



## Comparison between some Parametric Robust Methods for Estimating the Parameters of the Multiple Normal Distribution

M. Murtadha Mansour Abdullah

University of Wasit

College of Administration and Economics

mabdullah@uowasit.edu.iq

### **Abstract**

Estimating the parameters of the Multivariate Normal Distribution is very important process in many statistical Application like the Principal Component Analysis or Canonical Analysis. The paper aims at finding robust and efficient estimators for the parameters of the multivariate normal distribution by using parametric method which is the reweighted minimum vector method RMV and compare it with another robust estimation methods like S estimation method in case of deferent sample sizes and deferent contaminated ratios. Results shows that the reweighted minimum vector is the best method via simulation and real data was taken from waist sewerage directorate minimum esquire error (MMSE) was used as a comparison tool between the two estimation methods.

### **1. Introduction**

Estimating the parameters of the normal distribution is very important and main process [4] for many statistical applications and yet it will be more important when we deal with multivariate statistics as the data take deferent features. Estimating the mean vector end the covariance matrix of the multivariate normal distribution is a milestone for many important statistical analysis methods such as factor analysis and discernment analysis the classical methods of estimation location and parameters for the multivariate normal distribution give weak and non-efficient estimates especially when we have large dimension or we have any problem in

the assumption of linear regression model to avoid that we use here the RMV and S estimators as new methods for estimating the multivariate Normal Distribution parameters [7, 11]

## 2. The Multivariate Normal Distribution

The multivariate normal distribution function deals with random vectors of multivariate scale units with  $(n * 1)$  dimension in a sample of observations let

$\underline{X} = (X_1, X_2 \dots X_n), X_j \in R \forall j$  then the multivariate normal distribution function will be [ 11]

$$f(\underline{X}) = \frac{1}{(2\pi)^{\frac{p}{2}} |\Sigma|^{\frac{1}{2}}} \exp^{-\frac{1}{2}(\underline{X}-\underline{M})\Sigma^{-1}(\underline{X}-\underline{M})} \dots \dots \dots (1)$$

Where

$\underline{\mu}$  : is the mean veaor

$\Sigma$  : is the var –covariance matrix

## 3. The Reweighted Minimum Vector Method (RMV)

This method is considered as a parametric method that depends upon giving weights (0) for the extreme values which  $(d_{RMV}^2 > X_{n,0.025}^2)$  and (1) for the no extreme values then the RMV for the scale and location parameters will be [2,10]

$$\bar{X}_{RMVV}^{ROW} = \frac{\sum_{i=1}^M w_i X_i}{M} \dots \dots \dots (2)$$

$$S_{RMVV} = \phi_{m,n,p}^{*\alpha} K^*(M) \frac{\sum_{i=1}^M w_i (X_i - \bar{X}_{RMVV}^{ROW})(X_i - \bar{X}_{RMVV}^{ROW})^-}{M} \dots \dots \dots (3)$$

Where

$d_{RMV}^2$ : represents the maximum values of the parameters

M : represent the number of observation  $d_{RMVV}^2 \leq X_{n,0.025}^2$

$$W \begin{cases} 0 & \text{if } d_{RMVV}^2 (X_i - \bar{X}_{rmvv}) > X_{n,0.025}^2 \dots \dots \dots (4) \\ 1 & o.w \end{cases}$$

Which  $d_{RMVV}^2$  Will be equal to

$$d_{MVV}^2 = (X_i - \bar{X}_{mrv})^{-1} S^{-1} (X_i - \bar{X}_{mrv}) \dots \dots \dots (5)$$

Then efficient factor

$$K^*(M) = \frac{m/n}{P(X_{P+2}^2 < X_{P,\frac{m}{n}}^2)} \dots \dots \dots (6)$$

Last the quadratic distances will be

$$L_i^2(X_i, \bar{X}_{RMVV}, S_{RMVV}) = (X_i - \bar{X}_{rmrv})^{-1} S^{-1} (X_i - \bar{X}_{rmrv}) \dots \dots \dots (7)$$

#### **4. S Method**

This method considered one of the important methods to have a high robust estimators for the multivariate normal distribution function. Moreover it gives good efficiency prosperity then the S estimator for scale parameter is [3]

$$b = \frac{1}{n} \sum_{i=1}^n p\{d(X_i, \bar{X}, S)\} \dots \dots \dots (8)$$

Where

$(d(X_i, \bar{Y}, s))$  is the mahalanobis distance which is in form [1]

$$d^2 = (X_i - \bar{X})^{-1} S^{-1} (X_i - \bar{X}) \dots \dots \dots (9)$$

And for good consistency we have

$$b = E_{F_0}[P(\|X\|)] = E\left(P\left(\frac{d}{k}\right)\right) \dots \dots \dots (10)$$

Where

$$F_0 = N(0, I_p)$$

And x is cycle random variable

Here we chose Tukey function as weight for pin (8) because it's had a good properties as follow [3 , 10]

$$P(MD_i) = \begin{cases} 1 - [1 - (\frac{MD_i}{k})^2]^3 & \text{if } |MD_i| \leq k_0 \\ 1 & \text{if } |MD_i| > k_0 \end{cases}$$

Where [MDi] is the Mahalanobis distances and to get high breakdown point we determine the constant ( $K_0$ ) which gives central property called B-weight estimator and b depends on the value of  $k_0$  as follow

$$k_0 = \sqrt{p \left\{ \sqrt{\left(\frac{1}{9}\right) \left(\frac{2}{p}\right) c + \left(1 - \left(\frac{1}{9}\right) \left(\frac{2}{p}\right)\right)} \right\}^3 \dots \dots \dots (11)}$$

Then the consistency parameter b can be calculated in the following form in case of normal assumptions [5, 9]

$$b = \frac{p}{2} X_{p+2}^2(K_0^2) - \frac{p(p+2)}{2K_0^2} X_{p+4}^2(K_0^2) + \frac{p(p+2)(p+4)}{6K_0^4} X_{p+6}^2(K_0^2) + \frac{k_0^2}{6} (1 - X_p^2(K_0^2)) \dots \dots \dots (12).$$

**5. Minimum Men Square Error (MMSE)**

Here we use the MMSE as comparison tool to determine which the best estimator by using the Fibonacci distance as follow [7,11]

$$MMSE = E\{\|\varepsilon - \hat{\varepsilon}\|_F^2} = \frac{Trcce(\varepsilon - \hat{\varepsilon})^2}{n} \dots \dots \dots (13)$$

Then MMSE is good method of comparison among deferent estimators of the var-covariance matrix and for the mean vector. The MMSE for it will simply be the square Euclidean distance in the following from .

$$\|u - r\|^2 = (u_1 - r_1)^2 + (u_2 - r_2)^2 + \dots + (u_n - r_n)^2 \dots \dots (14)$$

The MSE for the mean vector will be [6]

$$MSE = \frac{1}{P} \|\underline{M} - \underline{\widehat{M}}\|^2 = \frac{1}{P} ((M_1 - \widehat{M}_1)^2 + (M_2 - \widehat{M}_2)^2 + \dots + (M_n - \widehat{M}_n)^2) \dots (15)$$

## 6. Simulation Results

Box- Muller method were used to generate multivariate normal data that contaminated with deferent ratios with deferent sample sizes by using MIATLAB

We calculate the contaminated distribution with certain ( $\alpha$ ) the ratio of contaminated and ( $1 - \alpha$ ) is the non- contaminated data as follow [1 , 8]

$$X = (1 - \alpha)N(M, \sigma^2) + \alpha N(M + t, \sigma^2) \quad t \neq 0 \quad \dots \dots \dots (16)$$

T is the contaminating degree

We repeat this experiment (1000) time for deferent sample sizes (25 , 70 , 150 , 300 ) and ( $\alpha = 0.6 , 0.2 , 0.35$  ) where we get the var – covariance matrix and the mean vector which is the theatrical parameters pop as follow

Then

**Table (1) covariance matrix ( $\alpha = 0.05$  ,  $n=25$  ) RMV Method**

	<b>X1</b>	<b>X2</b>	<b>X3</b>	<b>X4</b>	<b>X5</b>
<b>X1</b>	<b>11.52</b>	<b>0.68</b>	<b>-1.6</b>	<b>0.85</b>	<b>-4.2</b>
<b>X2</b>		<b>25.56</b>	<b>-3.38</b>	<b>-2.53</b>	<b>4.84</b>
<b>X3</b>			<b>84.91</b>	<b>2.38</b>	<b>-4.3</b>
<b>X4</b>				<b>110.17</b>	<b>7.31</b>
<b>X5</b>					<b>242.52</b>

**Table (2) covariance matrix ( $\alpha= 0.05$  ,  $n=25$  ) S Method**

	<b>X1</b>	<b>X2</b>	<b>X3</b>	<b>X4</b>	<b>X5</b>
<b>X1</b>	11.41	2.21	-0.55	-3.44	-3.74
<b>X2</b>		62.31	-5.34	-3.51	4.55
<b>X3</b>			74.89	3.65	-4.66
<b>X4</b>				120.11	4.32
<b>X5</b>					232.86

**Table (3) covariance matrix ( $\alpha= 0.05$  ,  $n =70$ ) RMV Method**

	<b>X1</b>	<b>X2</b>	<b>X3</b>	<b>X4</b>	<b>X5</b>
<b>X1</b>	11.37	0.78	-1.8	-0.3	-3.43
<b>X2</b>		52.2	-3.36	-3.93	5
<b>X3</b>			94.21	4.72	-4.58
<b>X4</b>				100.12	7.21
<b>X5</b>					264.63

**Table (4) covariance matrix ( $\alpha= 0.05$  ,  $n=70$  ) S Method**

	<b>X1</b>	<b>X2</b>	<b>X3</b>	<b>X4</b>	<b>X5</b>
<b>X1</b>	<b>11.97</b>	<b>0.68</b>	<b>-0.98</b>	<b>-0.72</b>	<b>-4.13</b>
<b>X2</b>		<b>52.4</b>	<b>-3.53</b>	<b>-4.1</b>	<b>-2.77</b>
<b>X3</b>			<b>94.16</b>	<b>2.92</b>	<b>-5.73</b>
<b>X4</b>				<b>110.21</b>	<b>8.7</b>
<b>X5</b>					<b>272.13</b>

**Table (5) covariance matrix for contaminated ( $\alpha=0.05$  ,  $n=150$  ) RMV Method**

	<b>X1</b>	<b>X2</b>	<b>X3</b>	<b>X4</b>	<b>X5</b>
<b>X1</b>	<b>13.26</b>	<b>0.55</b>	<b>2.7</b>	<b>0.41</b>	<b>4.71</b>
<b>X2</b>		<b>45.88</b>	<b>4.63</b>	<b>4.51</b>	<b>5.91</b>
<b>X3</b>			<b>94.61</b>	<b>4.22</b>	<b>4.88</b>
<b>X4</b>				<b>120.76</b>	<b>8.43</b>
<b>X5</b>					<b>262.21</b>

**Table (6) covariance matrix for contaminated ( $\alpha=0.05$  ,  $n=150$  ) S Method**

	<b>X1</b>	<b>X2</b>	<b>X3</b>	<b>X4</b>	<b>X5</b>
<b>X1</b>	<b>11.76</b>	<b>0.98</b>	<b>-0.37</b>	<b>-0.96</b>	<b>4.33</b>
<b>X2</b>		<b>62.51</b>	<b>2.75</b>	<b>5.22</b>	<b>4.22</b>
<b>X3</b>			<b>52.61</b>	<b>1.95</b>	<b>3.22</b>

<b>X4</b>				<b>120.72</b>	<b>10.55</b>
<b>X5</b>					<b>262.31</b>

**Table (7) covariance matrix for contaminated ( $\alpha=0.05$  ,  $n=300$  ) RMV Method**

	<b>X1</b>	<b>X2</b>	<b>X3</b>	<b>X4</b>	<b>X5</b>
<b>X1</b>	<b>12.36</b>	<b>0.86</b>	<b>-1.5</b>	<b>0.23</b>	<b>-4.36</b>
<b>X2</b>		<b>45.87</b>	<b>5.86</b>	<b>-4.35</b>	<b>5.98</b>
<b>X3</b>			<b>95.31</b>	<b>4.17</b>	<b>-4.85</b>
<b>X4</b>				<b>99.9</b>	<b>7.65</b>
<b>X5</b>					<b>242.6</b>

**Table (8) covariance matrix for contaminated ( $\alpha=0.05$  ,  $n=300$  ) S Method**

	<b>X1</b>	<b>X2</b>	<b>X3</b>	<b>X4</b>	<b>X5</b>
<b>X1</b>	<b>20.11</b>	<b>0.92</b>	<b>-0.77</b>	<b>-0.63</b>	<b>-5.21</b>
<b>X2</b>		<b>45.36</b>	<b>-5.93</b>	<b>-5.93</b>	<b>3.12</b>
<b>X3</b>			<b>51.97</b>	<b>2.55</b>	<b>-5.65</b>
<b>X4</b>				<b>110.37</b>	<b>10.63</b>
<b>X5</b>					<b>252.62</b>

**Table (9) mean vector estimator for contaminated normal distribution with ( $\alpha=0.05$  ,  
 $n=25$  )**



	$\hat{M}_1$	$\hat{M}_2$	$\hat{M}_3$	$\hat{M}_4$	$\hat{M}_5$
<b>RMVV</b>	<b>17.7</b>	<b>21.5</b>	<b>60.23</b>	<b>88.12</b>	<b>120.11</b>
<b>S</b>	<b>15.24</b>	<b>30.29</b>	<b>60.11</b>	<b>75.31</b>	<b>110.11</b>

**Table (10) mean vector estimator for contaminated normal distribution with ( $\alpha=0.05$  ,  
 $n=70$  )**

	$\hat{M}_1$	$\hat{M}_2$	$\hat{M}_3$	$\hat{M}_4$	$\hat{M}_5$
<b>RMVV</b>	<b>17.22</b>	<b>22.15</b>	<b>60.75</b>	<b>88.41</b>	<b>120.04</b>
<b>S</b>	<b>15.31</b>	<b>30.98</b>	<b>60.52</b>	<b>76.21</b>	<b>110.21</b>

**Table (11) mean vector estimator for contaminated normal distribution with ( $\alpha=0.05$  ,  
 $n=150$  )**

	$\hat{M}_1$	$\hat{M}_2$	$\hat{M}_3$	$\hat{M}_4$	$\hat{M}_5$
<b>RMVV</b>	<b>17.12</b>	<b>26.33</b>	<b>60.23</b>	<b>85.11</b>	<b>118.81</b>
<b>S</b>	<b>15.11</b>	<b>30.01</b>	<b>61.21</b>	<b>81.21</b>	<b>111.11</b>

**Table (12) mean vector estimator for contaminated normal distribution with ( $\alpha=0.05$  ,  
 $n=300$  )**

	$\hat{M}_1$	$\hat{M}_2$	$\hat{M}_3$	$\hat{M}_4$	$\hat{M}_5$
<b>RMVV</b>	<b>16.82</b>	<b>25.87</b>	<b>59.21</b>	<b>88.35</b>	<b>120.31</b>
<b>S</b>	<b>15.26</b>	<b>29.91</b>	<b>59.81</b>	<b>80.21</b>	<b>111.21</b>

**Table (13) MSE for covariance matrix for the methods ( $\alpha = 0.05$  )**

<b>%5</b>	<b>N</b>	<b>25</b>	<b>70</b>	<b>150</b>	<b>300</b>
	<b>RMVV</b>	<b>1.521</b>	<b>1.331</b>	<b>1.221</b>	<b>0.851</b>
	<b>S</b>	<b>3.551</b>	<b>3.124</b>	<b>2.957</b>	<b>2.225</b>

**Table (14) MSE for covariance matrix for the methods ( $\alpha=10\%$ )**

<b>%10</b>	<b>N</b>	<b>25</b>	<b>70</b>	<b>150</b>	<b>300</b>
	<b>RMVV</b>	<b>1.211</b>	<b>1.031</b>	<b>1.001</b>	<b>0.671</b>
	<b>S</b>	<b>7.521</b>	<b>4.231</b>	<b>5.211</b>	<b>6.321</b>

**Table (15) MSE for covariance matrix for the methods ( $\alpha=25\%$ )**

<b>%25</b>	<b>N</b>	<b>25</b>	<b>70</b>	<b>150</b>	<b>300</b>
	<b>RMVV</b>	<b>3.291</b>	<b>5.221</b>	<b>2.915</b>	<b>0.664</b>
	<b>S</b>	<b>12.551</b>	<b>28.121</b>	<b>11.951</b>	<b>9.541</b>

**Table (16) MSE for mean vector for the methods ( $\alpha=5\%$ )**

<b>%5</b>	<b>N</b>	<b>25</b>	<b>70</b>	<b>150</b>	<b>300</b>
	<b>RMVV</b>	<b>1.841</b>	<b>1.981</b>	<b>1.731</b>	<b>1.271</b>
	<b>S</b>	<b>3.662</b>	<b>4.321</b>	<b>4.388</b>	<b>2.707</b>

**Table (17) MSE for mean vector for the methods ( $\alpha=10\%$ )**

<b>%10</b>	<b>N</b>	<b>25</b>	<b>70</b>	<b>150</b>	<b>300</b>
	<b>RMVV</b>	<b>1.212</b>	<b>0.987</b>	<b>0.731</b>	<b>0.394</b>
	<b>S</b>	<b>1.725</b>	<b>1.425</b>	<b>1.265</b>	<b>1.131</b>

**Table (18) MSE for mean vector for the methods ( $\alpha=25$ )**

<b>%25</b>	<b>N</b>	<b>25</b>	<b>70</b>	<b>150</b>	<b>300</b>
	<b>RMVV</b>	<b>3.491</b>	<b>1.521</b>	<b>0.845</b>	<b>0.562</b>
	<b>S</b>	<b>5.211</b>	<b>3.744</b>	<b>3.211</b>	<b>2.451</b>

From results in the tables above - we can notice that MMSE get smaller with the large sample size for both estimation methods in all contaminated rates the results show that the RMV is the best method to estimate the parameters of the multivariate normal distribution via MSE and its better than S method .

## **7. Real data**

There is many pollution sources for river water and rainwater drainage network is the main sources of pollution. We collect data from the water pollution labrotary in Wasit sewerage directorate in Wasit state by taking samples of water from deferent spots in 51 week by helping hand from the employees of Wasit sewerage directorate we determine 5 variables as a pollution sources.

$X_1$ : O8G represent the oil and greases in water

$X_2$ : SO4 represent the sulfates in water

$X_3$ : AIK represent the alkaline in water

$X_4$ : PH represent the acidic in water

$X_5$ : TOS represent the total dissolved salts in water

Where the data tested to normal distribution and we estimate the location and the scale parameters for the multivariate normal distribution by using RMV and S methods by using MATLAB program and we calculate the MSE for covariance matrices and the mean vectors as follow

**Table (19) estimated scale parameters by using RMV**

	<b>X1</b>	<b>X2</b>	<b>X3</b>	<b>X4</b>	<b>X5</b>
<b>X1</b>	<b>100.25</b>	<b>-1.55</b>	<b>1.62</b>	<b>62.89</b>	<b>44.55</b>
<b>X2</b>		<b>210.51</b>	<b>31.55</b>	<b>-59.25</b>	<b>-110.21</b>
<b>X3</b>			<b>562.11</b>	<b>-259.51</b>	<b>151.22</b>
<b>X4</b>				<b>975.24</b>	<b>-582.11</b>
<b>X5</b>					<b>1251.22</b>

**Table (20) estimated scale parameters by using S method**

	<b>X1</b>	<b>X2</b>	<b>X3</b>	<b>X4</b>	<b>X5</b>
<b>X1</b>	<b>100.15</b>	<b>-1.76</b>	<b>1.5</b>	<b>64.21</b>	<b>40.22</b>
<b>X2</b>		<b>210.93</b>	<b>36.22</b>	<b>-63.96</b>	<b>-113.21</b>
<b>X3</b>			<b>562.21</b>	<b>-267.11</b>	<b>149.22</b>
<b>X4</b>				<b>975.19</b>	<b>-552.11</b>
<b>X5</b>					<b>1251.10</b>

**Table (21) estimated mean vector by using RMV and S method**

	$\hat{M}_1$	$\hat{M}_2$	$\hat{M}_3$	$\hat{M}_4$	$\hat{M}_5$
<b>RMVV</b>	<b>8.51</b>	<b>35.51</b>	<b>60.91</b>	<b>100.11</b>	<b>40.21</b>
<b>S</b>	<b>5.11</b>	<b>42.48</b>	<b>80.22</b>	<b>110.22</b>	<b>180.22</b>

**Table (22) MSE for estimated scale parameters by using RMV and S method**

<b>RMVV</b>	<b>S</b>
<b>1.229</b>	<b>7.521</b>

**Table (23) MSE for estimated mean vector by using RMV and S method**

<b>RMVV</b>	<b>S</b>
<b>0.421</b>	<b>2.551</b>

## **8. Conclusion**

We can see that the value of MSE for covariance matrix is decreasing as the sample size get larger and the results show that the RMV method is better than S method to estimate the covariance and the mean vector in small and large sample size we recommend the method of RMV as an efficient method to estimate location and scale parameters for the multivariate normal distribution

## **9. References**

- 1- Ahmed, A. B. (2008). Standard Modeling of National Energy Consumption in Algeria during the Period (1988: 10-2007: 03). (Master Unpublished). University of Algeria, Faculty of Economic Sciences and Management Sciences.
- 2- Olive, D. J. (2018) .Robust Multivariate analysis .springer.
- 3-Ali, H. Yahiya, s. s. s. & Omar, Z. (2014, June). The efficiency of reweighted minimum vector variance. In AIP conference proceedings (Vol. 1602, NO. 1, pp.1151-1156).AIP.
- 4- Bable, B., & Pawar, D. (2012). Vector time series: the concept and properties to the vector stationary time series. International Research Journal of Agricultural Economics and Statistics, 3(1), 84-95.
- 5- Verardi, V. & Mc Cathie, A. (2012). The S-estimator of multivariate Location and Scatter in Stata. Stata journal, 12(2), 299
- 6- Huang, S.-C. (2008). Combining wavelet-based feature extractions with SVMs for financial time series forecasting. Journal of Statistics and Management Systems, 11(1), 37-48.
- 7- Hossjer, O., Croux, C, & Rousseeuw, p. j. (1994) Asymptotic of generalized s-estimation journal of Multivariate Analysis, 51(1), 148-177.
- 8-Croux, C. Rousseeuw, p. J & Hossjer, O. (1994). "Generalized s-estimator", JASA, Vol. 89, NO. 428, 1271-1281.
- 9-Bigot, J, Biscay, R. J.loubes, J. M. & Muniz-Alvarez, L. (2011) .Group lasso estimation of high-dimensional covariance matrices. Journal of Machine learning Research, 12(Nov), 3187-3225
- 10-Chen, Y. , Wiesel, A. Eldar, Y. C. , & Hero, A. O. (2010). Shrinkage algorithms for MMSE covariance estimation. IEEE Transactions on signal processing, 58  
(10),5016-5029
- 11- Iman, T. M. (2014). Standard analytical study of family consumption of electricity - Study of the case of Sonelgaz Unit Al Buirra - during the period 2008:1 - 2013:12 (Vol. 18). Ministry of

Higher Education and Scientific Research University of Akli Mahnood Olhad/Faculty of  
Economic and Commercial Sciences and Management Sciences.