

Comparing Robust Wilks' statistics in Multivariate Multiple Linear Regression

Thamer Warda Hussein
 Department of Mathematics
 College of Science
 University of Basrah
 Basra, Iraq
thamer93work@gmail.com
[Orcid.org/0009-0006-0621-6373](https://orcid.org/0009-0006-0621-6373)

Abdullah A. Ameen
 Department of Mathematics
 College of Science
 University of Basrah
 Basra, Iraq
dr_abd64@yahoo.com

DOI: <http://dx.doi.org/10.31642/JoKMC/2018/110206>

Received Feb. 10, 2024. Accepted for publication Jul. 25, 2024

Abstract—In multivariate linear regression, the classical Wilks' statistic is the most used method to test hypotheses, which is extremely responsive to the effect of outliers. Many authors have examined the non-robust test statistic established on normal theories for various cases. In this study, we constructed a robust version of Wilks' statistic relying on a reweighed minimum covariance determinant estimator, that relies on the weight of the observations, with the weights being determined by using the Hampel weight function and Huber weight function. A comparison of the suggested statistics with the regular Wilks' statistic has been discussed. Monte Carlo studies are used to evaluate how well test statistics work with different datasets. So, this study looked at how two different test statistics perform under a normal distribution. The first is the classical Wilks' statistic, and the second is a proposed new one. For both of them, the rate of type I errors and the power of the tests were close to the expected significance levels. In situations where the distribution is contaminated, the proposed statistical method works best. If the data has been corrupted or affected somehow, this approach seems to perform the best out of the options

Keywords— Minimum Covariance Determinant Estimator, Outliers, Robustness, P-Value, Wilk's Statistic.

I. INTRODUCTION

Consider we have a q -variate dependent (predictor) $\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{iq})^T$ and a p -variate independent (response) $\mathbf{y}_i = (y_{i1}, y_{i2}, \dots, y_{ip})$. The multivariate multiple linear regression model is given by:

$$\mathbf{Y} = \mathbf{X}\mathbf{B} + \mathbf{\Xi} \quad (1)$$

where \mathbf{Y} where $\mathbf{Y} = (\mathbf{y}_1, \dots, \mathbf{y}_n)^T$, $\mathbf{X} = (\mathbf{1}_n, (\mathbf{x}_1, \dots, \mathbf{x}_n)^T)$, $\mathbf{1}_n$ is a n -dimensional vector whose all entries are 1, \mathbf{B} is $((q+1) \times p)$ slope matrix and $\mathbf{\Xi}$ is $(p \times n)$ errors matrix. Multivariate regression finds practical use in various fields such as engineering, biology, psychology, finance and many more fields. Recent research studies on the multivariate regression, Friedman and Breiman (1997) [1], McKean and Davis (1993) [2], Ollia and Koivunen (2003)[3]. For testing the null hypothesis H_0 there is no significant relationship between the set of dependent variables \mathbf{Y} and the set of independent variables \mathbf{X} in other words all population regression

coefficients are zero, meaning that H_0 contains that none of the independent variables collectively contribute to explaining the variation in the set of dependent variables. For testing H_0

various statistics have been used, Wilks' statistic Λ is the most widely used which is defined as:

$$\Lambda = \frac{|E|}{|E+H|}$$

(2)

Where, E and H are given by:

$$\mathbf{E} = \mathbf{Y}^T \mathbf{Y} - \hat{\mathbf{B}}^T \mathbf{X}^T \mathbf{Y} \quad (3)$$

$$\mathbf{H} = \hat{\mathbf{B}}^T \mathbf{X}^T \mathbf{Y} - n \bar{\mathbf{y}} \bar{\mathbf{y}}^T \quad (4)$$

where,

$$\hat{\mathbf{B}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y} \quad (5)$$

In case of $\Lambda \leq \Lambda_{\alpha, p, v_E, v_H}$ the null hypothesis H_0 will be rejected, where $\Lambda_{\alpha, p, v_E, v_H}$ is the critical values degrees of freedom p , $v_E = n - q - 1$ and $v_H = q$ and with level of significance α in Wilks' critical values table. When the Wilks' statistic yields significant, which means that its associated p -value falls below a predetermination significance threshold α , this gives evidence to reject the null hypothesis, which suggests that among the independent variables, there is at least one that has a noteworthy influence on the set of dependent variables as a group. Assuming that Y is distributed as multivariate normal distribution, the vast majority of classical statistics are intensely responsive to the impact of outliers (see

[4]). Many multivariate robust estimators of location a scatter which are resilient to feasible outliers in the data-set have been proposed, M-estimator Maronna 1976 [5], the minimum covariance determinant estimator (MCD) Rousseeuw (1984) [6], S-estimator (Davies 1987, Rousseeuw and Leroy 1987, Lopuhaa and Hendrik 1989) [7-9]. In high dimensions a robust estimator of location and scatter investigated by (Woodruff and Rocke 1994) [10]. The impact of outliers on the Wilks' statistic will be discussed in a simulation study outlined in section IV. There for, we present an alternative robust Wilks' statistic to the classical Wilks' statistic. The Minimum Covariance Determinant Estimator (MCD) that suggested by Rousseeuw in (1984) [6] which is highly robust estimator of scatter and location matrices, for this privilege they are used. To enhance efficiency while remain high robustness, one can utilizes re-weighted stops for MCD estimator, a summary of the MCD estimator is presented in section II. The proposed approximations have been constructed with accuracy examination which is shown in section III. In section IV a simulation study is used to assess the performance of the proposed statistics and to compare various test statistics in various cased considering factors such as robustness, the power of the test and significance level. In section V a real data is used to conduct a more in-depth assessment of the proposed robust statistics.

II. ROBUST ESTIMATOR

To construct the robust wilks' statistic, we need to estimate the multivariate parameters of the model. The MCD estimator of Rousseeuw (1984) [6] is extremely robust estimator of multivariate scatter matrix and location matrix. The MCD estimators seek to identify a subset of z observations that minimizes the sample covariance matrix determinant, where the subset size z is choose to be in the range from the half size of the sample to the full size of the sample. The covariance matrix of the subset gives the estimated scatter matrix \mathbf{C} of the MCD and the mean vector of the subset gives the estimated location \mathbf{L} of the MCD. The effective algorithms for computing the MCD estimates available in widely-used software programming languages such as *R*, *python*, *SAS* and *Matlab*. To intensify the effectiveness of the MCD, a robust re-weighted version has been used. To find the estimated covariance matrix for the MCD many methods have been suggested, (He and Fung 2000) [11] and (Huber and Van Driessen 2004) [12].

III. THE PROPOSED WILKS' APPROXIMATION STATISTICS

Due to the complicity of the classical Wilks' distribution which was introduced by Anerson (1958) [13], we will use Bartlett approach for the Wilks' statistics distribution which is defined by (see [14]):

$$-\left(v_E - \frac{1}{2}(p - v_H + 1)\right) \ln(\Lambda) \approx \chi_{pv_H}^2 \quad (6)$$

our proposed approximation is an alternative Wilks' statistic based on RMCD estimators which defined as:

$$\Lambda_R = \frac{|E_R|}{|E_R + H_R|} \quad (7)$$

where H_R and E_R are given by:

$$H_R = Y^T \left(W X (X^T W X)^{-1} X^T W - \frac{J_W}{\sum_{i=1}^n w_i} \right) Y \quad (8)$$

$$E_R = Y^T (W - W X (X^T W X)^{-1} X^T W) Y \quad (9)$$

where,

$$W = \text{diag}(w_i), \quad i = 1, 2, \dots, n, \\ J_w = \mathbf{w}^T \mathbf{w}, \quad \mathbf{w} = (w_1, w_2, \dots, w_n)^T.$$

In our present study, we propose the following recommendations:

Step-1. Compute the estimated location vector $\hat{\boldsymbol{\mu}}$ and the estimated covariance matrix $\hat{\boldsymbol{\Sigma}}$ by using the MCD.

Step-2. Calculate the weights w_i of the observation \mathbf{y}_i by using the Hampel weight function (see Campbell, (1980) [15]) and Huber weight function those defined as:

$$w_i = \begin{cases} 1, & MD(\mathbf{y}_i) \leq d_0 \\ \frac{d}{MD(\mathbf{y}_i)}, & MD(\mathbf{y}_i) > d_0 \end{cases} \quad (10)$$

where,

$$d = d_0 e^{-\frac{1}{2} \left(\frac{MD(\mathbf{y}_i) - d_0}{b_2} \right)}, \quad d_0 = \sqrt{p} + \frac{b_1}{\sqrt{2}}, \quad b_1 = 2, \quad b_2 = 1.25,$$

and $MD(Y_i)$ is the Mahalanobis distances which given by:

$$MD(\mathbf{y}_i) = \sqrt{(\mathbf{y}_i - \hat{\boldsymbol{\mu}}^0)^T (\hat{\boldsymbol{\Sigma}}^0)^{-1} (\mathbf{y}_i - \hat{\boldsymbol{\mu}}^0)} \quad (11)$$

And w_i with Huber weight function is given by:

$$w_i = \begin{cases} 1, & MD(\mathbf{y}_i) \leq \sqrt{\chi_{0.975}^2(p)}, \\ 0, & \text{wtherwise} \end{cases} \quad (12)$$

Step-3. For $j = j + 1$ calculate the weighted estimated location matrix $\hat{\boldsymbol{\mu}}$ and the weighted estimated covariance matrix $\hat{\boldsymbol{\Sigma}}$ as following:

$$\hat{\boldsymbol{\mu}}^j = \frac{1}{\sum_{i=1}^n w_i} \sum_{i=1}^n w_i \mathbf{y}_i \quad (13)$$

$$\hat{\boldsymbol{\Sigma}}^j = \frac{1}{\sum_{i=1}^n w_i} \sum_{i=1}^n w_i (\mathbf{y}_i - \hat{\boldsymbol{\mu}}^j) (\mathbf{y}_i - \hat{\boldsymbol{\mu}}^j)^T \quad (14)$$

Step-4. Repeat steps two and three until the following condition holds:

$$\frac{|\hat{\boldsymbol{\Sigma}}^{j+1}/p|^p}{\det(\hat{\boldsymbol{\Sigma}}^{j+1})} \leq \frac{|\hat{\boldsymbol{\Sigma}}^j/p|^p}{\det(\hat{\boldsymbol{\Sigma}}^j)}$$

Now we will present robust forms of Wilks' statistics which are much like Λ_R (7), namely Λ_{R_1} and Λ_{R_2} , but relied on re-weighted minimum covariance determinant estimator with Huber wight function and Hampel weight function and construct their approximate distribution.

$$-\left(v_{E_R} - \frac{1}{2}(p - v_{H_R} + 1)\right) \ln \Lambda_{R_W} \approx \chi_{pv_{H_R}}^2 \quad (15)$$

We can find the degrees of freedom v_{E_R} , and v_{H_R} to the robust Wilks' statistic Λ_R as follows:

$$v_{HR} = trace \left(WX(X^T WX)^{-1} X^T W - \frac{JW}{\sum_1^q w_i} \right) \quad (16)$$

$$v_{ER} = trace \left(W - WX(X^T WX)^{-1} X^T W \right) \quad (17)$$

Now we will use the QQ-plots technique to examine the accuracy of Λ_{R1} and Λ_{R2} by the simulation for $k = 3000$ samples of the multivariate normal distribution and many cases of p the number of dependent variables, q the number of the independent variables and n the sample size. The classical distribution of the k statistics will undergo comparison to the approximate distribution of Λ_{R1} and Λ_{R2} by using QQ-plots, we will insert some of the plots in Fig.1 and Fig .2. The standard cut-off values of a test, 0.95, 0.975 and 0.99 are presented in these plots as vertical lines. The plots show that the approximations are accurate for all the dimensions under the study $p \in \{2,3,4\}$ and $q \in \{2,3,5\}$, large and small sizes, and for medium and high correlation r and no correlation $r = 0$.

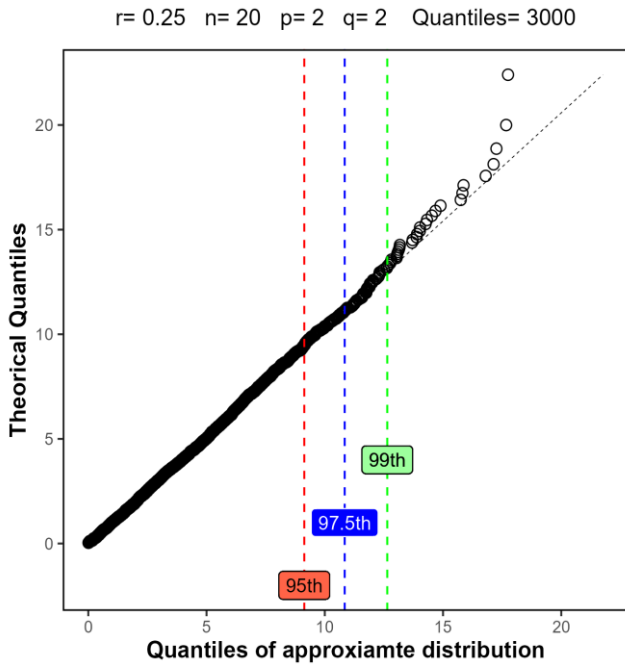


Fig.1: Λ_{R1} QQ-plots in case of $n=20, p=2, q=2$

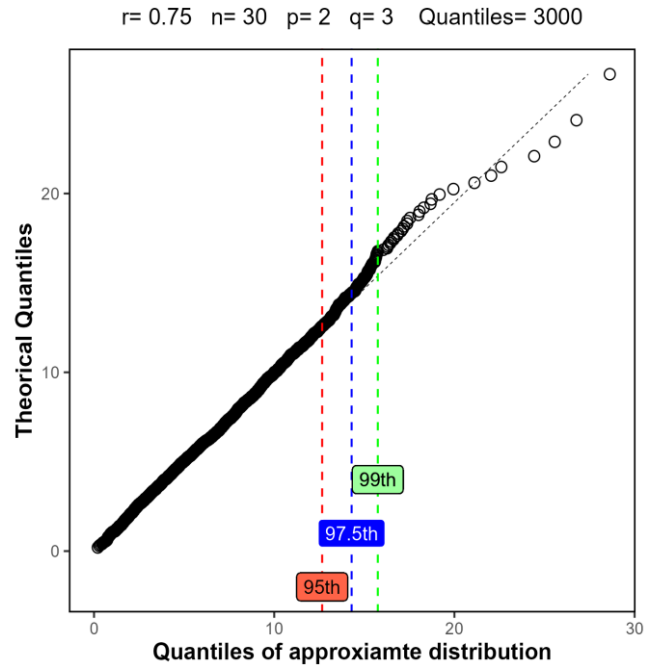


Fig.2: Λ_{R2} QQ-plots in case of $n=30, p=2, q=3$

IV. MONTE CARLO SIMULATION

Monte Carlo process is a valuable technique for assessing the performance goodness of the test statistics. The performance assessment of the test statistics based on type I error rate and the power of the test, using these two measures we will compare the behavior of the robust statistics and the classic Wilks' statistic, once in case of the data consist outliers and other in terms of the data is completely normally distributed. In order to study the type I error rate and the power of the test, we will consider number of cases, independent variables $p = 2, 3, 4$, dependent variables $q = 2, 3, 5$, sample sizes $n = 20, 30, 40$ and the correlation between the components of the dependent variable Y , no correlation $r = 0$, medium correlation $r = 0.5$, and strong correlation $r = 0.75$.

1. Significance Level

In order to evaluate and compare the rates of type I error $\hat{\alpha}$ of the test statistics under the study, we generate the observations of different sizes n once from the multivariate normal distribution $y_i \sim N_p(\mathbf{0}, \mathbf{I})$ and another whose contain outliers using the following model:

$$\begin{aligned} y_m &\sim N_p(\mathbf{0}, \mathbf{I}), & m &= 1, 2, \dots, [t], \\ y_s &\sim N_p(\mu^*, c\mathbf{I}), & s &= [t] + 1, \dots, n \end{aligned}$$

where, $t = \frac{80n}{100}$, $[t]$ is the largest positive integer that is not less than t , $\mu^* = v^2 \sqrt{\chi_{p,0.001}^2} \mathbf{1}_p^T$, $v = 5, c = 0.0625$ and $\mathbf{1}_p$ is a vector whose all entries are 1, under the null hypothesis $H_0 : B_1 = \mathbf{0}$, where B_1 is a matrix that consist all the rows of the coefficient matrix B except the first row. The classical Wilks' statistic Λ is compared to the Bartlett' χ^2 formula eq (7), the robust Wilks' statistics that we proposed are compared to the approximate distribution (Huber and Hampel) which given in section 3. In the simulation this process has been repeated for $k = 3000$ times to calculate $\hat{\alpha} = \frac{T(k)}{k}$ (where $T(k)$ is the number of the times that null hypothesis been rejected when it is true.) The values $\hat{\alpha}$ are considered to be the estimate of the threshold significance level when the simulated critical values are exceeded the threshold value of the significance level which is selected to be 0.05, we can get the nominal level interval from Fawcett and Salter standard error formula (1989) [16], $\alpha \pm 2 \times \sqrt{\frac{\alpha(1-\alpha)}{k}}$ gives the standard divination interval about the nominal level. We opt the plots of the P-value that suggested by Davidson and McKinnon (1998) [17], as it provides a more comprehensive representation of how the test statistics conform to the approximate distribution under the null hypothesis within the simulated samples. Fig.1 and Fig.2 show the plots of the P-value, the statistics Λ , Λ_{R_1} and Λ_{R_2} are close to the 45° line.

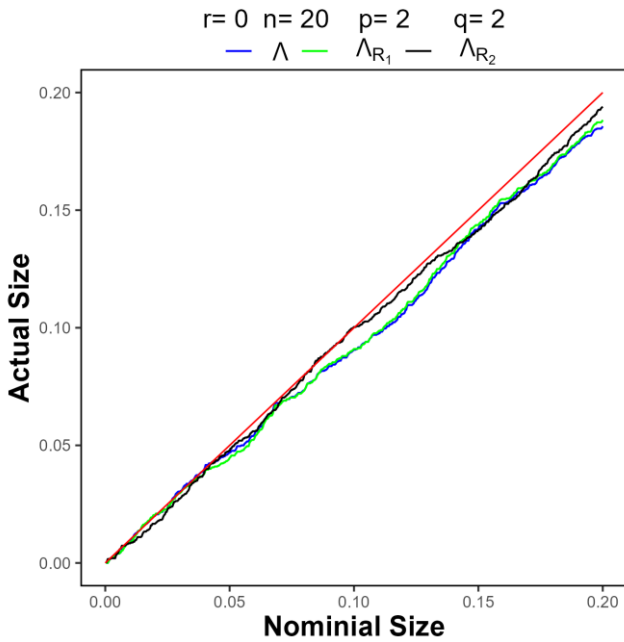


Fig.3: P-value in case of normal data

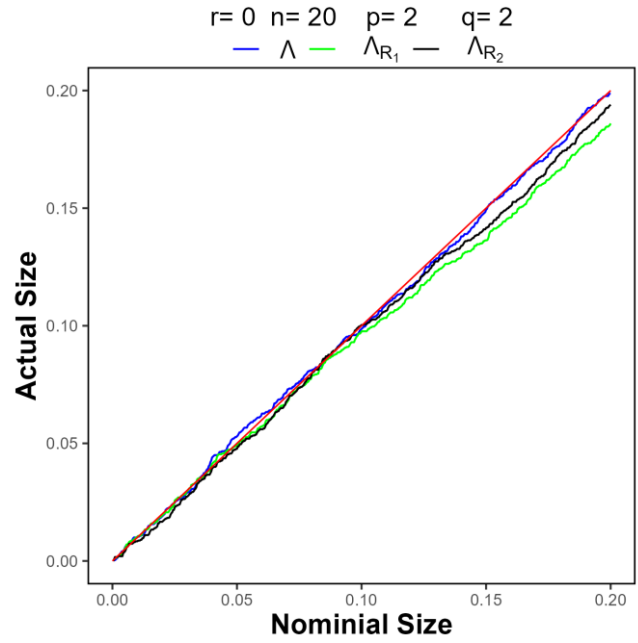


Fig.4: P-value in case of normal data

2.Power of the test

To find the power of the test of the statistics under the study and do comparisons between them, we generate the observations of different sizes n once from the multivariate normal distribution $\mathbf{y}_i \sim N_p(\mathbf{0}, \mathbf{I})$ and another whose contain outliers using the model:

$$\begin{aligned} \mathbf{y}_m &\sim N_p(\mathbf{0}, \mathbf{I}), & m = 1, 2, \dots, [t], \\ \mathbf{y}_s &\sim N_p(\mu^*, c\mathbf{I}), & s = [t] + 1, \dots, n \end{aligned}$$

where $t = \frac{80n}{100}$, $[t]$ is the largest positive integer that is not less than t , $\mu^* = v^2 \sqrt{\chi_{p,0.001}^2} \mathbf{1}_p^T$, $v = 5, c = 0.0625$ and $\mathbf{1}_p$ is a vector whose all entries are 1, under the alternative hypothesis. We used Davidson and MacKinnon (1998) [18] method to compare the resulting power size curves. Fig.5, Fig.6, Fig.7 and Fig.8 show the statistics Λ , Λ_{R_1} and Λ_{R_2} are close the each other in term of size power and with respect to the correlation between the dependent variables, where the data-sets are normally distributed and contains no outliers. In case the data has been corrupted and in case of no-correlation and low-correlation between the dependent variables the statistics Λ_{R_1} and Λ_{R_2} works way better than Λ as shown in Fig.9, Fig.10 and Fig.11, in case of high correlation the statistics Λ , Λ_{R_1} and Λ_{R_2} are close to each other as shown in Fig.12.

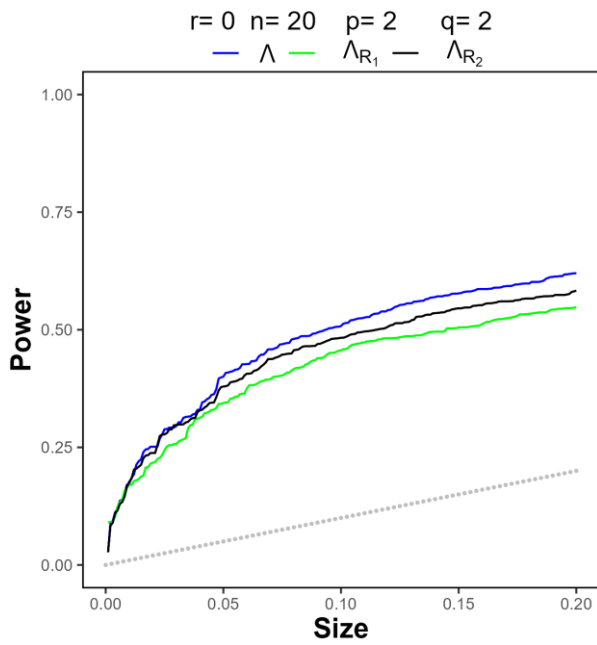


Fig 5: Curves of size power normal data

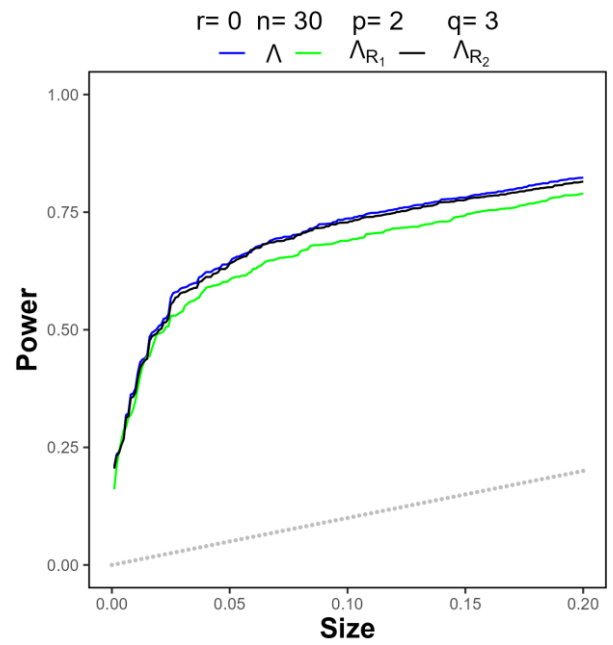


Fig 6: Curves of size power normal data

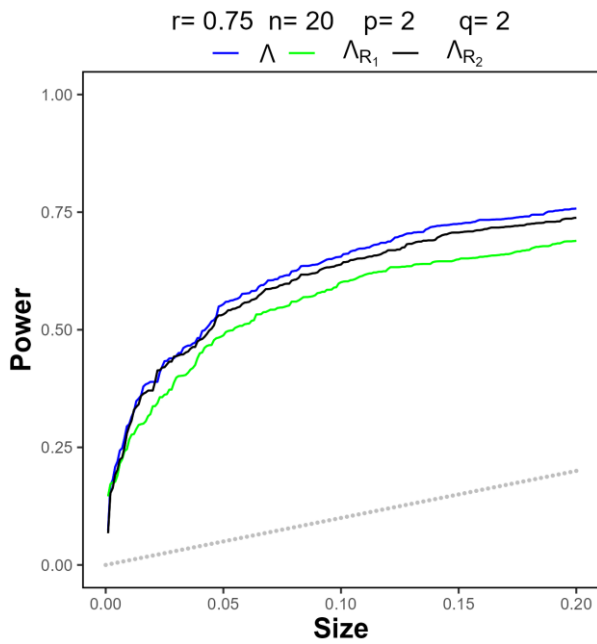


Fig 7: Curves of size power, normal data

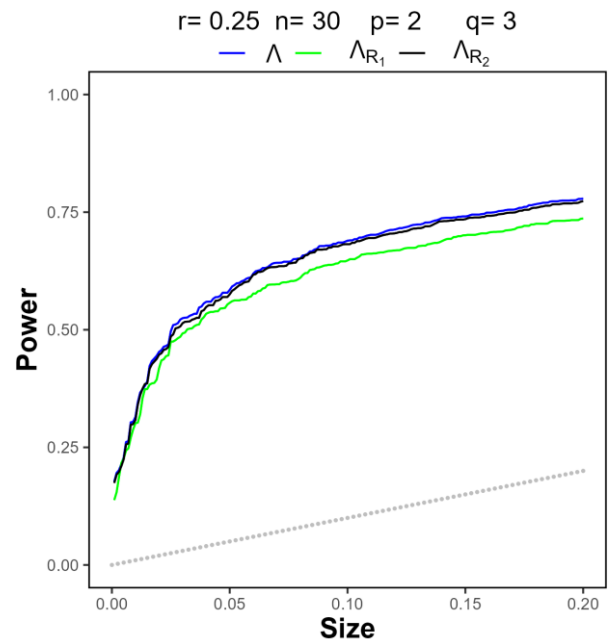


Fig 8: Curves of size power, normal data

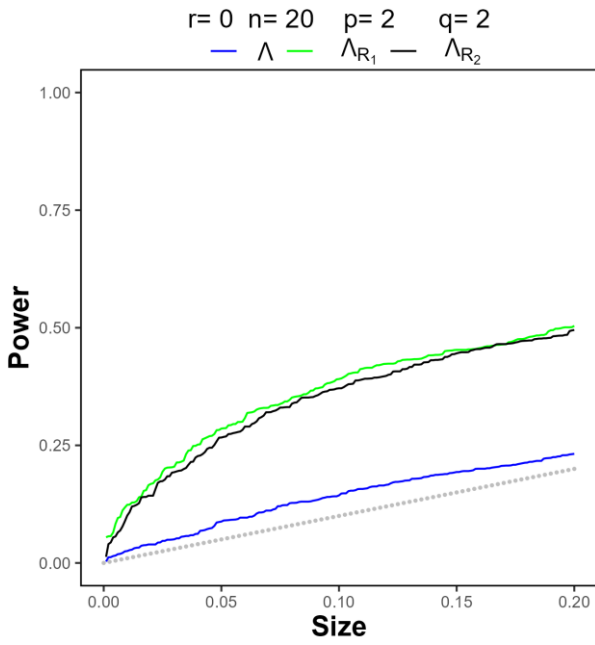


Fig 9: Curves of size power, corrupted data

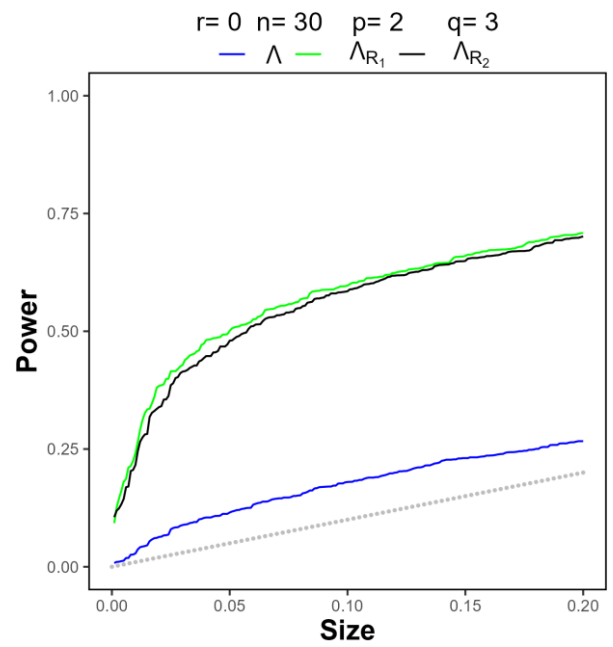


Fig 10: Curves of size power, corrupted data

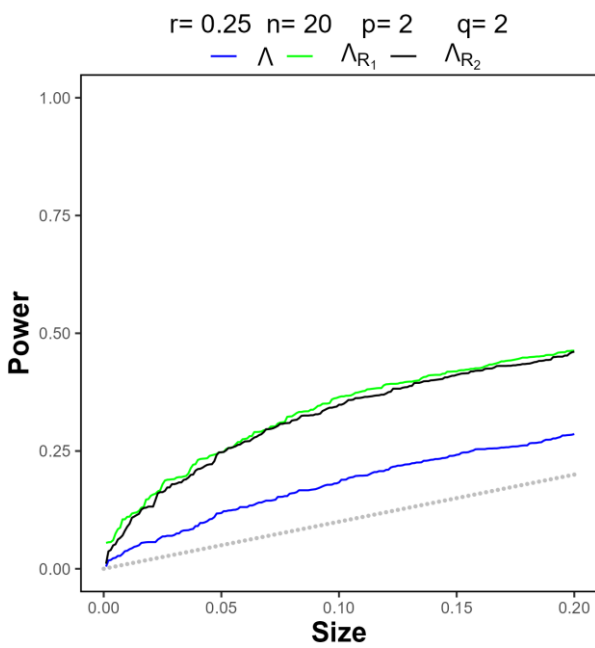


Fig.11: Curves of size power, corrupted data

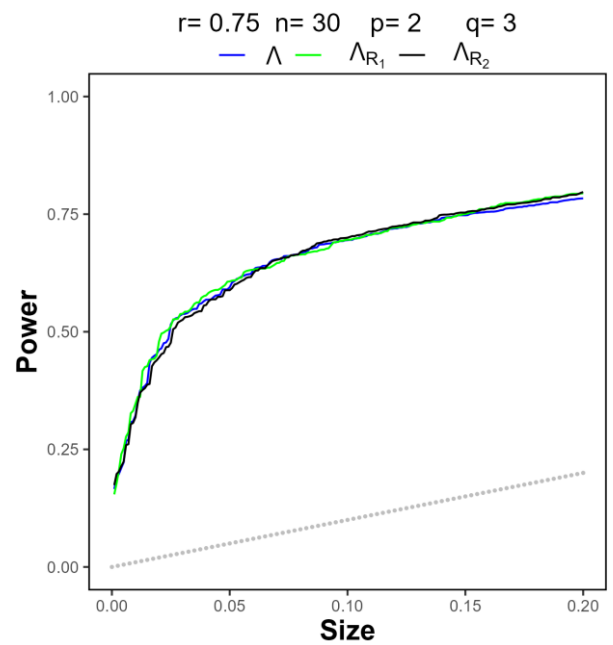


Fig.12: Curves of size power, corrupted data

V. REAL DATA

In this section we used a real data set (see [14] chapter 3), normal patients and diabetics where under the study to measure their glucose intolerance, insulin response to oral glucose, insulin resistance, relative weight and fasting plasma glucose, table I below shows the first 10 observations.

TABLE I. REAL DATA OF DIABETES PATIENT (FIRST 10 OBSERVATIONS ONLY)

n	y_1	y_2	x_1	x_2	x_3
1	0.81	80	356	124	55
2	0.95	97	289	117	76
3	0.94	105	319	143	105
4	1.04	90	356	199	108
5	1	90	323	240	143
6	0.76	86	381	157	165
7	0.91	100	350	221	119
8	2.1	85	301	186	105
9	0.99	97	379	142	98
10	0.78	97	296	131	94

Where, the independent variables are:

- x_1 = glucose intolerance,
- x_2 = insulin response to oral glucose,
- x_3 = insulin resistance,

and the dependent variables are:

- y_1 = relative weight,
- y_2 = fasting plasma glucose.

To test the normality of the data set we will use:

A. *Mardia Multivariate Normality Test*: The Mardia statistical test is a statistical examination that using to assess if a given multivariate dataset follows the multivariate normal distribution (see [18]). This test named after Jagdish K.Mardia, who introduced it in his book “Measures of Multivariate Skewness and Kurtosis with Applications.”[19].

TABLE II. MARDIA TEST OF MULTIVARIATE NORMALITY

Test	Statistic	P value	Result	MVN Result
Mardia Skewness	101.10593	5.71975e-21	NO	NO
Mardia Kurtosis	13.53029	0.0	NO	

B. *The Royston Multivariate Normality Test*: Often referred to as the Royston Test, is another statistical examination used to determine whether a given multivariate dataset follows the multivariate normal distribution or not. This test was introduced by Martin Royston [20-21], which is an extension of Shapiro test.

TABLE III. MULTIVARIATE NORMALITY ROYSTON TEST

Test	H	P value	MVN Result
Royston	35.61727	1.844201e-8	NO

Testing of the hypothesis was under Λ , Λ_{R_1} and Λ_{R_2} . The test cannot be rejected according to the corresponding P-value which is shown in table 4 at $\alpha = 0.05$, the proposed statistics Λ_{R_1} and Λ_{R_2} give powerful evidence to reject the testing hypotheses.

TABLE IV. P-VALUE OF CLASSICAL WILKS’ STATISTIC, Λ_{R_1} AND Λ_{R_2} METHODS.

Method	P value
Λ	0.703493688
Λ_{R_1}	0.017180954
Λ_{R_2}	0.017180938

VI. CONCLUSIONS

We introduced enhanced robust forms of Wilks’ statistic based on re-weighted minimum covariance determinant estimator Λ_{R_1} and Λ_{R_2} and formulated their approximated distributions. The results indicate that the proposed statistics are close to the classic Wilks’ statistic Λ in case of the data is normally distributed, while the proposed statistics are best of the classic Wilks’ in case of the contaminated distribution.

REFERENCES

- [1] Leo Breiman and Jerome H Friedman. Predicting multivariate responses in multiple linear regression. Journal of the Royal Statistical Society Series B: Statistical Methodology, 59(1):3–54, 1997.
- [2] James B Davis and Joseph W McKean. Rank-based methods for multivariate linear models. Journal of the American Statistical Association, 88(421):245–251, 1993.
- [3] Esa Ollila, Hannu Oja, and Visa Koivunen. Estimates of regression coefficients based on lift rank covariance matrix. Journal of the American Statistical Association, 98(461):90–98, 2003.
- [4] S Frosch Møller, Jürgen von Frese, and Rasmus Bro. Robust methods for multivariate data analysis. Journal of Chemometrics: A Journal of the Chemometrics Society, 19(10):549–563, 2005.
- [5] Ricardo Antonio Maronna. Robust m-estimators of multivariate location and scatter. The annals of statistics, pages 51–67, 1976.

-
- [6] Peter J Rousseeuw. Least median of squares regression. *Journal of the American statistical association*, 79(388):871–880, 1984.
- [7] P Laurie Davies. Asymptotic behaviour of s-estimates of multivariate location parameters and dispersion matrices. *The Annals of Statistics*, pages 1269–1292, 1987.
- [8] Peter J Rousseeuw and Annick M Leroy. A robust scale estimator based on the shortest half. *Statistica Neerlandica*, 42(2):103–116, 1988.
- [9] Hendrik P Lopuhaa. On the relation between s-estimators and m-estimators of multivariate location and covariance. *The Annals of Statistics*, pages 1662–1683, 1989.
- [10] David L Woodruff and David M Roche. Computable robust estimation of multivariate location and shape in high dimension using compound estimators. *Journal of the American Statistical Association*, 89(427):888–896, 1994.
- [11] Xuming He and Wing K Fung. High breakdown estimation for multiple populations with applications to discriminant analysis. *Journal of Multivariate Analysis*, 72(2):151–162, 2000.
- [12] Mia Hubert and Katrien Van Driessen. Fast and robust discriminant analysis. *Computational Statistics & Data Analysis*, 45(2):301–320, 2004.
- [13] Curtis A Parvin. *An Introduction to Multivariate Statistical Analysis*, 3rd ed. T.W. Anderson. Hoboken, NJ: John Wiley amp; Sons, 2003, 742 pp., 99.95, hardcover. ISBN0 – 471 – 36091 – 0. *Clinical Chemistry*, 50(5): 981 – –982, 052004.
- [14] Sons Rencher, John Wiley. *Methods of Multivariate Analysis*, Second Edition. Brigham Young University, 2002.
- [15] Norm A Campbell. Robust procedures in multivariate analysis i: Robust covariance estimation. *Journal of the Royal Statistical Society Series C: Applied Statistics*, 29(3):231–237, 1980.
- [16] KC Salter and RF Fawcett. A robust and powerful rank test of treatment effects in balanced incomplete block designs. *Communications in Statistics-Simulation and Computation*, 14(4):807– 828, 1985.
- [17] Russell Davidson and James G MacKinnon. Graphical methods for investigating the size and power of hypothesis tests. *The Manchester School*, 66(1):1–26, 1998.
- [18] Selcuk Korkmaz, Din,cer G“oks“ul“uk, and GOKMEN Zararsiz. Mvn: An r package for assessing “ multivariate normality. *R JOURNAL*, 6(2), 2014.
- [19] Kanti V Mardia. Measures of multivariate skewness and kurtosis with applications. *Biometrika*, 57(3):519–530, 1970.
- [20] J Patrick Royston. An extension of shapiro and wilk’s w test for normality to large samples. *Journal of the Royal Statistical Society: Series C (Applied Statistics)*, 31(2):115–124, 1982.
- [21] Patrick Royston. Remark as r94: A remark on algorithm as 181: The w-test for normality. *Journal of the Royal Statistical Society. Series C (Applied Statistics)*, 44(4):547–551, 1995.