

UNIVERSITY OF TARTU  
FACULTY OF MATHEMATICS AND COMPUTER SCIENCE  
Institute of Computer Science

Stenver Jerkku

# Paralell Wilcoxon Signed-rank tests

Bachelor's Thesis (6 ECTS)

Supervisor: Sven Laur, PhD

Tartu 2012

## **Abstract**

Wilcoxon Signed-rank test is a paired statistical test that is used when measurement values are not normally distributed. The test is used by BIIT(Bioinformatics, Algorithmics and Data mining group) research group for gene regulation, gene expression data analysis, biological data mining and others. BIIT is joint research group between the Department of Computer Science (University of Tartu), Quretec, and the Estonian Biocenter.

The current implementations of the Wilcoxon Signed-rank tests are slow and unoptimized. This project will look into the foundations of wilcoxon signed-rank test, its current implementations and how to optimize it. In order to make the implementation more accurate, the relationship between Wilcoxon statistic and Guassian approximate will be observed. In order to make the implementation faster, some dynamic programming methods will be used to save computation time. The goal of optimizing is to make it more accurate and speed up the test running.

The end goal of this project is to create an accurate and fast wilcoxon test implementation in C++ shared library. In the scope of this project, the library will be integrated with two tools -command line and Cron-R. Due to the nature of shared library, it will be easy integrate the library with any other tools one might desire.

# Contents

<b>1</b>	<b>Introduction</b>	<b>4</b>
1.1	NetCDF file format . . . . .	4
1.2	Current BIIT procedures . . . . .	6
1.2.1	GNU R . . . . .	6
1.2.2	Web interface . . . . .	6
1.2.3	Command line . . . . .	6
<b>2</b>	<b>Wilcoxon signed-rank test</b>	<b>7</b>
2.1	Statistical Tests . . . . .	7
2.1.1	Test Statistic . . . . .	7
2.2	Hypothesis . . . . .	8
2.2.1	Null hypothesis . . . . .	8
2.2.2	Alternative hypothesis . . . . .	8
2.2.3	P-value . . . . .	9
2.3	Wilcoxon test . . . . .	10
2.4	Wilcoxon test Assumptions . . . . .	10
2.5	How to compute . . . . .	11
2.6	Example . . . . .	12
2.7	P-value assignment . . . . .	12
2.7.1	V and P example . . . . .	14
2.7.2	Exact computation of P-value . . . . .	14
2.7.3	The V and P table formulas . . . . .	15
2.7.4	The Gaussian approximation of P table . . . . .	15
2.7.5	Code sample . . . . .	15
2.7.6	Bonferroni correction . . . . .	17
<b>3</b>	<b>Approximating the p-value</b>	<b>18</b>
3.1	Motivation . . . . .	18
3.2	When can we approximate . . . . .	19
3.3	Investigating the $P_{N,k}$ and $P_{approxN,K}$ similarities . . . . .	24
3.4	Finding out when can we approximate p-value . . . . .	27
<b>4</b>	<b>Further optimizations</b>	<b>28</b>
4.0.1	dose-response curve . . . . .	29
4.1	Linear approximation . . . . .	29
4.1.1	Relative values between P tables . . . . .	30
4.1.2	Linear approximation algorithm . . . . .	30
4.1.3	Linear interpolation . . . . .	32
4.2	Optimization summary . . . . .	32

4.3	Further optimizations . . . . .	33
<b>5</b>	<b>The implementaion</b>	<b>34</b>
5.1	RcppWilcoxonTest . . . . .	34
5.2	TerminalWilcoxonTest . . . . .	34
5.3	WilcoxonTestLibrary . . . . .	34
5.4	WilcoxonVTable . . . . .	35
5.5	Seminar_paper . . . . .	35
<b>6</b>	<b>Conclusion</b>	<b>36</b>
<b>7</b>	<b>Eestikeelne pealkiri</b>	<b>37</b>

# 1 Introduction

The BIIT research group has biologist who conduct a variety of experiments on a control group. They try to measure, compare and test different attributes of a living organism. These attributes might be, but not limited to blood pressure, proteine amount in the blood, RNA amount, purity of urine, brainwaves etc. These attributes can be repeadedly measured on the same subject and the results are never exactly the same. There is always a little noise and varience between the samples. Biologist often experiment on the control group and try to affect these observed attributes. The take samples before and after simulating the observable attribute with some stimulant. This is called the case-control study.

However, since even without external stimulation the two results are never the same. Also, different experiments can have wildy different assumptions, so you cannot simply use a constant error margin on all your tests. For example, the blood pressure measuring before and after jogging could have a lot bigger changes then brainwave changes. Therefore the biologist need to use statistical tests to find out which case-control studies are actually significant and which are simply noise.

There are many case-control tests available. For example paired Students t-test, t-test for matched pairs and Wilcoxon signed-rank test. The Wilcoxon signed-rank test is used when the measurement values are not normally distributed when the experimental conditions are fixed. In practise, this means that the histogram of measurements is either asymmetrical or it does not resemble the bell shape.

It is common to assume that properly scaled gene expression measurements follow Gaussian distribution, while quantities of proteins and metabolites are assumed to have non-gaussian distribution.

## 1.1 NetCDF file format

The BIIT group holds statistical data in NetCDF file format. NetCDF is an open source standard of a set of data formats, programming interfaces, and software libraries that help read and write scientific data files. [http://www.unidata.ucar.edu/software/netcdf/docs/what\\_is\\_netcdf.html](http://www.unidata.ucar.edu/software/netcdf/docs/what_is_netcdf.html)<sup>2</sup> Data can be held in a structured manner in a NetCDF file. You can define dimensions, name them, put variables in the dimensions and later retrieve or change them. This is useful, for example, to hold matrix like data in a file and have fast lookups. In addition, you can define helper dimensions for the matrix for extra information. NetCDF file format is used by the BIIT researh group to hold gene data and will be given as an input file for the program.

Currently the BIIT group holds the data in the NetCDF file in the following format:

## Data matrix

- double data(m, n) - the data matrix of m genes (variables) and n samples (experiments, patients, ...)
- string gene[m] - gene IDs
- string array[n] - sample IDs

## Non-track metadata

- string \_\_FileFormat - EVDF format identifier string. Current value: "ExpressView 1.1".
- string \_\_DatasetType - A string identifying the dataset type. This is used to extract data-specific parameters from a global configuration file. See datatypes.conf (ExpressViewConfigFiles) for allowed values.
- string[n] Organism - Organism (per column)
- string[n] MetadataOrder - Sometimes metadata (including column tracks) may be in a different order than columns in data. This variable is a permutation of the array variable, specifying the order of column variables. This field is deprecated and all column metadata is expected to be in array order in new datasets.
- string InvestigationTitle - short title for the dataset
- string ExperimentDescription - long description for the dataset
- string DatasetID - Optional Source-specific dataset identifier
- string DatasetLink - Optional A URL associated with the dataset
- string \_\_VariableLabel - Optional Custom variable type name for the dataset. Overridden by variable\_label in project configuration. (ExpressView) Names of any other (optional, non-track) metadata variables must be prepended by two underscores. Additional variables may be mandated by a data type specification, see ExpressViewConfigFiles. Known specific variables:
- string \_\_Organism - ArrayExpress chip organism (ArrayExpress; as opposed to the source of biological material stored in Organism)

## **Tracks**

Tracks are all other variables whose name begins with an uppercase character, that are arrays of either string or double, with the size of their first dimension either m or n. The former are referred to as row tracks and the latter as column tracks. These variables encode some information about either the rows or columns of the data matrix; e.g. ExpressView displays column tracks above the heatmap.

Examples:

string Chemotherapy[n] string Local Relapse[n] double RelativeVariance[m]

## **1.2 Current BIIT procedures**

Currently the BIIT group has 3 ways of analysing its data.

### **1.2.1 GNU R**

They can use GNU R project to make custom and complex analyses. For that they need to manually read NetCDF file to the R and then call the R built in function `wilcox.test` to use Wilcoxon test. The problem with R wilcoxon test is that it is slow. If you run thousands of tests in a row, the results might come after hours of computer processing. One of the goals of this paper and project is to make Wilcoxon test run on command line in mere seconds.

### **1.2.2 Web interface**

They can use a web interface which can take in certain arguments and command line script. It then will use whatever command line tools it has available and visualizes the output.

### **1.2.3 Command line**

They can also use the command line directly which has a variety of tools installed. The wilcoxon test will ultimately call R wilcoxon function.

So in the end - no matter how the biologist analyse their data, they will always use wilcoxon test in the GNU R project.

## 2 Wilcoxon signed-rank test

In this section, we discuss a particular statistical test called Wilcoxon signed-rank test. Wilcoxon test can be used to make sure that something has an effect on the measured samples. For example, biologist could measure 3 patients white blood cell levels, create some external stimulation on the patients and measure the white blood cell levels again. The wilcoxon signed-ranked test shows wether the external stimulation had any statistical effect on the white blood cell levels.

### 2.1 Statistical Tests

The statistical tests are usually done, because initial research has raised a hypothesis about the data. It is desirable to know if the hypothesis is true or not.

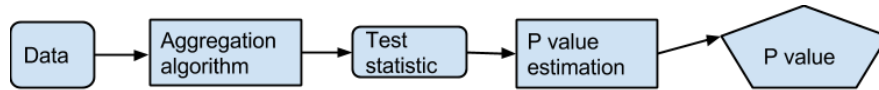


Figure 1: Graphical representation of statistical test flow

Statistical tests are usually built by first gathering the necessary raw data. For example, make 80 persons run for 30 minutes and measure blood pressure before and after running. The data is then organized by some kind of aggregation algorithm which also might filter out samples that are not related to the hypothesis in question. For example, sorting and grouping the blood pressure results by gender, weight and in ascending order and then filter out any group that had too few samples. Then a test statistic is chosen and calculated. After that a p-value can be estimated which tells if the observable feature is interesting or not. For example, if running plays a role in raising a persons blood pressure.

#### 2.1.1 Test Statistic

The test statistic is a scalar function of all the observations, which summarizes the data by a single number. This value is used to estimate a p-value and accept or reject a null hypothesis. Commonly, the test statistic is either one-sample, two-sample or paired statistic, depending on the tests.

A one-sample test tests whether the data collectively satisfies some hypothesis, e.g. the data is generated by gaussian distribution.

A two-sample test compares two sub-populations of data, e.g. persons with high blood persons and persons with low blood pressures.

A paired test characterises the change of before and after measurements, e.g. does some new medical drug treatment work. The differences between measured



observable case and control experiment pairs comes from the same distribution. The distribution of differences are also symmetrical. A positive difference  $x$  is equally likely as a negative difference  $-x$ .

## 2.2 Hypothesis

### 2.2.1 Null hypothesis

In statistics, a null and alternative hypothesis are raised to find out if the test bears any interesting results. A null hypothesis in general terms means that the test results are boring or that there is no effect, or in more formal terms, general position. It means there is no relationship between the two observed features. For example, when measuring blood pressure before and after running, the null hypothesis would state that blood pressure will not change when running.

The null hypothesis also defines the data distribution. For example, it might define that data is generated by a Gaussian distribution. Since null hypothesis is the general position, then generally the data will be distributed according to the null hypothesis. Since the null hypothesis defines the data distribution, it will also define the test statistic distribution. Finally, since the null hypothesis defines the data distribution and test statistic distribution, then it will also play a role in estimating the p-value for the test statistic.

For example, if the observations  $x_1, x_2, \dots, x_n$  are generated by coin flip, then the test statistic  $T = x_1 + x_2 + \dots + x_n$  has a binomial distribution  $Binomial_{n, \frac{1}{2}}$ . If test statistic is defined as  $T = x_1 * (1 - x_1) + x_2 * (1 - x_2) + \dots + x_n * (1 - x_n)$ , then it is constantly zero. The exact form of the test statistic determines the properties of the test.

### 2.2.2 Alternative hypothesis

The alternative hypothesis is the opposite to the null hypothesis - in simple terms, one could think of it as true or interesting result. It means that the two measured samples have a relationship between them, either in direct or indirect ways. A good example to explain these two is the court verdict. If null hypothesis is kept, then the suspect has not been found guilty. There was not enough evidence. On the contrary, if null hypothesis was rejected, the suspect was found guilty, because there was enough evidence to make that decision. Although the alternative hypothesis could be anything that is not a null hypothesis, some alternatives to the null hypothesis are easier to distinguish from the null hypothesis.

When the null hypothesis is true, then the test statistic has a well-defined probability distribution. This means that one can predict pretty well the possible outcome of a measurable subset of samples. For example, the alternative

hypothesis can be separated from the null hypothesis if the distribution of the test statistic is very different from the null hypothesis. The alternative hypothesis corresponding to the null hypothesis causes the test statistic to be at the edges of a probability distributions, which makes them hard to predict.

### 2.2.3 P-value

P-value is estimated from the test statistic. It is used to find out if the test data is interesting. A certain threshold, called significance level (usually 0.05) is taken for the p-value. If the p-value is within significance level, then null hypothesis is rejected. If the p-value is not within significance level, then null hypothesis is kept and the data falls into distribution as defined by null hypothesis. A variety of algorithms have been thought of for p-value, depending on the type of data, goals of the tests and even the amount of data. Statistics need to make the right choice between the alorithms by taking all these assumptions into an account.

The p-value can be either one sided or two sided. One-sided p-value is the probability that you get a certain value  $T$  that is larger or equal to the computed test statistic from the observed data in the distribution that the null-hypothesis defines. This can be seen visually in Figure 2.

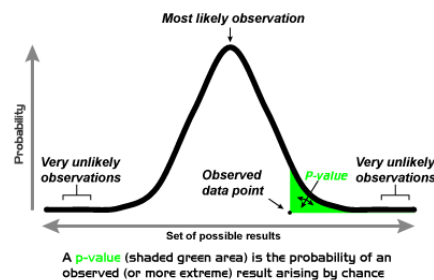


Figure 2: One sided p-value representation (Author Repapettilto, Wikipedia Foundation Inc)

Two-sided p-value shows the probability that you get a certain value  $T$  that is either larger or smaller than the computed test statistic from the observed data in the distribution that the null-hypothesis defines. This can be seen visually in the Figure 3.

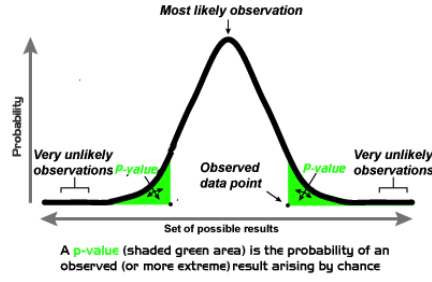


Figure 3: Two sided p-value representation (Modified by author, source from Wikipedia Foundation inc, Repapettilto)

## 2.3 Wilcoxon test

The Wilcoxon signed-rank test is a statistical hypothesis test that is used when comparing two related samples, matched samples, or repeated measurements on a single sample. This means that with Wilcoxon test you can look at some measurable feature, e.g. blood pressure. If two or more experiments in different conditions have been made on it, e.g. gene experiments and their confirmation results, Wilcoxon test can be used to see if the changes of the measurement levels are relevant. The null hypothesis  $H_0$  means that the difference between pairs is zero. The alternative  $H_1$  means it is not.

It was popularized by Sidney Siegel in his book "Nonparametric statistics - for the behavioral sciences". Sidney used the symbol  $T$ , to denote test statistic. Because of this, the test is sometimes referred to as the Wilcoxon T test and the test statistic is reported as a value of  $T$ .

## 2.4 Wilcoxon test Assumptions

1. Two measurements can always be compare and thus ordered. For example, measurement is a real value or ordered categorical values
2. Measurements can be naturally paired into case-control or before-after experiments. For example, measure patients white blood cell levels before and after treatment.
3. All measurement pairs are independent from other pairs. For example, measurement of one individual does not influence the other. In medicine patients in the study are sampled randomly. There is no evident procedure in selecting patients or no planned action to organize a study so one could get the "desired" result.

## 2.5 How to compute

Let  $N$  be the sample size, the number of pairs. Then we can use the following variable to compute the  $W$  test. For  $i = 1, \dots, N$ , let  $x_{1,i}$  and  $x_{2,i}$  denote the measurements. Let  $W$  be the wilcoxon test statistic. Let  $z - score$  be the wilcoxon test standard score. Let  $N_0$  be the number of pairs from which onward we can approximate the p-value. This value is currently undefined and is the focus of this paper. It will be discussed in the future chapters.

The computation of the statistic  $W$  is organised into five steps. The last step diverges into two possible paths and the choice of the path that is chosen depends whether  $N \geq N_0$ .

1. For  $i = 1, \dots, N$ , calculate  $|x_{2,i} - x_{1,i}|$  and  $\text{sign}(x_{2,i} - x_{1,i})$ , where  $\text{sign}(x)$  is the sign function.
2. Order the  $N$  pairs from smallest absolute difference to largest absolute difference,  $|x_{2,i} - x_{1,i}|$ .
3. Rank the pairs, starting with the smallest as 1. Ties receive a rank equal to the average of the ranks they span. Let  $R_i$  denote the rank.
4. Calculate the test statistic  $W$ , the absolute value of the sum of the signed ranks.

$$W = \left| \sum_{i=1}^N [\text{sign}(x_{2,i} - x_{1,i}) * R_i] \right| \quad (1)$$

5. As  $N$  increases, the sampling distribution of  $W$  converges to a gaussian distribution. Where  $N_0$  depends on how accurate you want your results to be.

- (a) For  $N \geq N_0$ , a  $z - score$  can be calculated as

$$z = \frac{W - 0.5}{\sigma_w}, \sigma_w = \sqrt{\frac{N(N+1)(2N+1)}{6}} \quad (2)$$

If  $Z > Z_{critical}$  then reject  $H_0$ , where  $Z_{critical}$  is calculated with Gaussian distribution.

- (b) For  $N < N_0$ ,  $W$  is compared to a p-value that can be calculated from enumeration of all possible combinations of  $W$  given  $N$ .

If  $W \geq P_{critical}$ ,  $N$  then reject  $H_0$

## 2.6 Example

Given pairs  $(6, 8), ((2, -3), (-3, 3), (1, 3))$ .

1. Calculate absolute values and signs:  
Absolute values  $(2, -5, 6, 2)$   
Signs  $(1, -1, 1, 1)$
2. Order the pairs. Ties receive the rank equal to average of the ranks they span  
 $(2 \Rightarrow 1.5), (5 \Rightarrow 3), (6 \Rightarrow 4), (2 \Rightarrow 1.5)$
3. Calculate the test statistic  $W$   
 $W = 1 * 1.5 + 1 * 1.5 + (-1 * 3) + 1 * 4 = 4.0$
4. Since  $N_r$  is very small, use a table to look up the  $p$  value. The calculation of  $p$  table is in a later chapter.  
 $P(5, 4) = 0.3125$   
Since  $0.3125 > 0.05$ , reject  $H_1$

The calculations can also be seen on Table 1 and Table 2.

i	$x_{2,i}$	$x_{1,i}$	sign	abs
1	6	8	1	2
2	2	-3	-1	5
3	-3	3	1	6
4	1	3	11	2

Table 1: Initial data

i	$x_{2,i}$	$x_{1,i}$	sign	abs	$R_i$	sign * $R_i$
1	6	8	1	2	1.5	1.5
4	1	3	11	2	1.5	1.5
2	2	-3	-1	5	3	-3
3	-3	3	1	6	4	4

Table 2: After ordering by absolute difference

## 2.7 P-value assignment

Let us consider Wilcoxon test if we have 6 measurements. Then there are three differences as shown in Figure 4. As the test considers only the sign and order of the measurements, then only 8 configurations are possible as shown in Figure 5 and Table 3.

The formula to calculate that result was  $2^n$ . Due to the symmetry of the possible configurations under the distribution of the null hypothesis, all 8 configurations can be proven equally probable.

Now if the  $W$  value in our data is  $+2$ , then there are 3 configurations that have larger or equal  $W$  value. This can be seen from the Table 3. Then the one-sided

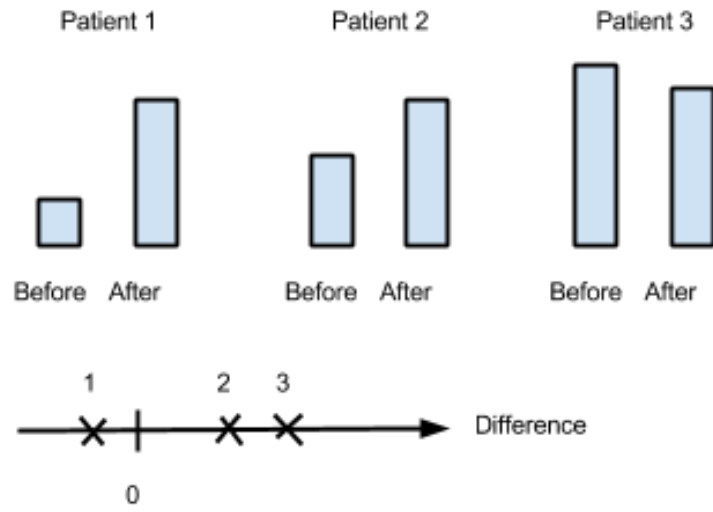


Figure 4: Patient example

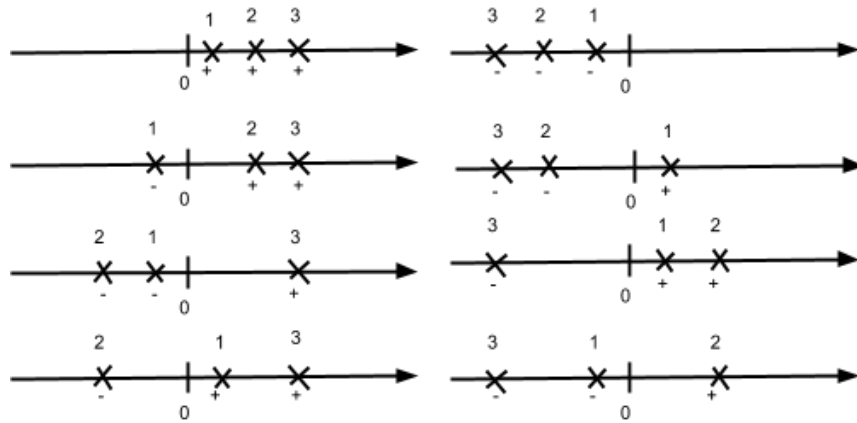


Figure 5: Two sided p-value representation

p-value is  $\frac{3}{8}$ . The two-sided p-value is  $\frac{3}{4}$ . The same procedure is needed to assign p-value.

$W = 1 + 2 + 3 = 6$	$W = -1 - 2 - 3 = -6$
$W = -1 + 2 + 3 = 4$	$W = 1 - 2 - 3 = -4$
$W = -1 - 2 + 3 = 0$	$W = 1 + 2 - 3 = 0$
$W = 1 - 2 + 3 = 2$	$W = -1 + 2 - 3 = -2$

Table 3: All eight configurations. Values of the test statistic corresponding to the configuration.

So in the end, since Wilcoxon test neglects the actual value of the tests, then to estimate the p-value, we only need to know the  $W$  value and number of tests. No matter the actual value of the test cases,  $W$  statistic will always create the same configurations, given a certain  $N$ .

### 2.7.1 V and P example

Given  $N$  is 2, then ranks are 0, 1, 2, one can combine them into  $0+1+2$  or  $0+1-2$  or  $0-1+2$  or  $0-1-2$ . This gives us one way to get 3, one way to get 1, one way to get  $-1$  and one way to get  $-3$

It can now be calculated that the probability that our there is a 0.25% that our value is 3 or lower and 0.5% that our value is 1 or lower.

This distribution is the most accurate table you can compare wilcoxon test results on - it shows the exact probability that your tests falls into a certain gaussian distribution range. However, the table is expensive to calculate.

### 2.7.2 Exact computation of P-value

The problem of assigning a  $p$ -value can be reduced to the following problem. Given a  $N$  ranks, we assign  $+$  or  $-$  sign to each rank randomly. Since the sign distribution is random, then either sign has equal probability to be assigned. For fixed number of ranks  $N$ , we can consider all  $2^n$  possible sign assignments and find out in how many ways we can assign the ranks to get a certain sum  $k$ . If we do this for every possible  $k$ , then a  $V$ -array is formed. It signifies all the possible random signed rank combination sums in an  $N$  sized ranked array, where  $k$  shows how many times a certain sum was achieved. If we calculate  $V$ -array for every possible number of ranks up to the fixed  $N$ , then a  $V$ -table is formed. We can use the  $V$ -table to look up the how many times a certain sum  $k$  was achieved, given  $N$  ranks. Let  $V_{N,k}$  be the number of times a certain sums was achieved.

From this table, the probabaility that  $W \geq k$  can be calculated. This value will show that given  $N$  ranks with random signs, how probable it is that a sum  $k$  is the result. If this process of calculating values is done repeatedly for every possible  $N$  and  $k$ , then a  $P$ -table forms. The  $P$ -table shows the probability of

every possible sum  $k$  appearing given  $N$  ranks with random signs. Let  $P_{N,k}$  be the function to get the desired  $P$  value from that table.

If the  $P$ -table gets large enough, it will start taking the shape of Gaussian distribution. Because of that, given enough ranks, a gaussian distribution can be used to approximate the  $P$ -table.

### 2.7.3 The V and P table formulas

Given  $N$  which shows us the number of ranks starting from 0 and  $k$  which shows us the sum that we are interested in. Recursively apply this algorithm until you reach to the  $N = 1$ .

$$V_{N+1,k} = V_{N,k-1} + V_{N,k+N+1} \quad (3)$$

There are some additional conditions to the formula:

$$V_{N,k} = 0 \quad \text{if} \quad k < -N\frac{N+1}{2} \quad \text{or} \quad k > N\frac{N+1}{2} \quad (4)$$

$$V_{N,k} = 1 \quad \text{if} \quad k = -N\frac{N+1}{2} \quad \text{or} \quad k = N\frac{N+1}{2} \quad (5)$$

With these formulas, the probability that  $W$  is a certain value can be calculated. Given  $V$  table,  $N$  and  $k$ , we can calculate the  $P$  by

$$P(T) = \frac{1}{2^n} \left| \sum_{K=T}^{\infty} W_{N,K} \right| \quad (6)$$

To get the P table, we go through all  $N$  and  $K$  values that we are interested in.

### 2.7.4 The Gaussian approximation of P table

As mentioned above, if the number of ranks  $N$  gets large enough, then Gaussian distribution can be used. To use that, the cumulative density function  $F(x)$  is usually used. A simple Gaussian function  $F(x, 1)$ , where  $x$  is the statistic  $z$  - value, eg standard score and 1 is the sigma, will give approximately the same results as the accurate  $P$  table.

### 2.7.5 Code sample

A python code to calculate the V table is as follows:



```

def V(n, k, vTable):
    if((n == 1 and (k == 1 or k == -1)) or (n == 0 and k == 0)):
        return 1
    if((k < -(n * (n + 1) / 2)) or (k > (n * (n + 1) / 2))):
        return 0
    n = n - 1

    kIndex = k + kSize

    leftValue = 0 if kIndex - n - 1 < 0 or \
        kIndex - n - 1 >= len(vTable[n]) else vTable[n][kIndex - n - 1]
    rightValue = 0 if kIndex + n + 1 < 0 or \
        kIndex + n + 1 >= len(vTable[n]) else vTable[n][kIndex + n + 1]

    return leftValue + rightValue

def calculateVValues():
    vTable = []
    for n in range(nSize):
        tableRow = []
        for T0 in range(kSize * 2 + 1):
            tableRow.append(V(n, T0 - kSize, vTable))
        vTable.append(tableRow)
    return vTable

```

The python code to calculate the P table is as follows:

```

def calculatePValues(vTable):
    pTable = []
    for n in range(nSize):
        tableRow = []
        sumValue = 0
        powerValue = 2**n
        for T0 in range(kSize * 2, kSize-1, -1):
            sumValue += vTable[n][T0]
            p = sumValue / powerValue
            tableRow.insert(0, p)
        pTable.append(tableRow)
    return pTable

```

You can find a python implementation of a code sample in <https://github.com/stenver/wilcoxon-test/blob/master/WilcoxonVTable/pAccurateTable.py>

### 2.7.6 Bonferroni correction

If you make multiple hypothesis on a test, then you increase the risk in which you reject null hypotheses when its actually true. The Bonferroni test helps to counteract this in a simple and naive way. For each hypothesis you make on a test, you should use a significance level  $M$  times lower than before. This ensures that the the increased risk in which we can reject null hypothesis will not raise, no matter the number of hypothesis on a test. So for example, if you make  $M$  hypothesis and want want a significance level  $\alpha$ , then you should run each test at a significance level of  $\frac{\alpha}{M}$ .

If you want to use Bonferron correction with Wilcoxon signed-rank test, then you need to keep in mind that the approximated  $P$  value significance level needs to be much more accurate, since the bonferron test divides it by the amount of features tested. This means higher the  $N_0$  value, the more features you add in the test.

## 3 Approximating the p-value

### 3.1 Motivation

The current implementations of Wilcoxon signed-rank test usually assume that  $p$ -value distribution is sufficiently close to Gaussian distribution, so that we can use  $Z$  to calculate the  $P$  value. They draw this conclusion from the assumption that the number of samples  $N$  is sufficiently large to use  $Z$ . When  $N = 1$ , then approximation is quite unaccurate. As number of samples increases, then the approximation becomes more and more accurate on absolute scale. When the  $N \rightarrow \infty$ , then the  $P_{approx}$  approximate region will almost equal the  $P$  accurate region, which also means that the region will be almost perfect Gaussian Distribution. This means that when  $N \rightarrow \infty$ , then there is no reason not to use Gaussian Distribution, as it will be almost completely accurate, except at the very edges. This allows them to use Gaussian Distributions to approximate the  $P$  value and the accurate  $P$  table calculation is often left unoptimized.

In the BIIT research group, and many others, the  $N$  is usually not sufficiently large. This means that Gaussian distribution cannot be used. Also, since BIIT uses multiple hypothesis testing, the approximation must be good even at the very edges of distribution. As mentioned above, the Bonferroni correction will lower the significance level significantly as  $N$  rises. Furthermore, it is not known where exactly the  $N$  limit is, so an arbitrary number is used for that purpose. The textbooks [RLow11] currently suggest a sufficient  $N$  is 10. In this work we studied the question whether this assumption is justified in detail.

The approximation is computationally cheap and the preferred method to be used over accurate  $P$  table, which is computationally expensive. Because of this, one focus of the paper is finding out how good the approximation of Gaussian Distribution is and finding out the exact  $N$  where we can use approximation instead on accurate  $P$  table while maintaining a high enough accuracy. Let  $N_0$  be high enough  $N$ , which allows us to use approximation accurately.

The bottlenecks of the other implementations of the Wilcoxon test is, when  $N$  is low, but the number of tests is high. In this case they calculate the  $P$  and  $V$  tables every time the test is run, thus if you run thousands of tests in a row, a lot of complex recomputation is done. This eventually becomes a performance issue. To speed this up - this project will calculate the entire  $P_{N,k}$  for every value that is possible where  $N < N_0$ . Then the  $P_{N,k}$  will be hardcoded inside the program for quick lookup of the value.

### 3.2 When can we approximate

The relative and absolute error behave very differently. The absolute error between accurate table values  $P$  and approximate table values  $P_{approx}$  first falls and then rises until at the centre of the distribution, but gets stable and almost constant from the middle to the edge of the distribution. It can be seen on the Figure 6. The Relative error gets smaller, the higher the number of samples  $N$  is. However, the relative error looks like a slide. It is stable at the centre of the distribution, but starts rising sharply at the middle. At the edge of the distribution, the relative error is nearly 100%. It can be seen at the Figure 7

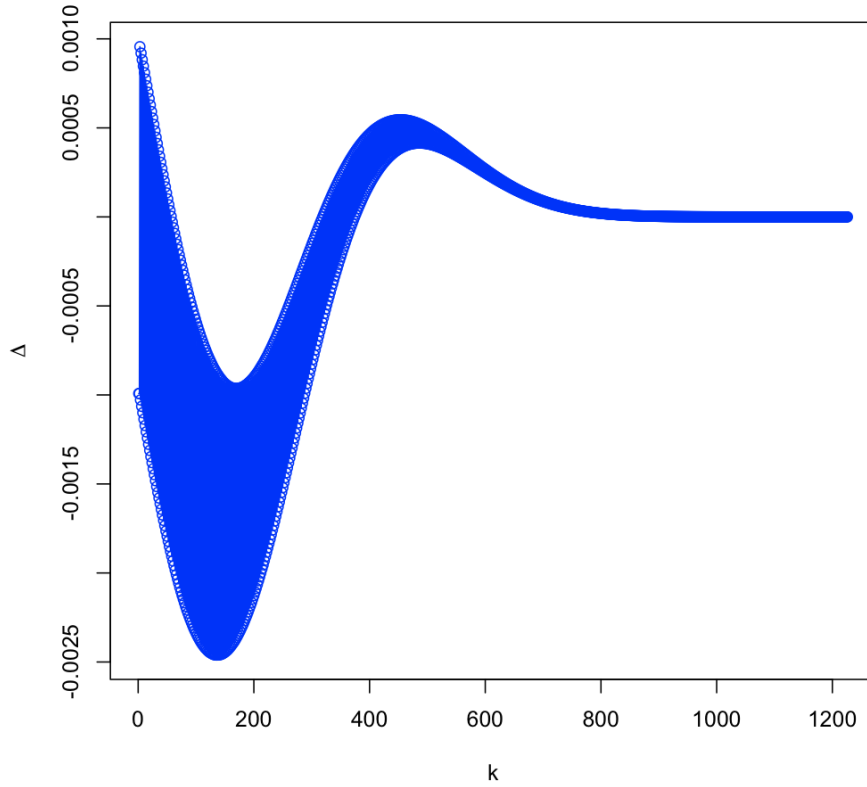


Figure 6: Absolute error between accurate  $P$  and approximate  $P_{approx}$  values.

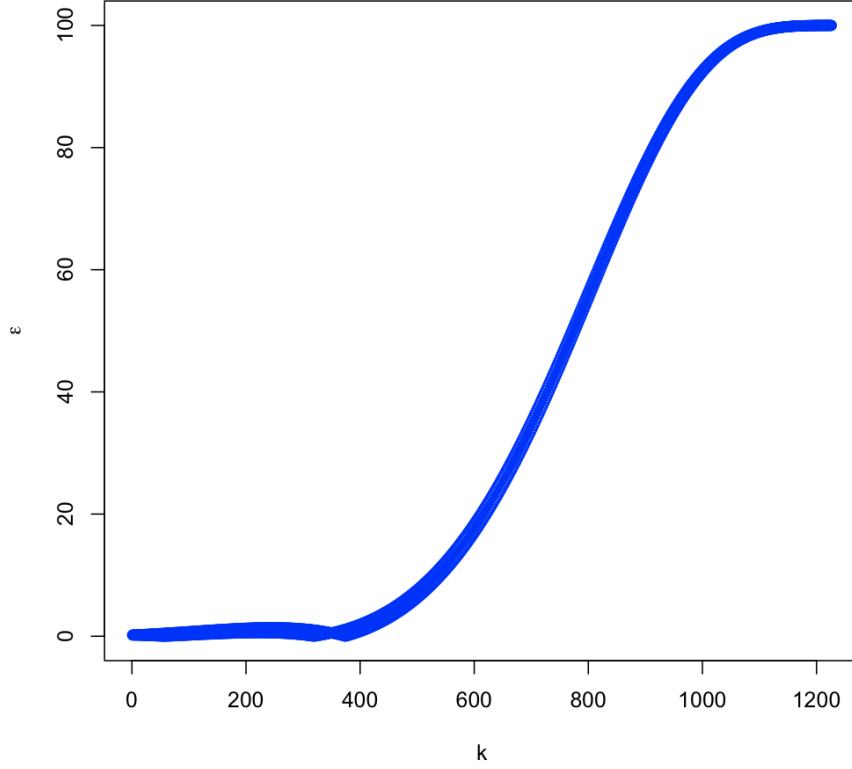


Figure 7: Relative error between accurate  $P$  and approximate  $P_{approx}$  values.

The absolute error is calculated simply as follows:

$$\Delta = P - P_{approx} \quad (7)$$

The relative error is calculated as follows:

$$\epsilon = \frac{\Delta}{P_{approx}} * 100\% \quad (8)$$

It is important to have low relative error on the region - i.e. good region, because this ensures that the test results stay accurate. If the relative error is high, then this will also increase the possibility of a false positive, that is, rejecting  $H_0$  when we actually should not.

First we wanted to confirm that as the number of samples grows, accurate table values  $P$  and approximate table values  $P_{approx}$  will get become more similar. To do this, we used the following equation for each  $N$ :

$$k_1 = \max_k \{k : p(k) > 0\} \quad (9)$$

And compared them to the results of:

$$k_0 = \max_k \{k : \epsilon_k \leq 5\%\} \quad (10)$$

The results of the table can be seen on Figure 8.

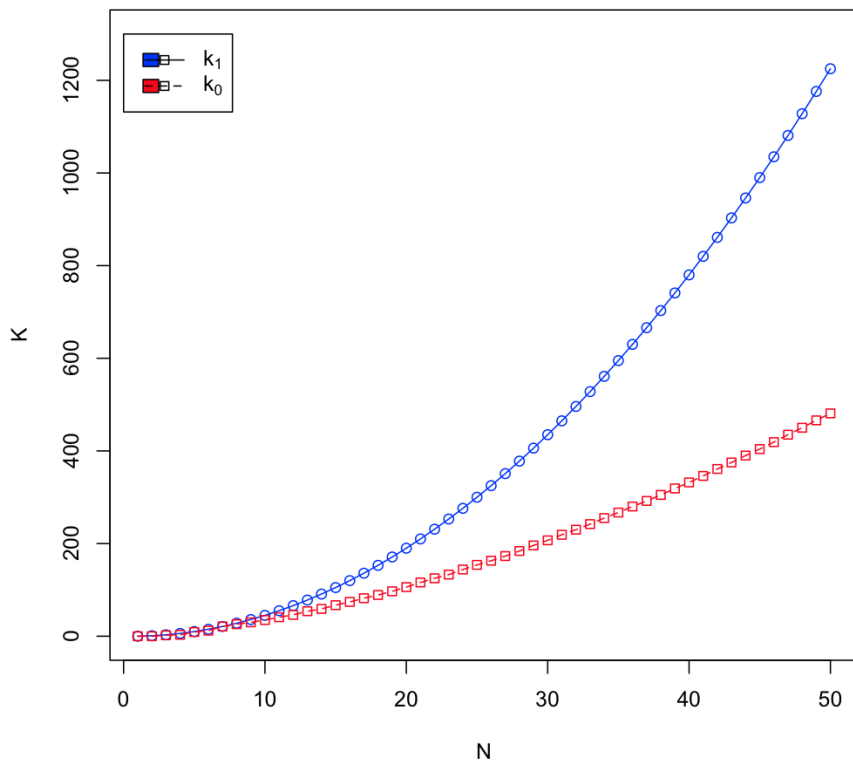


Figure 8: Region between low relative error and actual distribution edge increasing as sample size increases.

Much to our surprise, as  $N$  grew, the gap between approximate and accurate values grew bigger. This meant that as  $N$  grows, the relative approximate error of Gaussian distribution will get more inaccurate toward the tails of the distribution. Previously we thought that with the growth of  $N$ , Gaussian would surely get more and more accurate. From Figure 9 we can't really make a difference between accurate  $P$  and approximate  $P_{approx}$  values, but if we put the values to logarithmic scale, the difference is easier to spot. From Figure 10 we can see that approximated  $P_{approx}$  values are a little bit bigger toward the edge of distribution - as can be seen when  $K$  grows.

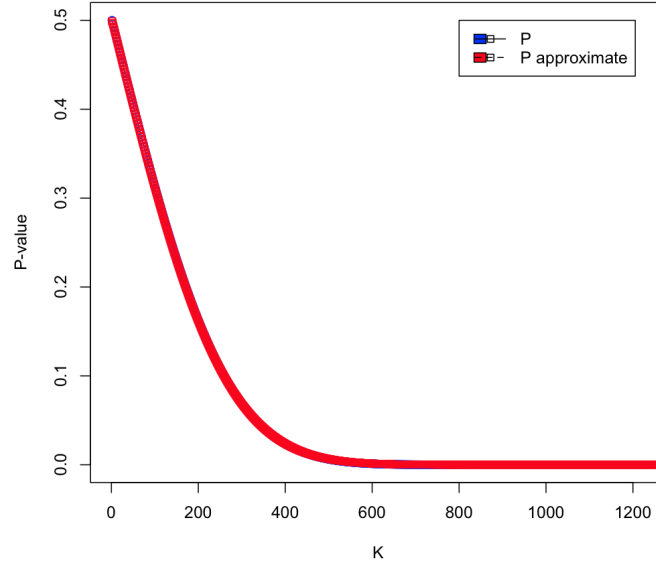


Figure 9: Accurate  $P$  and apporximate  $P_{approx}$  values as  $k$  increases.

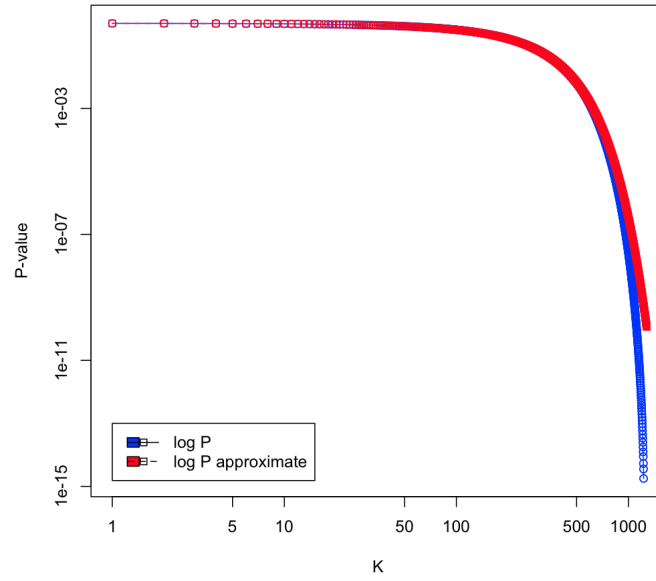


Figure 10: Logarithmic accurate  $P$  and apporximate  $P_{approx}$  values as logarithmic  $K$  increases.

To further investigate this finding, we decided to find out the minimum  $P$  value that you can get for each  $N$  while maintaining a certain relative error  $\epsilon$

threshold. The thresholds chosen were  $\epsilon = 5\%$ ,  $\epsilon = 10\%$ ,  $\epsilon = 20\%$ ,  $\epsilon = 50\%$ . From Figure 11 we can see that as the  $N$  grows, the graph becomes stable, meaning that the distribution becomes a Gaussian distribution at around  $N > 10$  and  $N < 25$ . When  $N < 10$ , then the approximation is so random that it cannot be used reliably. We can see that the minimum  $P$  value under error threshold does get smaller, as  $N$  increases, so this further proves that Gaussian distribution does get more accurate as  $N$  grows. This is because as  $N$  increases, you can take a  $P$  value closer to the tail of the distribution and still get an accurate value.

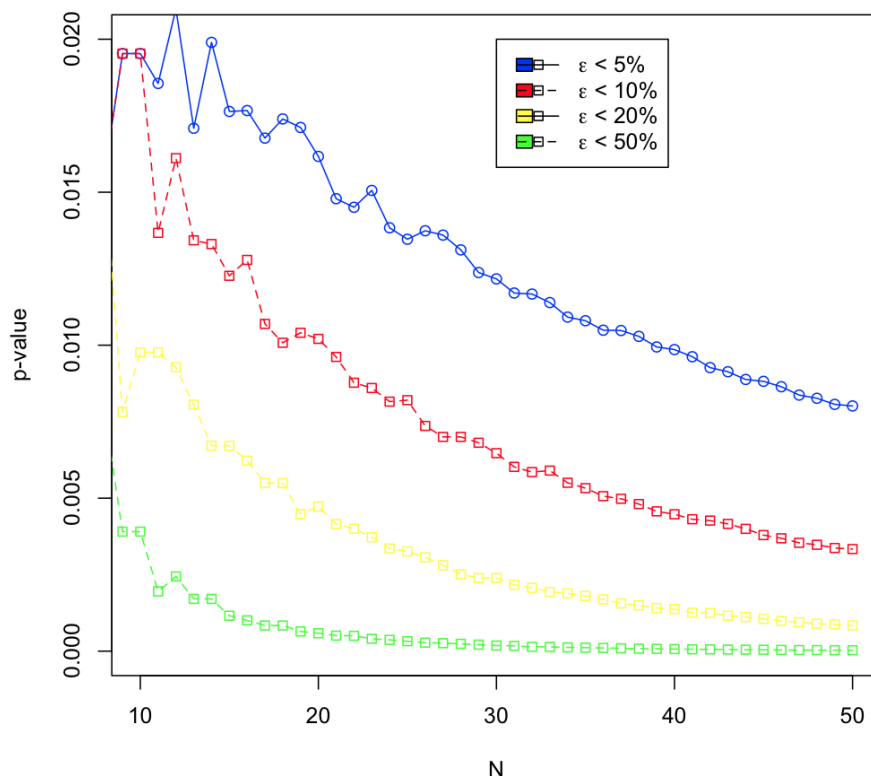


Figure 11: Relative error  $\epsilon$  as sample size  $N$  increases

The Figure 12 shows us the relative error as  $P$  grows when  $N = 20$ . The  $P$  values are logged on this graph. There are two noteworthy things here, however. First, around  $P = 0.05$ , the relative error actually get a little small tip toward accuracy. This is evidently true, since the the accurate  $P$  and approximate  $P_{approx}$  values area is equal.

$$\left| \sum_{k=1}^{\infty} P(k) \right| = \left| \sum_{k=1}^{\infty} P_{approx}(k) \right| \quad (11)$$



Since the approximation  $P_{approx}$  value is higher at the edge of the distribution, then obviously it is smaller at the centre of distribution. Second, it can be seen that there are two accuracy paths that the error takes, depending on whether in  $P_{approxN,K}$  the  $K$  is even or odd number. This is a computational artefact caused by the recurrence of the  $V_{N,K}$  values and we will not make any conclusions of that. To understand this better, one must familiarize himself with the algorithm defined in the Chapter 2.7.3.

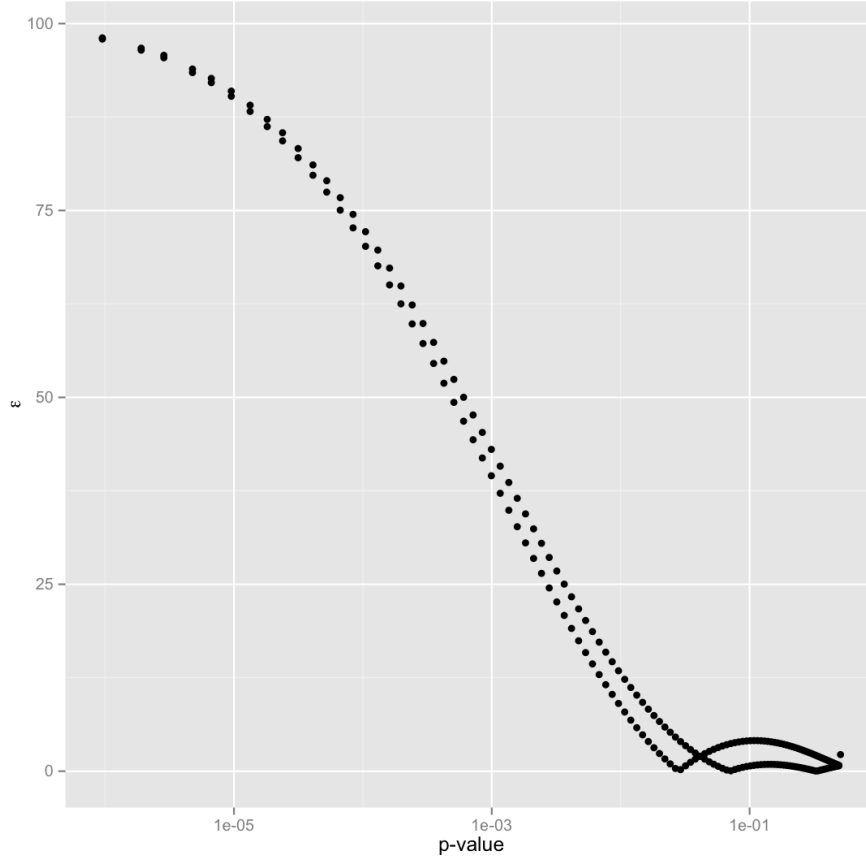


Figure 12: Relative error  $\epsilon$  decreasing as p-value grows, when  $N = 20$

### 3.3 Investigating the $P_{N,k}$ and $P_{approxN,K}$ similarities

To further prove that we can start using Gaussian Distribution at a certain  $N$ , we needed to prove that  $P_{approx}$  will get more accurate as the  $N$  increases toward the tail of the distribution as well. To achieve this, we found out the largest relative

error  $\epsilon$  between  $P$  and  $P_{approx}$  for each  $N$  and for each p-value where:

$$p - value = \{(\frac{1}{10^0}; \frac{1}{10^1}), (\frac{1}{10^1}; \frac{1}{10^2}), \dots, (\frac{1}{10^6}; \frac{1}{10^7})\} \quad (12)$$

As can be seen from Figure 13 and Figure 14, the relative error increases fast. When p-value =  $(10^{-5}, 10^{-6})$  the relative error is massive. When  $N = 200$ , then the error is still around 40%. The bigger the  $N$  gets, the bigger the error at the tail. However, it can also be seen that the relative error decreases steadily as  $N$  grows for each of the p-value ranges chosen. This means that even though the error of the distribution tail edge does increase as  $N$  grows, the approximate  $P_{approx}$  and accurate  $P$  values do get more similar as the  $N$  grows, i.e. the theory under question does not lie.

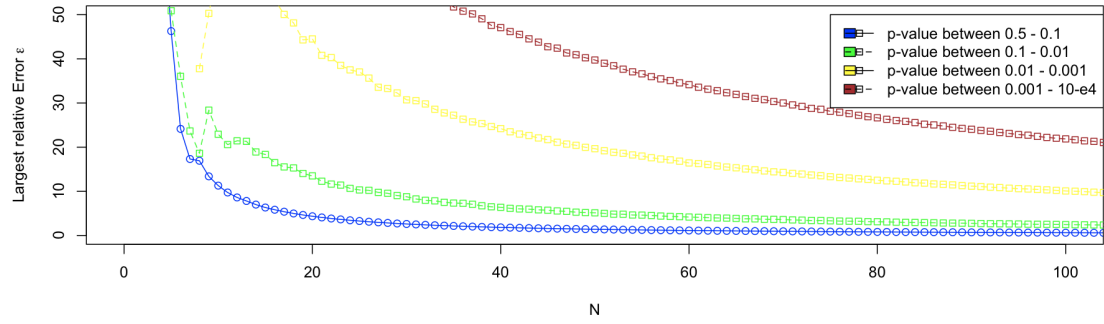


Figure 13: Showing the maximum relative error  $\epsilon$  decreasing in probability ranges, as sample size  $N$  increases

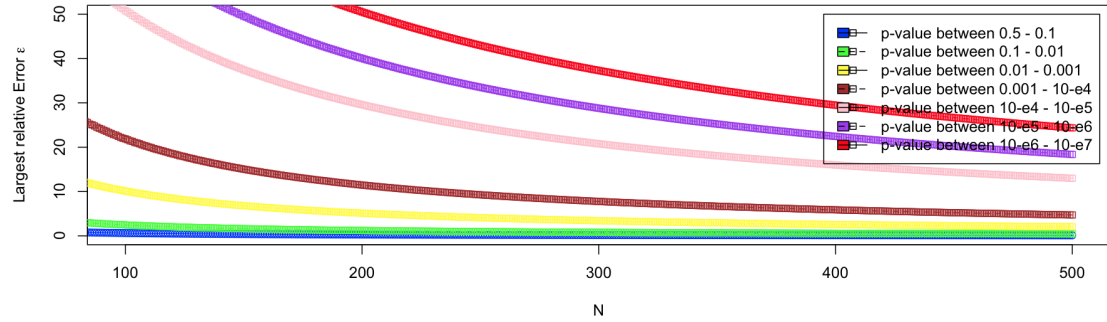


Figure 14: Showing the maximum relative error  $\epsilon$  decreasing in probability ranges, as sample size  $N$  increases

To investigate the relation between accurate  $P$  and approximate  $P_{approx}$  values, we wanted to see the difference between all p-values values when  $N = 50$  and when  $N = 25$ . As can be seen from Figure 15 or Figure 16, the  $P_{approx}$  is constantly a little bit larger. When  $N = 25$ , the graph is a little less smooth than when  $N = 50$ . This follows the same story as before - the bigger the  $N$ , the more similar the p-values.

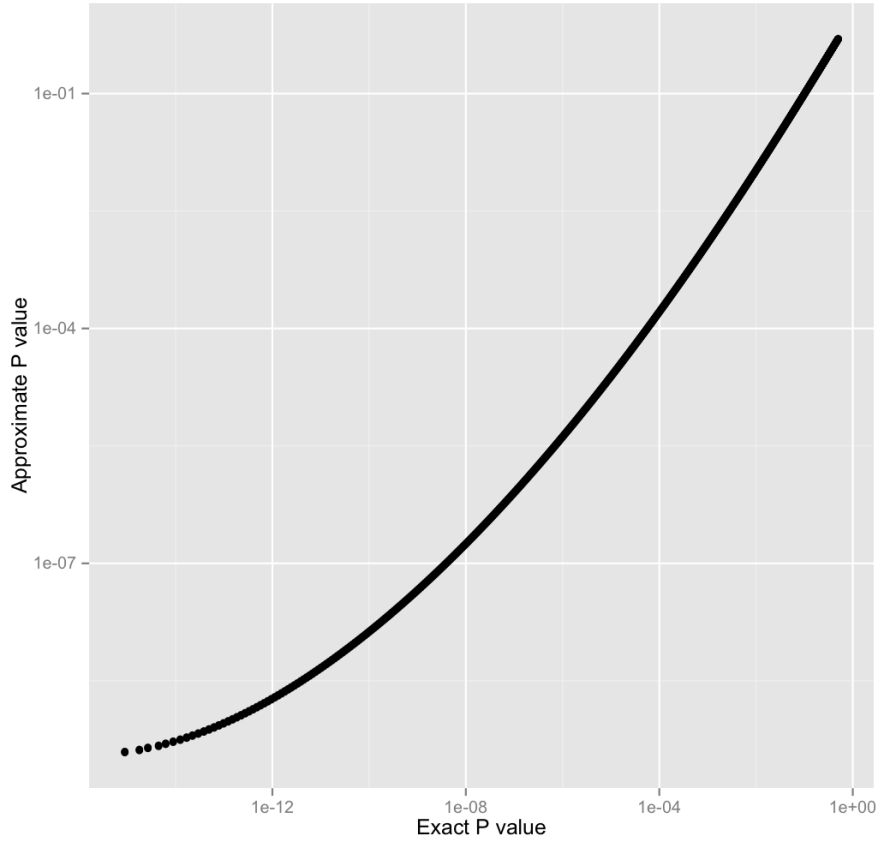


Figure 15: Relative value of all accurate  $P$  and approximated  $P_{approx}$  values when sample size  $N$  is 50

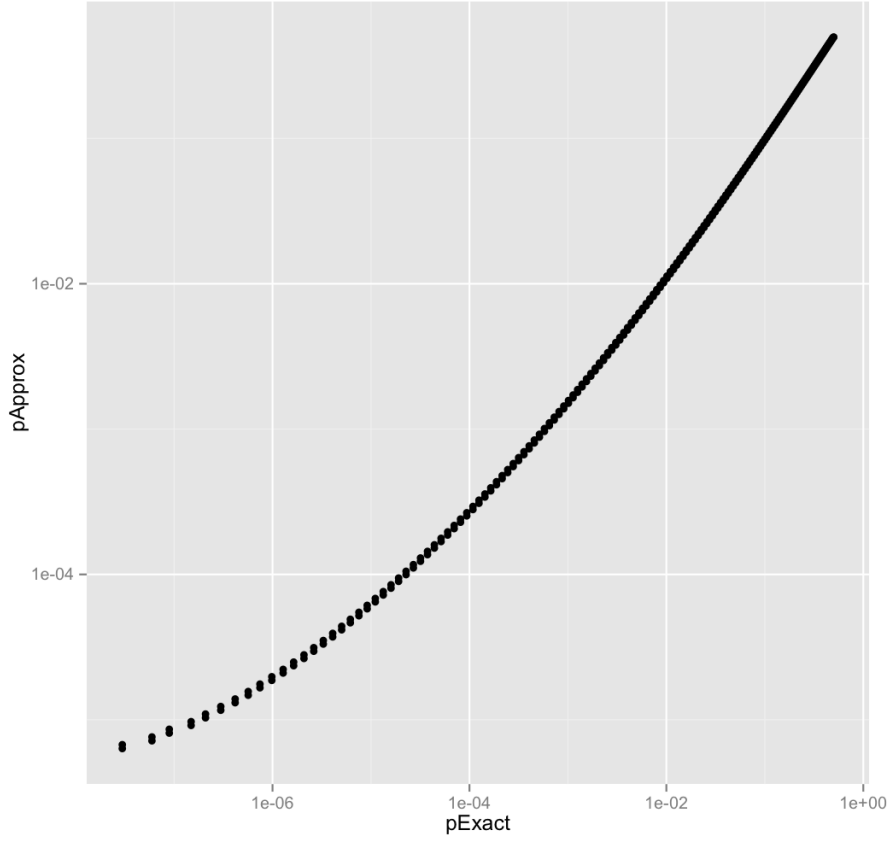


Figure 16: Relative value of all accurate  $P$  and approximated  $P_{approx}$  values when sample size  $N$  is 25

### 3.4 Finding out when can we approximate p-value

Judging by all the data that has been gathered, it be can clearly see that as  $N$  grows, accurate  $P$  and approximate  $P_{approx}$  values grow apart at the edges of distribution. Nevertheless, overall the distribution gets more and more similar. The most helpful graphs to help us determine the necessary number of samples  $N_0$  we need to approximate p-value is Figure 14. Table 4 illustrates the number of samples  $N$  needed for a certain relative accuracy. So the number of samples  $N_0$  we need to approximate depends on the number of measurements and the required accuracy.

Table 4: Showing the minimal required sample size for a certain relative error and number of measurements

Measurements	Error	Required N
10	5%	25
100	5%	250
1000	5%	500
10000	5%	> 500
100000	5%	> 500
10	10%	80
100	10%	200
1000	10%	> 500
10000	10%	> 500
100000	10%	> 500
10	20%	50
100	20%	100
1000	20%	280
10000	20%	500
100000	20%	> 500
10	50%	25
100	50%	50
1000	50%	100
10000	50%	150
100000	50%	200

At the beginning of this chapter, we raised the goal of optimizing our implementation on Wilcoxon test by hardcoding the accurate accurate  $P$  value table  $P_{N,K}$ , so the program would not need to recompute it all the time. When given  $N = 80$ , then the table  $P_{N,K}$  size was 2.5 MB. As pointed out in the beginning, when the  $N$  grows, the table  $P_{N,k}$  size grows with it with the speed  $O(n^2)$ . At  $N = 500$ , the file size was over 10 GB. This raised the need to optimize the table  $P_{N,k}$  further.

## 4 Further optimizations

To further optimize the library, it was needed to approximate the accurate  $P$  table. Two algorithms were chosen for that matter dose-response curve and linear approximation.

#### 4.0.1 dose-response curve

When using the dose-reponse curve algorithm, the data range as shown in Figure 17 changes sharply toward the edge of the range. Further tests also showed that the algorithm was completely inaccurate at the end of the far edge, so unfortunately this algorithm could not be used.

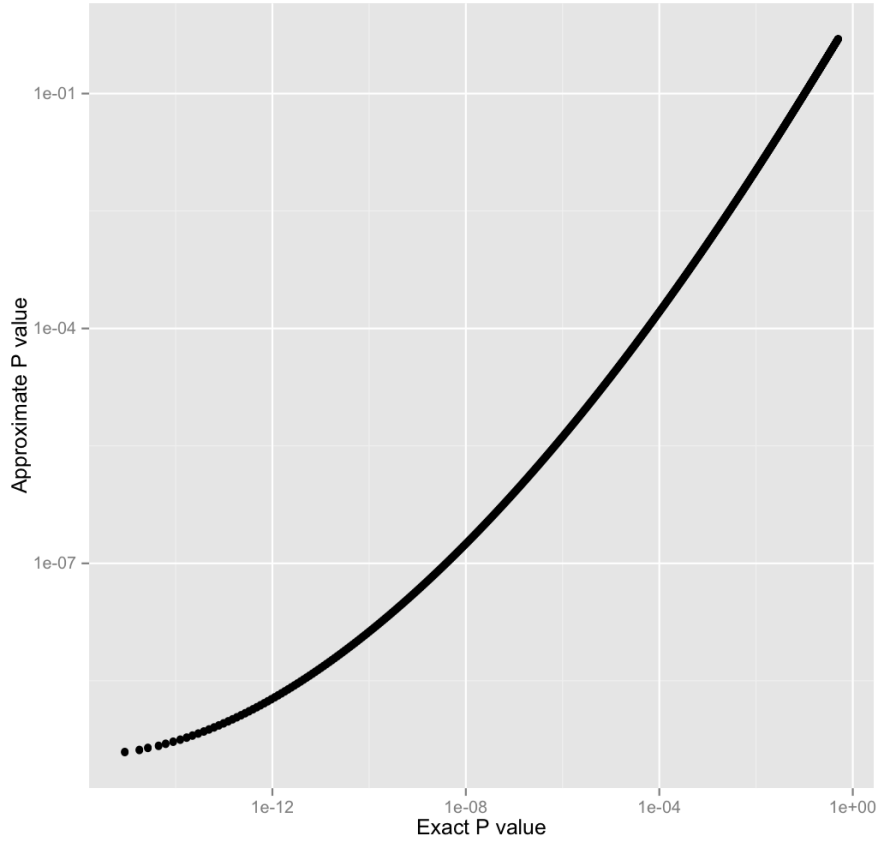


Figure 17: P range with DRC algorithm when  $N=500$

#### 4.1 Linear approximation

A more straight forward algorithm to use was linear approximation. The tests showed that using relative  $P$  values between the  $P$  tables gave less approximated data points, then using straight forward accurate  $P$  table.

#### 4.1.1 Relative values between P tables

The algorithm to get a relative value for every position between the P tables was pretty straight forward. In Python it looked like this:

```
def getRelativeValuesInRow(n, pTable, gaussianPTable):
    pRelativeValues = []
    for j in range(len(pTable[n])):
        relativeValue = pTable[n][j] / gaussianPTable[n][j]
        pRelativeValues.append(relativeValue)
        if(pTable[n][j] == 0):
            break
    return pRelativeValues

def calculateRelativeTable(pTable, gaussianPTable):
    relativeTable = []
    for n in range(len(pTable)):
        relativeTable = getRelativeValuesInRow(n, pTable, gaussianPTable)
    return relativeTable
```

This returned a new table which had the same number of elements as the accurate P table or Guassian P rable, but for every position there was a relative p-value.

#### 4.1.2 Linear approximation algorithm

The linear approximation algorithm works by starting from data point in the data range and then one by one moving to the next data point, until the points between them cannot be approximated with enough accuracy. Then the last accurate point is saved, picked as the next starting point and the process is repeated until end of range has been reached. This guarantees approximated data range within error margin.

The to approximate the data points in python look like this

```

def linearInterpolate(approximatePoint,
                      endPoint,
                      startPoint,
                      endValue,
                      startValue):
    return startValue +
        ((endValue - startValue) * (approximatePoint - startPoint) /
         (endPoint - startPoint))

def canApproximate(actualPoints, currentPoint, nextPoint):
    startValue = actualPoints[currentPoint]
    endValue = actualPoints[nextPoint]

    for i in range(1, nextPoint - currentPoint):
        approximatePoint = currentPoint + i
        trueValue = actualPoints[approximatePoint]
        approximateValue = linearInterpolate(approximatePoint,
                                              currentPoint,
                                              nextPoint,
                                              startValue,
                                              endValue)

        if approximateValue == 0:
            return True

        relativeAccuaracy = trueValue / approximateValue

        if relativeAccuaracy < 0.9 or relativeAccuaracy > 1.1:
            return False
    return True

def getNextApproximatedPoint(actualPoints, currentPoint):
    nextPoint = currentPoint + 1

    while(canApproximate(actualPoints, currentPoint, nextPoint)):
        if(nextPoint == len(actualPoints) - 1):
            break
        nextPoint += 1
    return nextPoint

def calculateApproximateRow(actualPoints):

```



```

approximatedPoints = []
currentPoint = 0
approximatedPoints.append([currentPoint, actualPoints[currentPoint]])

while(currentPoint != len(actualPoints) - 1):
    nextPoint = getNextApproximatedPoint(actualPoints, currentPoint)
    approximatedPoints.append([nextPoint, actualPoints[nextPoint]])
    currentPoint = nextPoint
return approximatedPoints

```

### 4.1.3 Linear interpolation

To reverse the algorithm, a very simple formula - linear interpolation is used. Given the approximate  $x$ , start point  $x_1$ , end point  $x_2$ , start point  $y_1$ , end point  $y_2$ , you can easily approximate the  $y$  of that point.

$$y = y_1 + \frac{(y_2 - y_1) * (x - x_1)}{(x_2 - x_1)}$$

In python the code looks like this:

```

def linearInterpolate(approximatePoint,
                      endPoint,
                      startPoint,
                      endValue,
                      startValue):
    return startValue +
        ((endValue - startValue) * (approximatePoint - startPoint) /
         (endPoint - startPoint))

```

## 4.2 Optimization summary

In the end, thanks to linear approximation, the approximated table, when  $N = 500$  could be created. The error margin for the chosen approximated table was chosen to be 20%, since that was deemed as accurate enough. Its still a lot more accurate then Gaussian P table at the edges of the table.

The approximated table is only 1.6mb large. Smaller then when  $N = 80$  the accurate P table was.

As for the library itself - it could run over 20000 wilcoxon tests with 120 samples in under 0.7 seconds. Compared to the vanilla R implementation which took many seconds to run 20 tests, this is a significant improvement.

### 4.3 Further optimizations

A lot of further optimizations could be done on this subject, for example

1. The algorithm to approximate accurate P table could be improved.
2. The plain text file could be compressed.
3. The shared library algorithm implementations could be improved and expanded.

## 5 The implementaion

Repository that contains everything about our findings can be found in

[https://bitbucket.org/stenver/wilxoni-astaku-test/src/870cc6112d0d?](https://bitbucket.org/stenver/wilxoni-astaku-test/src/870cc6112d0d?at=default)

at=default repository

The repository contains a number of folders.

### 5.1 RcppWilcoxonTest

An interface that connects our implementation of the optimized Wilcoxon algorithm to the R. Note that you must have installed our implementation to use it. The interface must be compiled with separate R commands from the command line. You need Rcpp packages for R to use it. It can be installed in R console by running:

---

```
$install.packages("Rcpp")
```

---

After that, the package can be installed to R by running the command in the command line in the folder:

---

```
$R CMD INSTALL .
```

---

You can now load our library in R by calling the

---

```
>library('RcppWilcoxonTest')
```

---

and you invoke the function by calling

---

```
>RcppWilcoxonTest::WilxTest(dataMatrix, testIndexes, controlIndexes)
```

---

### 5.2 TerminalWilcoxonTest

An interface that connects our implementation of the optimized Wilcoxon algorithm to the R. Note that you must have installed our implementation to use it. The interface can be compiled by going to the folder, compiling and running help for further help.

---

```
$make  
WilcoxonTest help
```

---

Currently only supports NetCDF file as input data.

### 5.3 WilcoxonTestLibrary

Implementation for the optimized wilcoxon test. The library can be installed by going to the folder and running

---

```
make install
```

---

## 5.4 WilcoxonVTable

Python program that can calculate  $V$ ,  $P$  and approximated  $P$  tables, print them, create files of the tables and create a number of graphs on the tables.

## 5.5 Seminar\_paper

The folder that contains this paper and all images attached to it. It also contains R programs to create those images.

## 6 Conclusion

It was found out that Gaussian distribution tail edge grows larger apart from the  $P_{N,k}$  as  $N$  increases while the overall distribution grows more similar. In addition, it was confirmed that if thousands of parallel tests are ran, then using approximation is not reliable and instead an accurate table of p-values must be used. Calculating that table is very expensive and thus it is advised to precalculate the table for the program. However, the research shows that if 10 000 measurements are run under 5% relative error, then the hardcoded  $N$  should be over 1000. Table of that size would take many gigabytes of space and is impractical to use.

To get around that, approximation of the accurate  $P$  table was used. This allowed to significantly reduce the size of the file and thus increase the hardcoded  $N$  that was implemented inside the library.

For the time being, the chosen  $N$  was 500, with 20% error margin on the approximation. This is still significantly more accurate then Gaussian approximated table at the edges of the distribution. It will allow to make 1000 parallel measurement with 5% error margin.

The library itself could run over 20000 parallel tests with 100 samples in under a second. It is a shared C++ library with currently implemented interfaces in Cron R and Terminal.

Overall the research included some surprising results and can be considered a success.

## 7 Eestikeelne pealkiri

BakalaureusetÕÕ(6 EAP)

Eesnimi Perekonnanimi

ResÕijmee

## References

- [FWil45] Wilcoxon, F. Individual Comparisons by Ranking Methods. *Biometrics Bulletin*, Vol. 1. Pages 80-83. 1945.
- [SSieg56] Siegel, S. Nonparametric statistics for the behavioral sciences. Pages 75-83. 1956.
- [RGSH11] Rew, R., Davis, G., Emmerson, S. and Daview, H. The NetCDF Users Guide. Pages 5-7, 17-23. 2011.
- [RLow11] Lowry, R. Concepts & Applications of Inferential Statistics. <http://vassarstats.net/textbook/ch12a.html>. ch.12a. 2011.
- [MTrio01] Triola, M. Elementary statistics (8 ed.). Pages 388. 2001.
- [JsCr13] Ritz, C. and Strebig, J. Dose-reponse curve. <http://cran.r-project.org/web/packages/drc/index.html>. Pages 111. 2013.
- [ATSD14] DebRoy, S. and Trapletti, A. Writing R extensions. <http://cran.r-project.org/doc/manuals/R-exts.html>. ch 5-6. 2014.
- [RFDE14] Eddelbuettel, D. and Francois, R. Seamless R and C++ integration. <http://cran.r-project.org/web/packages/Rcpp/vignettes/Rcpp-introduction.pdf>. Pages 3-15. 2014.
- [RFDE14] Wikipedia Foundation Inc. z-score. [http://en.wikipedia.org/wiki/Standard\\_score](http://en.wikipedia.org/wiki/Standard_score). 16 January 2014
- [RFDE14] Wikipedia Foundation Inc. linear interpolation. [http://en.wikipedia.org/wiki/Linear\\_interpolation](http://en.wikipedia.org/wiki/Linear_interpolation). 10 January 2014.
- [RFDE14] Wikipedia Foundation Inc. Gaussian function. [http://en.wikipedia.org/wiki/Gaussian\\_function](http://en.wikipedia.org/wiki/Gaussian_function). 8 January 2014.
- [RFDE14] Wikipedia Foundation Inc. Test Statistic. [http://en.wikipedia.org/wiki/Test\\_statistic](http://en.wikipedia.org/wiki/Test_statistic). 10 January 2014.

**Non-exclusive licence to reproduce thesis and make thesis public**

I, Stenver Jerkku (date of birth: 10th of October 1990),

1. herewith grant the University of Tartu a free permit (non-exclusive licence) to:
  - 1.1 reproduce, for the purpose of preservation and making available to the public, including for addition to the DSpace digital archives until expiry of the term of validity of the copyright, and
  - 1.2 make available to the public via the web environment of the University of Tartu, including via the DSpace digital archives until expiry of the term of validity of the copyright,

Paralell Wilcoxon Signed-rank tests

supervised by Sven Laur

2. I am aware of the fact that the author retains these rights.
3. I certify that granting the non-exclusive licence does not infringe the intellectual property rights or rights arising from the Personal Data Protection Act.

Tartu, 14.05.2014