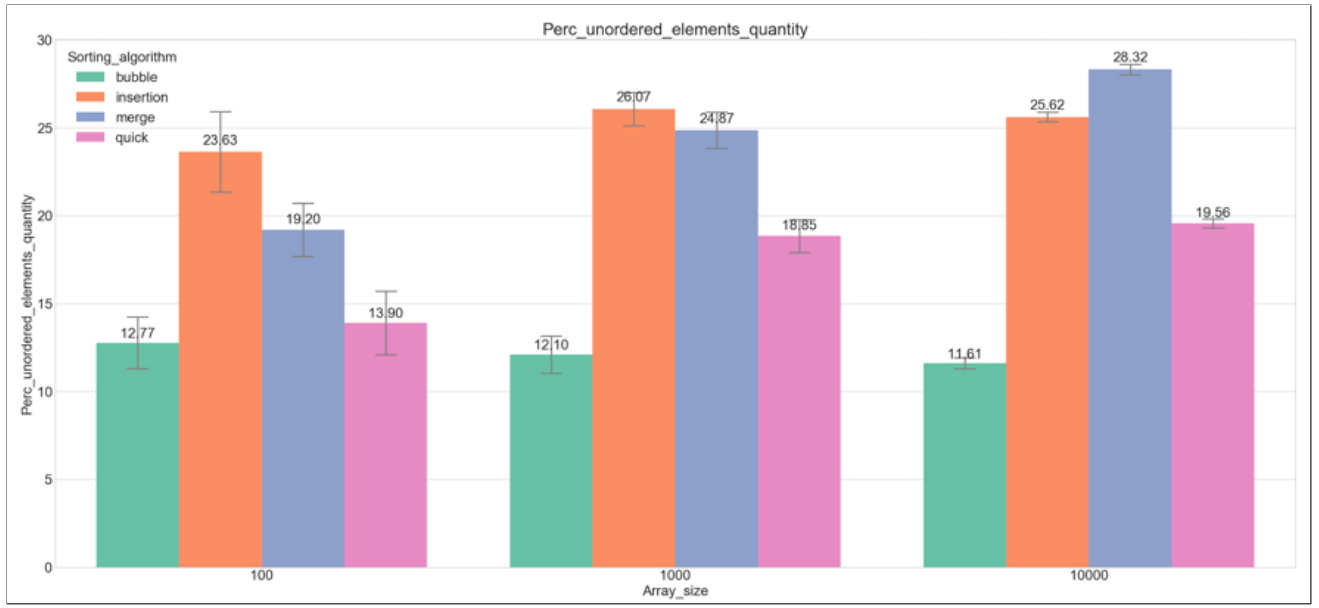


Fig. 14: Data for *probability of failure* of 5%.

Considering the data in the graphs, we can conclude that Bubblesort generates lower mean value for variable %UEQ for all defined combinations (*probability of failure* X *array size* X *sorting algorithm*). These data demonstrate that Bubblesort was the less affected algorithm by the memory faults simulated in this experiment. We believe that this fact happened because of the operation of the algorithm itself, where all elements order are verified each iteration, causing an element incorrectly positioned in an iteration (because of a memory fault) can have its location corrected in an upcoming iteration.

Related to Bubblesort yet, we could identify that despite different sizes of input arrays (100, 1000, and 10000), there was not a significant impact on the mean value obtained for %UEQ variable, proving what we verify with ANOVA when doing hypothesis 3 testing for this algorithm.

V. LINEAR REGRESSION

The linear regression analysis consists of statistical analysis to verify the existence of a functional relationship between a dependent variable with one or more independent variables. In other words, this analysis consists of obtaining a linear equation that tries to explain the variation from the dependent variable by variation from the independent variable level.

The statistical model for this case should be:

$$y = a.x + b \quad (1)$$

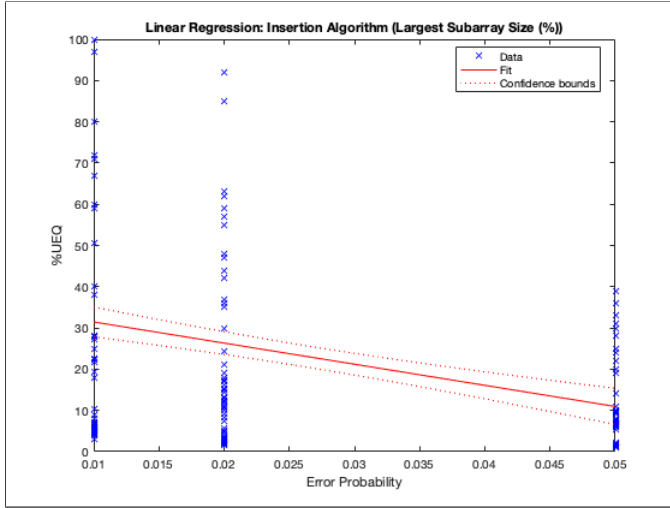
Where y represents the dependent variable, a represents the line slope of the linear model, and b represents the y-axis intercept. The regression line is a method to calculate the parameters a and b .

The determination coefficient, frequently called R^2 , or simply r^2 , for the linear regression case, provide auxiliary information to the analysis of the result of regression variance, as a way to verify if the proposed model is suitable or not to describe the phenomenon.

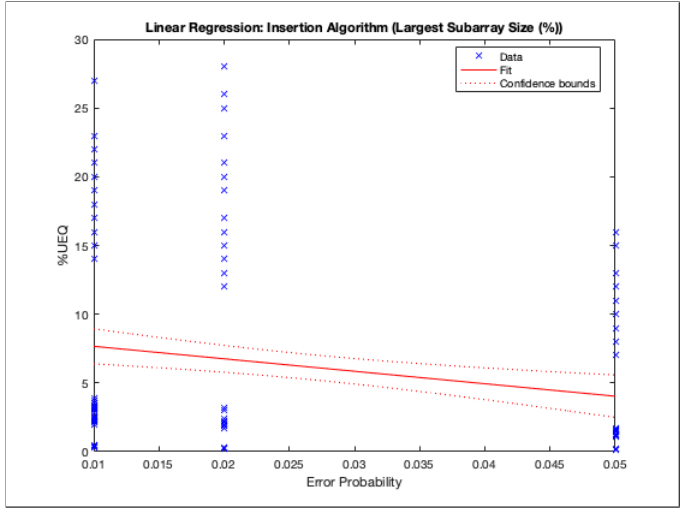
In this work, we verified how close is the relationship between the *probability of failure* independent variable and the other dependent variables: %LSS and %UEQ to a linear equation.

To discover an equation that represents the studied phenomenon, we could build a graph, called *Scatter Plot*, to verify by visualization how is the variation between the dependent and independent variables.

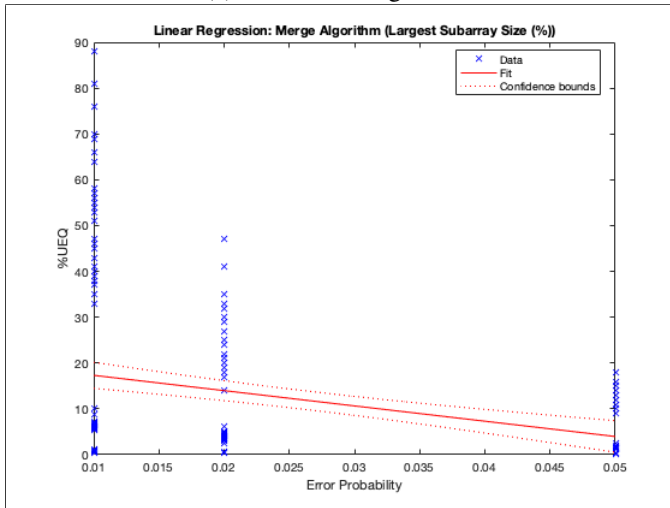
To discover an equation that represents the studied phenomenon, we could build a graph, called scatter plot, to verify by visualization how is the variation between the dependent and independent variables. Figures 15 and 16 are the scatter plots for %LSS and %UEQ for each algorithm, respectively. In these figures, we can see the line from the linear equation from the estimated model that used the minimum mean square error as criteria. The red dotted line represents the bounds for 95% of the confidence interval.



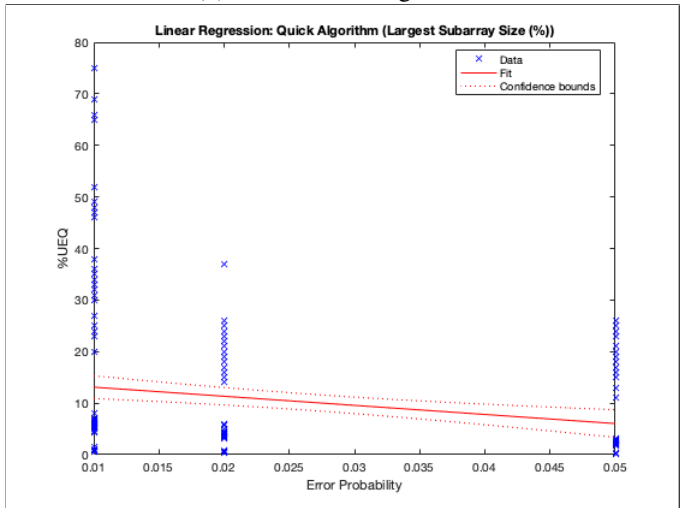
(a) Bubblesort algorithm.



(b) Insertionsort algorithm.

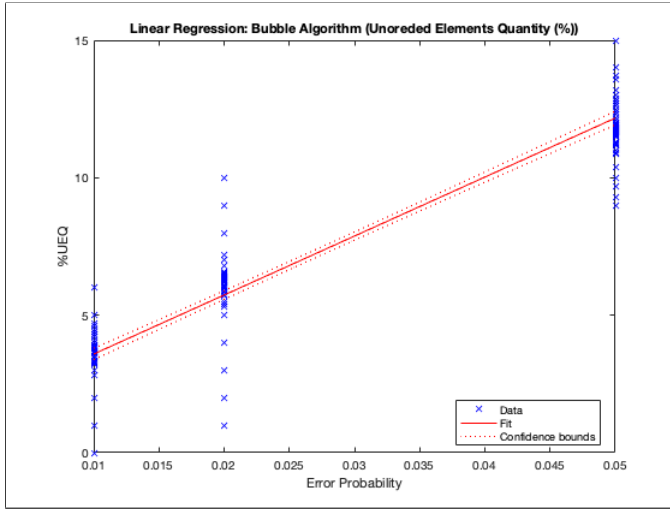


(c) Mergesort algorithm.

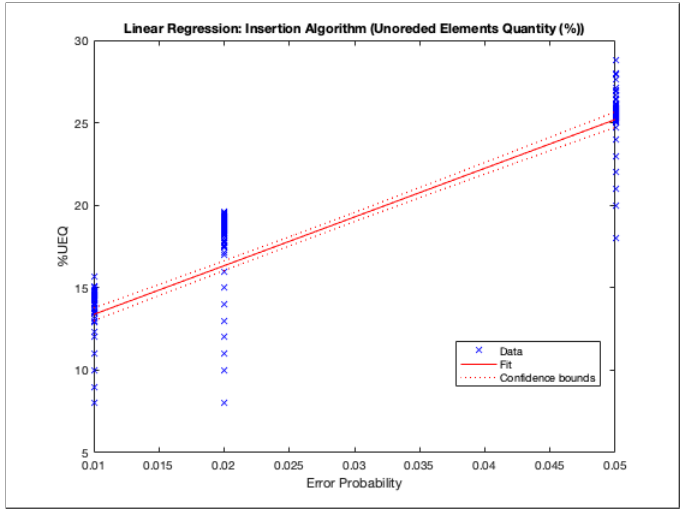


(d) Quicksort algorithm.

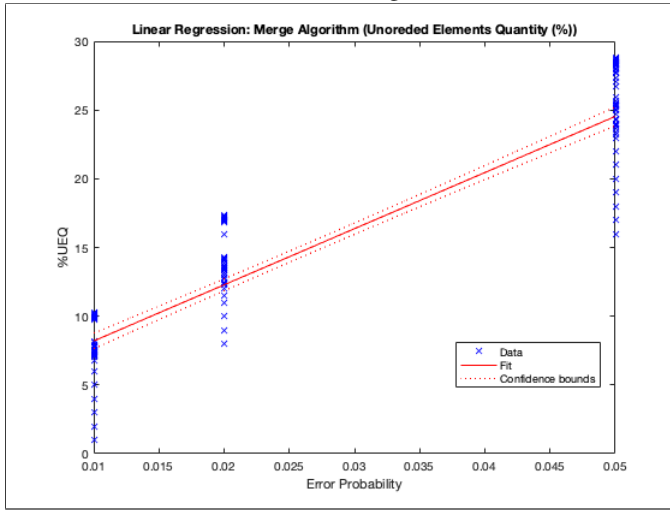
Fig. 15: Scatter plot and the linear regression for %LSS versus Probability of Failure.



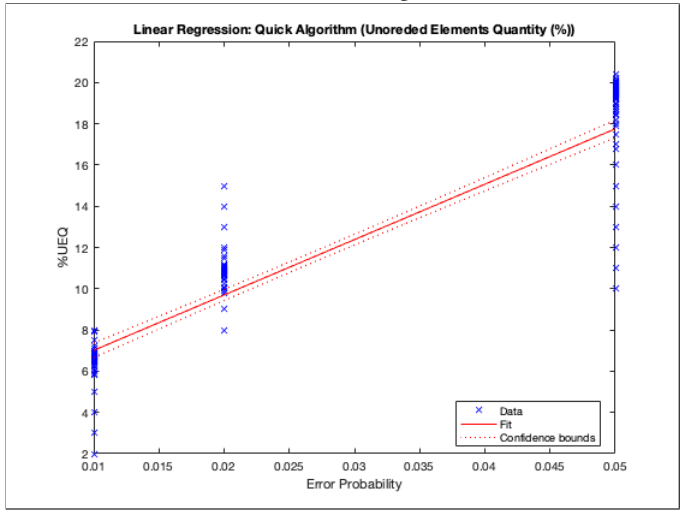
(a) Bubblesort algorithm.



(b) Insertionsort algorithm.



(c) Mergesort algorithm.



(d) Quicksort algorithm.

Fig. 16: Scatter plot and the linear regression for %UEQ versus Probability of Failure.

However, we can verify that the points at scatter plots, don't have perfect adjust to the proposed math model line. There are, in major part, a great distance between the points at the plot and the model line. This distance happens because, in each probability of failure, we have many values of the dependent variable and with large dispersion. This situation explains why the determination coefficient r^2 has a low value in most of the cases.

The Tables XIII and XIV below shows the parameters for linear regression of both variables %LSS and %UEQ. For %UEQ variable, we found better (from 0.809 to 0.899) r^2 than %LSS (from 0.0499 to 0.141) variable for all algorithms.

TABLE XIII: Parameters from linear regression for the %LSS variable.

Parameter	Bubblesort	Quicksort	Mergesort	Insertionsort
Intercept (b)	36.576	14.86	20.651	8.5707
Slope (a)	-512.32	-176.4	-333.64	-90.711
r^2	0.141	0.0499	0.101	0.0399
Root Mean Squared Error (RMSE)	21.5	13.1	17	7.59

TABLE XIV: Parameters from linear regression for the %UEQ variable.

Parameter	Bubblesort	Quicksort	Mergesort	Insertionsort
Intercept (b)	1.4429	4.3258	4.1165	10.426
Slope (a)	-214.56	268.59	408.4	295.39
r^2	0.899	0.823	0.809	0.826
Root Mean Squared Error (RMSE)	1.23	2.12	3.38	2.31

VI. CONCLUSION

This work proposes a study and discussion of how sorting algorithms, particularly Quicksort, Mergesort, Insertion Sort, and Bubblesort, are affected by memory faults. To achieve this, we define the experimental setup, showing the dependent and independent variables, then the used dataset with its characteristics. Next, we performed data analysis after the execution of sorting algorithms over the dataset. As explained in this work, we analyzed only the dependent variables *percentage of the largest subarray size* (%LSS) and *percentage of unordered elements quantity* (%UEQ). These variables, because they are a percentage value, already were normalized (i.e., the same order of magnitude) related to dependent variable *array size*.

We ran a normality test and, after that, we determined that only the dependent variable %UEQ (*percentage of unordered elements quantity*) has a normal distribution related to mean for all algorithms, so we chose to test just the hypothesis associated with the dependent variable %UEQ.

After data analysis and discussion, we can conclude that, in this experiment, the worse algorithm related to %UEQ variable was Insertion sort. This algorithm had the highest mean values in almost all combinations of study. These values were much higher than other algorithms when tested with *probability of failure* of 1% (0.01). Finally, our tests show that, in general, Quicksort algorithm was better than Mergesort.

For linear regression, our results showed that for dependent variable %UEQ variable we get a linear model that explains at least 80% of relation between this variable and independent variable probability of failure.

We can, then, with our results, elaborate a ranking of the performance of algorithms when considering memory faults:

- 1) Bubblesort
- 2) Quicksort
- 3) Mergesort
- 4) Insertion sort

All data used for this work can be downloaded at <https://bit.ly/392LZKR>.

REFERENCES

- [1] I. Finocchi and G. F. Italiano, "Sorting and searching in the presence of memory faults (without redundancy)," in *Conference Proceedings of the Annual ACM Symposium on Theory of Computing*, no. January, pp. 101–110, 2004.

- [2] S. Hamdioui, Z. Al-Ars, A. J. Van De Goor, and M. Rodgers, "Dynamic faults in random-access-memories: Concept, fault models and tests," in *Journal of Electronic Testing: Theory and Applications (JETTA)*, vol. 19, pp. 195–205, 2003.
- [3] Nitin Arora and S. Kumar, "A Novel Sorting Algorithm and Comparison with Bubble sort and Insertion sort," *International Journal of Computer Applications*, vol. 45, no. 1, p. 2.
- [4] V. P. Nelson, "Fault-Tolerant Computing: Fundamental Concepts," *Computer*, vol. 23, no. 7, pp. 19–25, 1990.
- [5] C. A. R. Hoare, "Quicksort," *The Computer Journal*, vol. 5, no. 1, pp. 10–16, 1962.
- [6] X. Wang, "Analysis of the time complexity of quick sort algorithm," *Proceedings - 2011 4th International Conference on Information Management, Innovation Management and Industrial Engineering, ICIII 2011*, vol. 1, pp. 408–410, 2011.
- [7] T. Cormen, C. Leiserson, R. Rivest, and C. Stein, *Introduction to algorithms*. MIT Press, 2009.
- [8] D. Abhyankar and M. Ingle, "A Novel Mergesort," vol. 1, no. 3, pp. 17–22, 2011.
- [9] K. Mansotra, V. Sourabh, "Implementing Bubble Sort Using a New Approach," pp. 1–6, 2011.
- [10] P. Prajapati, N. Bhatt, and N. Bhatt, "Performance Comparison of Different Sorting Algorithms," vol. VI, no. Vi, pp. 39–41, 2017.
- [11] U. Ferraro-Petrillo, I. Finocchi, and G. F. Italiano, "The price of resiliency: A case study on sorting with memory faults," *Algorithmica (New York)*, vol. 53, no. 4, pp. 597–620, 2009.
- [12] S. S. Shapiro and M. B. Wilk, "An analysis of variance test for normality (complete samples)," *Biometrika*, vol. 52, no. 3/4, pp. 591–611, 1965.