

## 1 Brief introduction

While the effect of outliers on the parameter estimates of a normal linear regression has been studied, it is of interest to investigate the effect of outliers and contamination of the error terms on the variability of the parameter estimates. Should outliers and contamination have significant effects on the variances of the estimates, then the values of often used test statistics such as  $t$ - and Wald-statistics, the corresponding significance levels of the estimates, and ultimate conclusions drawn based on these outcomes can change drastically. Hence, circumventing such problems is of interest, since it is crucial for actors such as policymakers, organizations and scientists to minimize the probability of drawing incorrect conclusions and/or making undesirable policy decisions. With this paper we aim to investigate the effects on the variability of parameter estimates through inclusion of outliers and contamination of the error terms by heavy-tailed distribution by means of a simulation experiment.

## 2 Data

We implement a simulation study in which we estimate our models on  $n = 1000$  data points in each replication. We randomly generate a set of explanatory variables  $(x_1, x_2, x_3)$  from a zero-mean multivariate normal distribution with  $cov(x_i, x_j) = (1/2)^{|i-j|}$  for  $i = 1, 2, 3$ . Furthermore, we generate the linearly dependent variable  $y$  from a normal distribution with mean  $x'\beta = x_1 + 0.5x_2 + 0.25x_3$  and covariance matrix  $\Sigma_{ij} = cov(x_i, x_j)$ , with the aforementioned covariances. Subsequently, we consider two types of distortions of the simulated data and we examine their effects on the variability of the parameter estimates  $(\beta_1, \beta_2, \beta_3)$ . These two types concern outliers, and contamination of the error terms by a heavy-tailed distribution.

### Contamination with bad leverage points

First, we introduce contamination with outliers in the form of bad leverage points. That is, with probability  $\epsilon = 0.05$  original data points of the above described uncontaminated data set  $(y, x_1, x_2, x_3) := (y - 2, x_1 + 2, x_2 + 2, x_3 + 2)$ .

### Contamination with t-distribution

Next, we introduce contamination of the error terms by a heavy-tailed distribution in the form of a  $t(3)$ -distribution. Error terms are contaminated with probability  $\epsilon = 0.05$ .

## 3 Methodology

We make use of a set of estimators to study the effect of outliers on the variability of the parameter estimates in the case of a normal linear regression. As a benchmark model we apply the Ordinary Least Squares (OLS) regression,

which is expected to produce strongly inflated estimates of the parameters' variability in case of contamination. Subsequently, we consider the Mallows-type M-estimator and the MM-estimator, which are established in order to provide more robust estimates. We investigate the performance of our estimators by replicating the analysis  $R = 500$  times.

### 3.1 Mallows-type M-estimator

The Mallows-type M-estimator is obtained by solving the following equation:

$$\int \Psi(z; T) dF = 0 \quad (1)$$

where  $\Psi(z; T) = \psi_H(y - x'\beta)w(x)$  is the score function with  $w(x)$  representing the weighting function based on the Mahalanobis distance which bounds the influence of leverage outliers, while  $\psi$ , which is the Huber function, takes care of the vertical outliers.

### 3.2 MM-estimator

The MM-estimating comprises of the following two steps:

1. S-estimation In this step the residual scale estimate  $\hat{\sigma}_S$  is retained where  $\hat{\beta}_S = \arg \min_b \min \hat{\sigma}_M^2(b)$
2. M-estimation In this step the residual scale  $\hat{\sigma}_S$  is fixed and is the obtained by solving  $\frac{1}{n} \sum_{i=1}^n \rho(\frac{x_i}{\hat{\sigma}_S}) = \delta$  where  $\delta$  is equal to the expected value of the loss function.

## 4 Results

### 4.1 Uncontaminated data

Figure 1 shows four boxplots containing all estimators on the uncontaminated data belonging to  $\beta_0, \beta_1, \beta_2, \beta_3$ , respectively. In Figure 1 through 3, 1 expresses the OLS estimator, 2 expresses the Mallows-type M-estimator and 3 expresses the MM-estimator. It can be concluded that the MM estimator produces the most accurate estimates of the parameters. In particular, the variability of the  $\beta_2$  and  $\beta_3$  estimates is visibly smaller for MM than for OLS and Mallows-type M. Moreover, the boxplots for OLS and Mallows-type M are of similar scale.

### 4.2 Contaminated data with t-distribution

Figure 2 shows four boxplots containing all estimators on the contaminated data with t-distribution on the error terms belonging to  $\beta_0, \beta_1, \beta_2, \beta_3$ , respectively. It can be concluded that all estimators consistently estimate the parameters despite the contamination. Furthermore, it is observed that, especially for  $\beta_1, \beta_2$

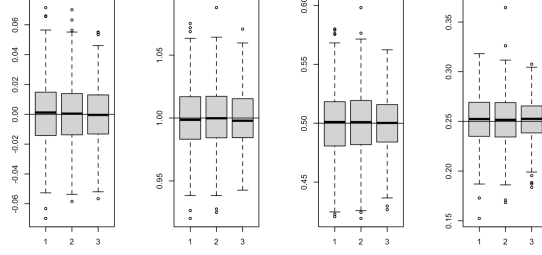


Figure 1: Boxplots for beta parameters for the uncontaminated data

and  $\beta_3$ , the estimated variances of the parameters clearly increase for Mallows-type M compared to OLS. An opposite conclusion can be drawn regarding MM, which produces more ambiguous results: whereas the boxplots for  $\beta_0$  and  $\beta_2$  indicate more precise estimates of MM than of OLS, the boxplots of  $\beta_1$  and  $\beta_3$  imply slightly less precise estimates.

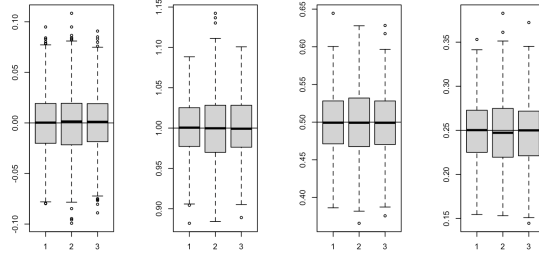


Figure 2: Boxplots for beta parameters for the contaminated data with t-distribution.

### 4.3 Contaminated data with bad leverage points

Figure 3 displays four boxplots containing the aforementioned estimators on the contaminated data with bad leverage points belonging to  $\beta_0, \beta_1, \beta_2, \beta_3$ , respectively. It is evident that the OLS estimator fails to correctly estimate the parameters, since the medians of the OLS estimates are considerably distant from the true values of the parameters. It is also observable that the Mallows-type M-estimator is more accurate in estimating the parameters, while the MM-estimator is the most precise in estimating the parameters. It can also be observed that the MM-estimates have the smallest variance, followed by Mallows-type M, and ultimately OLS.

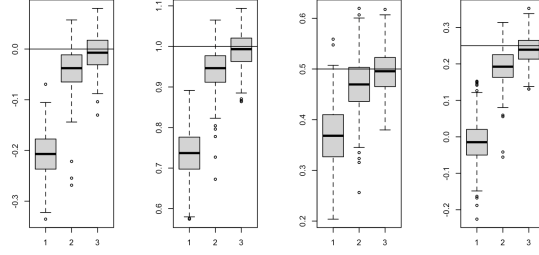


Figure 3: Boxplots for beta parameters for the contaminated data with bad leverage points.

#### 4.4 Uncontaminated versus Contaminated

Figures 4 through 6 in the Appendix compare the variability of the estimated parameters for each of the estimation methods by displaying boxplots for results based on uncontaminated data and contaminated data where the error terms are contaminated with the  $t$ -distribution. Likewise, Figures 7 through 9 show these results but based on contaminated data with bad leverage points. It is clear that across both types of contamination, all estimation methods, and all parameters the variability of the parameters increases when the data set is contaminated with respect to when the data set stays uncontaminated. Moreover, contamination with the  $t$ -distribution does not influence the consistent estimation of all three estimation methods, whereas contamination with bad leverage points leads to highly significantly inconsistent estimates of OLS (Figure 7) and Mallows-type M (Figure 8).

Table 1 summarizes the averages of the estimated coefficient, their standard errors and significance levels based on the null hypothesis that the estimated parameters are equal to the true parameters. From Table 1 we can induce a slightly different conclusion: the OLS coefficients are significantly different from their true values, whereas Mallows-type M still produces unbiased results. However, with respect to the uncontaminated data,  $t$ -statistics increased a bit albeit without any change on  $p$ -values and hence conclusions about significance of the values.

	Uncontaminated				Contamination with t-distribution				Contamination with bad leverage			
	$\beta_0$	$\beta_1$	$\beta_2$	$\beta_3$	$\beta_0$	$\beta_1$	$\beta_2$	$\beta_3$	$\beta_0$	$\beta_1$	$\beta_2$	$\beta_3$
OLS	0.002 (0.021)	1.000 (0.025)	0.501 (0.028)	0.249 (0.025)	-0.000 (0.038)	1.000 (0.043)	0.498 (0.049)	0.251 (0.043)	-0.208 (0.046)	0.733*** (0.052)	0.364** (0.059)	-0.010*** (0.052)
Mallows-type M	0.002 (0.032)	1.000 (0.038)	0.501 (0.042)	0.239 (0.038)	0.000 (0.035)	1.000 (0.043)	0.497 (0.048)	0.251 (0.043)	-0.041 (0.035)	0.938 (0.042)	0.469 (0.048)	0.193 (0.042)
MM	0.001 (0.019)	1.000 (0.022)	0.500 (0.024)	0.250 (0.022)	0.000 (0.038)	1.001 (0.038)	0.496 (0.042)	0.251 (0.038)	-0.007 (0.034)	0.987 (0.040)	0.494 (0.044)	0.242 (0.040)

Table 1: Simulation results for both uncontaminated and contaminated data sets, each of the three estimation methods, and all parameters.

## 5 Appendix

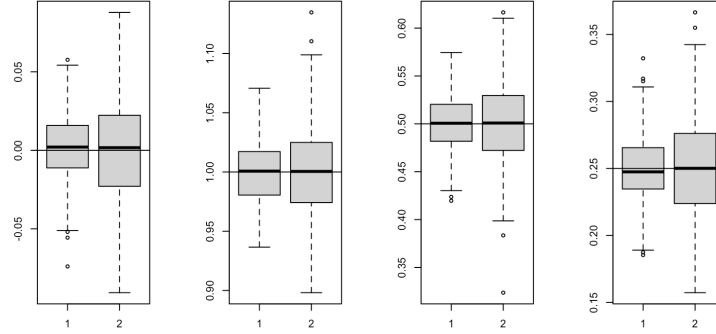


Figure 4: Boxplots for beta parameters for OLS to compare uncontaminated versus contaminated with t-distribution.

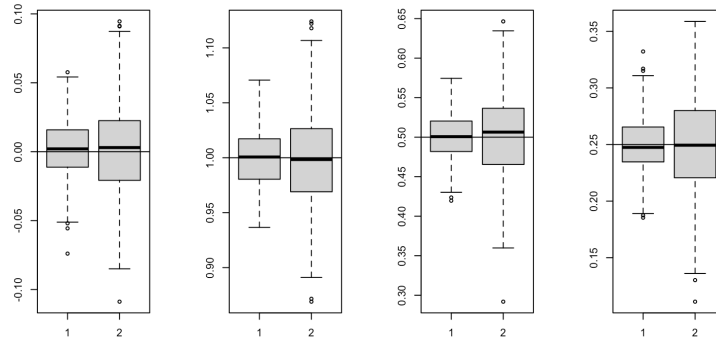


Figure 5: Boxplots for beta parameters for Mallows-type M to compare uncontaminated versus contaminated with t-distribution.

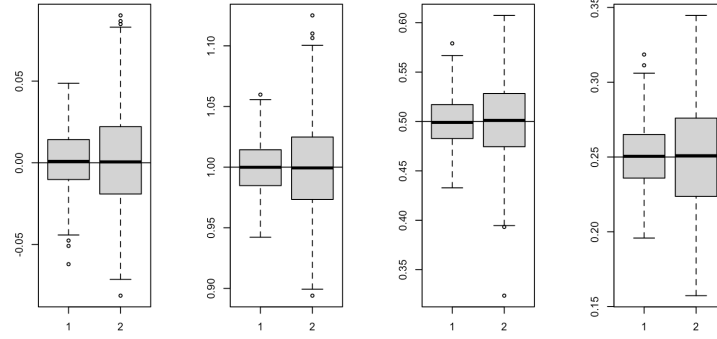


Figure 6: Boxplots for beta parameters for MM to compare uncontaminated versus contaminated with t-distribution.

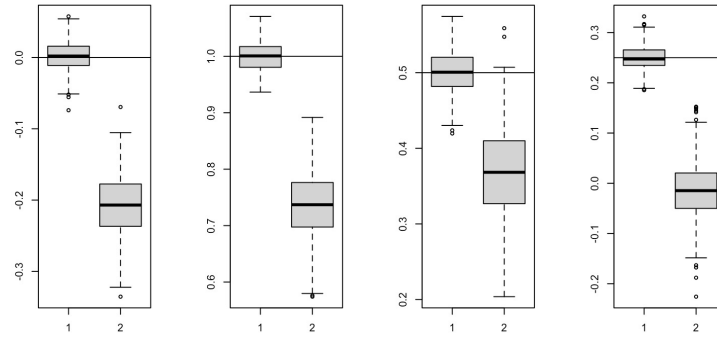


Figure 7: Boxplots for beta parameters for OLS to compare uncontaminated versus contaminated with bad leverage points.

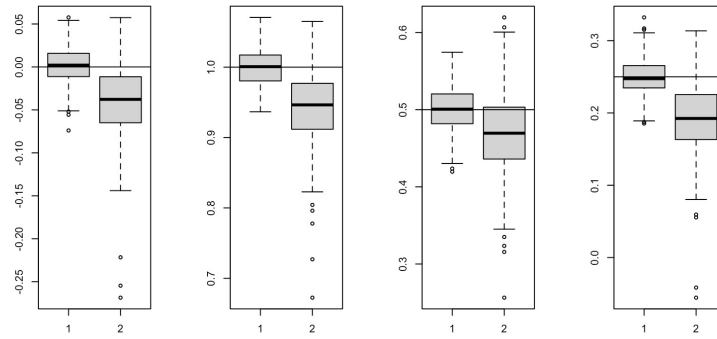


Figure 8: Boxplots for beta parameters for Mallows-type M to compare uncontaminated versus contaminated with bad leverage points.

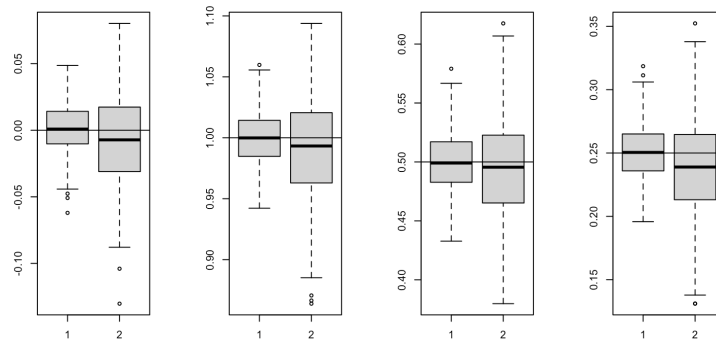


Figure 9: Boxplots for beta parameters for MM to compare uncontaminated versus contaminated with bad leverage points.

```

1 install.packages("mvtnorm")
2 library("mvtnorm")
3 install.packages("robustbase")
4 library("robustbase")
5 install.packages("MASS")
6 library("MASS")
7
8 # Control Parameters
9 set.seed(123)
10 R <- 500
11 n <- 1000
12 mu <- rep(0,3)
13 sgmX <- matrix(c(1,0.5, 0.25, 0.5, 1, 0.5, 0.25, 0.5, 1),3,3)
14 bt <- c(1, 0.5, 0.25)
15
16 # Uncontaminated data
17 resultsUncontaminated <- replicate(R, {
18   # Generate data
19   x <- rmvnorm(n, mean = mu, sigma = sgmX)
20   LP <- x%*%bt
21   y <- rnorm(n, LP, sgmX)
22
23   # Ordinary Least Squares (fit, param., var.)
24   fit.ols <- lm(y~x)
25   coef.ols <- coef(lm(y~x))
26   var.ols <- coefficients(summary(fit.ols))[2,]
27   tstat.ols <- (coef.ols-c(0,bt))/var.ols
28   sign.ols <- coefficients(summary(fit.ols))[4,]
29
30   # Mallows-type Estimator (fit, param., var.)
31   fit.mallows_est <- glmrob(y~x, family = gaussian(link = identity), method = "Mqle", weights.on.x = "hat")
32   coef.mallows_est <- coefficients(summary(fit.mallows_est))[1,]
33   var.mallows_est <- coefficients(summary(fit.mallows_est))[2,]
34   tstat.mallows_est <- (coef.mallows_est - c(0,bt))/var.mallows_est
35   sign.mallows_est <- coefficients(summary(fit.mallows_est))[4,]

```

Figure 10: R code part 1.

```

37 # MM-estimator
38 fit.mm_est <- lmrob(y~x)
39 coef.mm_est <- coefficients(summary(fit.mm_est))[,1]
40 var.mm_est <- coefficients(summary(fit.mm_est))[,2]
41 tstat.mm_est <- (coef.mm_est - c(0,bt))/var.mm_est
42 sign.mm_est <- coefficients(summary(fit.mm_est))[,4]
43
44 c(coef.ols, var.ols, tstat.ols, sign.ols, coef.mallows_est, var.mallows_est, tstat.mallows_est, sign.mallows_est, coef.mm_est, var.mm_est, tstat.mm_est, sign.m
45 })
46 estim_var <- rowMeans(resultsUncontaminated)
47 par(mfrow = c(1, 4))
48 boxplot(resultsUncontaminated[,1], resultsUncontaminated[,17], resultsUncontaminated[,33,]); abline(h = 0)
49 boxplot(resultsUncontaminated[,2], resultsUncontaminated[,18], resultsUncontaminated[,34,]); abline(h = bt[1])
50 boxplot(resultsUncontaminated[,3], resultsUncontaminated[,19], resultsUncontaminated[,35,]); abline(h = bt[2])
51 boxplot(resultsUncontaminated[,4], resultsUncontaminated[,20], resultsUncontaminated[,36,]); abline(h = bt[3])
52

```

Figure 11: R code part 2.

```

53 # Contaminated data
54 resultsContaminated <- replicate(R, {
55   # Generate data
56   x <- rmvnorm(n, mean = mu, sigma = sgmX)
57   LP <- x[,"%bt]
58   ei <- rnorm(n,0,1)
59   eps <- rbinom(n, 1, 0.05) # Contamination with probability 0.05
60
61   # OPTION 1: CONTAMINATE DATA WITH T-DISTRIBUTION ERROR
62   eps <- rbinom(n, 1, 0.05) # Contamination with probability 0.05
63   ei <- (1-eps)*ei + eps*rt(n,3) # Simulate contamination from t distribution with 3 d.o.f.
64   y_contam <- LP + ei
65   x_contam <- x
66
67   # OPTION 2: OUTLIERS (BAD LEVERAGE POINTS)
68   eps <- rbinom(n, 1, 0.05) # Contamination with probability 0.05
69   y <- LP + ei
70   y_contam <- (1-eps)*y + eps*(y-2)
71   x_contam <- (1-eps)*x + eps*(x+2)
72
73   # Ordinary Least Squares (fit, param., var.)
74   fit.ols <- lm(y_contam~x_contam)
75   coef.ols <- coef(lm(y_contam~x_contam))
76   var.ols <- coefficients(summary(fit.ols))[,2]
77   tstat.ols <- (coef.ols-c(0,bt))/var.ols
78   sign.ols <- coefficients(summary(fit.ols))[,4]
79
80   # Mallows-type Estimator (fit, param., var.)
81   fit.mallows_est <- glmrob(y_contam~x_contam, family = gaussian(link = identity), method = "Mqle", weights.on.x = "covMcd")
82   coef.mallows_est <- coefficients(summary(fit.mallows_est))[,1]
83   var.mallows_est <- coefficients(summary(fit.mallows_est))[,2]
84   tstat.mallows_est <- (coef.mallows_est - c(0,bt))/var.mallows_est
85   sign.mallows_est <- coefficients(summary(fit.mallows_est))[,4]

```

Figure 12: R code part 3.

```

87 # MM-estimator
88 fit.mm_est <- lmrob(y_contam~x_contam)
89 coef.mm_est <- coefficients(summary(fit.mm_est))[,1]
90 var.mm_est <- coefficients(summary(fit.mm_est))[,2]
91 tstat.mm_est <- (coef.mm_est - c(0,bt))/var.mm_est
92 sign.mm_est <- coefficients(summary(fit.mm_est))[,4]
93
94 c(coef.ols, var.ols, tstat.ols, sign.ols, coef.mallows_est, var.mallows_est, tstat.mallows_est, sign.mallows_est, coef.mm_est, var.mm_est, tstat.mm_est, sign.m
95 })
96 estim_var <- rowMeans(resultsContaminated)
97 par(mfrow = c(1, 4))
98 boxplot(resultsContaminated[,1], resultsContaminated[,17], resultsContaminated[,33,]); abline(h = 0)
99 boxplot(resultsContaminated[,2], resultsContaminated[,18], resultsContaminated[,34,]); abline(h = bt[1])
100 boxplot(resultsContaminated[,3], resultsContaminated[,19], resultsContaminated[,35,]); abline(h = bt[2])
101 boxplot(resultsContaminated[,4], resultsContaminated[,20], resultsContaminated[,36,]); abline(h = bt[3])

```

Figure 13: R code part 4.