

Chapter 8

Monte Carlo

Introduction

The case studies of Chapter 7 offer a comparison of the competing regression methods by viewing their respective performance in the analysis of datasets that have been well documented and studied in the literature. While exceptional performance of an estimator regarding a specific dataset is, of course, a desirable result, it is more substantial to demonstrate comparable or superior performance across a multitude of datasets derived under identical conditions and, furthermore, under a wide variety of conditions. To this end, Monte Carlo simulation becomes the focal point of this chapter (another such study is found in You (1999)). By viewing several different underlying data structures, insight into the general strengths and weaknesses of the competing methods is offered. The regression methods under comparison are

- OLS,
- LTS (computed using 2000 random subsets),
- M1S (using LTS as the initial estimator),
- S1S (using LTS as the initial estimator),
- BI (using OLS as an initial estimator for a Huber ψ -function BI estimator and, finally, a bisquare ψ -function BI estimator (scale is not iterated in either IRLS procedure in order to alleviate convergence issues)) and
- CBI (reported using both the original v^2 and revised v_w^2 scale estimates).

Monte Carlo simulation results will be discussed relative to issues of (1) unbiasedness of the coefficient estimates, (2) unbiasedness of the scale estimate, (3) expected standard errors for

the coefficient estimates and (4) high breakdown capability via coefficient stability. Certain notation is now introduced in order that the simulation results may be presented in a compact, interpretable form. First, the $p \times 1$ vector $\hat{E}[\hat{\boldsymbol{\beta}}]$ represents the *estimated* expected value of the regression estimator $\hat{\boldsymbol{\beta}}$. Here, $\hat{\boldsymbol{\beta}}$ may represent $\hat{\boldsymbol{\beta}}_{OLS}$, $\hat{\boldsymbol{\beta}}_{LTS}$, $\hat{\boldsymbol{\beta}}_{MIS}$, $\hat{\boldsymbol{\beta}}_{SIS}$, $\hat{\boldsymbol{\beta}}_{BI}$ or $\hat{\boldsymbol{\beta}}_{CBI}$. Second, the scalar $\hat{E}[\hat{\sigma}^2]$ represents the *estimated* expected value of the scale variance estimator $\hat{\sigma}^2$. While $\hat{\sigma}^2$ represents the mean square error in OLS, $\hat{\sigma}^2$ would also represent $\hat{\sigma}_{LTS}^2$ (i.e. $\hat{\sigma}_0^2$ in Chapter 4) or the squared MAD scale estimate in either BI or CBI analyses. Third, the scalar $\hat{E}[v^2]$ represents the *estimated* expected value of the scale variance estimator v^2 , as defined previously for use with either BI or CBI regression analyses. Fourth, $\hat{E}[v_w^2]$ represents the *estimated* expected value of the revised scale variance estimator v_w^2 , computed near the conclusion of a CBI regression analysis. Finally, the $p \times 1$ vector $\hat{E}[se[\hat{\boldsymbol{\beta}}]]$ represents the *estimated* expected value of the standard error for each of the regression coefficient estimates comprising $\hat{\boldsymbol{\beta}}$. The phrase *estimated* has been emphasized to call attention to the fact that due to the finite simulation process of a Monte Carlo study, reported numerical values are susceptible to the inherent sampling error due to random generation.

Each Monte Carlo study has a summary table of results, each with the same structural format as shown in Table 8.1. Each table utilizes the left-most column to provide the expression labels for those associated rows. Six column headers are found in the top row of each table and refer to one of the six competing regression methods. However, make note that the right-most partial column is also associated with the CBI estimator. Table entries may be viewed as falling into one of four zones.

The first zone corresponds to the expression $\hat{E}[\hat{\boldsymbol{\beta}}]$ and has p rows associated with it, one for each regression parameter (in standard order). These entries address coefficient unbiasedness.

Table 8.1: Format of a Monte Carlo simulation summary table. An “x” represents a numerical entry. Shaded cells are not applicable. The following is representative of a $p = 3$ regressor situation.

	OLS	LTS	MIS	S1S	BI	CBI	
$\hat{E}[\hat{\beta}]$	x	x	x	x	x	x	
	x	x	x	x	x	x	
	x	x	x	x	x	x	
$\hat{E}[\hat{\sigma}^2]$	x	x			x	x	$\hat{E}[v_w^2]$
$\hat{E}[v^2]$					x	x	x
$\hat{E}[se[\hat{\beta}]]$	x		x	x	x	x	x
	x		x	x	x	x	x
	x		x	x	x	x	x
$\hat{\beta}$: Range Minimum Maximum IQR	x	x	x	x	x	x	
	x	x	x	x	x	x	
	x	x	x	x	x	x	
	x	x	x	x	x	x	
	x	x	x	x	x	x	
	x	x	x	x	x	x	
	x	x	x	x	x	x	

The second zone addresses unbiasedness in variance estimation. Two rows, corresponding to $\hat{E}[\hat{\sigma}^2]$ and $\hat{E}[v^2]$ define this zone. Cell entries that are inapplicable are shaded (e.g. MIS and S1S share a common scale estimate with LTS). The right-most partial column has a special form, however, having the header $\hat{E}[v_w^2]$ displayed and relating to the entry found directly beneath it.

The third zone relates to the expression $\hat{E}[se[\hat{\beta}]]$ and has p rows associated with it, one for each regression parameter (in standard order). No standard errors for LTS are computed, so these table cells are not utilized. Viewing these estimated expected standard errors is useful in the comparison of regression estimators since smaller standard errors are an indicator of parameter tests with higher power. Yet this concept must be tempered by the understanding that

a scale estimate of poor quality may distort the magnitude of the standard errors or that a repeated breakdown in coefficient estimation may lead to an excessive false alarm (type I, or α , error) rate.

The fourth zone of the Monte Carlo summary table corresponds to a descriptive summary of the observed regression estimators themselves and constitutes the last p rows of the table. Each cell consists of the range, the minimum, the maximum and the interquartile range (IQR) observed during the course of the simulation. These values are useful in determining the stability of a particular regression estimator and may indicate the effectiveness of the method to handle data contamination (i.e. breakdown point capability), provided that estimator bias is not evident.

Eight Monte Carlo studies are presented in this chapter, providing a variety in the number of regressor variables, sample sizes, outlier structure and level of contamination. A total of 2,500 simulated datasets were randomly generated for each Monte Carlo study. Below is a brief overview of these eight scenarios:

- ❶ $(n, p) = (40, 3)$, with no contamination. This is the baseline for classical regression analysis.
- ❷ $(n, p) = (10, 2)$, having a well-defined trend and a high influence cluster of size 2 (the cluster always possesses a positive residual with respect to the general trend).
- ❸ $(n, p) = (26, 4)$, utilizing the Pendleton-Hocking regressors and contamination structure. Specifically, the original four outlying observations are fixed throughout the study.
- ❹ $(n, p) = (75, 4)$, utilizing the Hawkins, Bradu and Kass regressors and contamination structure. The high influence cluster always possesses a positive residual with respect to the general trend.
- ❺ $(n, p) = (80, 4)$, with 40% contamination that is scattered in the regressor space along a secondary linear trend.

- ⑥ $(n, p) = (80, 4)$, with 40% contamination entirely contained within a high leverage cluster region. The high influence cluster always possesses a positive residual with respect to the general trend.
- ⑦ $(n, p) = (80, 4)$, with 40% contamination entirely contained within a high leverage cluster region. There is a random sign for direction of the outlier cluster with respect to the general trend computed for each simulation replicate.
- ⑧ $(n, p) = (100, 6)$, with 40% contamination in a complex configuration. There are 5% low leverage outliers, 5% spurious high leverage outliers and a 30% high leverage outlier cluster along a secondary linear trend.

The third and fourth Monte Carlo studies each involved literature data, with the provided regressor data being fixed throughout the simulation. However, the other six Monte Carlo studies each had a prescribed regressor space generation rule. Each of these studies was conducted under two distinct and separate simulations by utilizing (1) a fixed regressor space format as well as (2) a random regressor space format. In a fixed format, a regressor space was generated once and subsequently held constant throughout the rest of the simulation. Conversely, under a random format, a new regressor space was generated for each simulation run. As the underlying linear model studied throughout the course of this research assumes a fixed regressor space, a fixed regressor space format was used. However, research indicated potential issues with respect to internal instability of the current high breakdown regression methodologies. Thus, a random regressor format allowed for the investigation of the performance of the competing regression methods under a wide variety of regressor space configurations. As will be demonstrated throughout the chapter, the results for a fixed regressor case and its corresponding random regressor case were generally quite consistent.

For reference purposes, the fixed regressor case format simulations (aside from studies 3 and 4) are denoted by the suffix A, while the random regressor case format simulations are denoted by the suffix B. Each simulation study will now be presented in succession. As an aid

Table 8.2: True parameter values for the Monte Carlo simulation studies.

Sim.	σ^2	β_0	β_1	β_2	β_3	β_4	β_5
1-A, B	16.0	-25.0	10.0	1.0			
2-A, B	9.0	100.0	-5.0				
3	0.5	20.0	3.0	-2.0	0.0		
4	0.5	0.20	-0.15	0	0.10		
5-A, B	1.0	100.0	-10.0	5.0	20.0		
6-A, B	1.0	50.0	5.0	-10.0	1.0		
7-A, B	25.0	50.0	5.0	-10.0	0.0		
8-A, B	1.0	50.0	5.0	-10.0	1.0	0.0	0.0

for interpreting the numerous summary tables, Table 8.2 provides the true parameter values from which the general trend data was generated during each Monte Carlo study. Each study section contains the specific details regarding data generation as well as the discussion and interpretation of the results.

§8.1 Study #1: Uncontaminated Data with 2 Regressors

One of the goals of the research is that the CBI regression estimator performs better than the current high breakdown methods when the classical normal theory assumptions are valid. A secondary goal would be for the CBI regression estimator to approach the performance of M and BI regression under said condition. OLS is the best linear unbiased estimator in this setting (Myers, 1990). In this first Monte Carlo study, the competing regression methods are compared via the analysis of datasets each consisting of $n = 40$ observations in two regressor variables that is quite well-behaved.

The two regressor variables were randomly generated independently via

$$x_{1i} \sim N[\mu_{x_1} = 5, \sigma_{x_1}^2 = 0.0625]$$

and

$$x_{2i} \sim N[\mu_{x_2} = 25, \sigma_{x_2}^2 = 16].$$

The response variable was then generated according to the linear model

$$y_i = -25 + 10x_{1i} + x_{2i} + \varepsilon_i,$$

with random errors following

$$\varepsilon_i \sim N[\mu_\varepsilon = 0, \sigma_\varepsilon^2 = 16].$$

For Monte Carlo study #1A the regressor space was generated just once, with a new response vector generated for each simulation run. This regressor space is presented in Appendix B. For Monte Carlo study #1B the regressor space was generated anew for each simulation run (with a new response vector generated for each simulation run as well).

§8.1.1 Results for Monte Carlo Study #1A (Fixed Regressor Space)

Table 8.3 provides the simulation summary for the fixed regressor space scenario of the first Monte Carlo study. Again, the parameters to be estimated are $\beta' = [-25 \ 10 \ 1]$ and $\sigma_\varepsilon^2 = 16$. A discussion of these results is now presented in four subsections.

Unbiasedness of the Regression Estimator

All entries corresponding to $\hat{E}[\hat{\beta}]$ were, aside from random sampling variability, in agreement with the true underlying model parameters. Therefore, each of the six regression procedures demonstrated the ability to produce an unbiased regression estimator when no data contamination was present.

Scale Estimation

With a true scale of $\sigma_\varepsilon^2 = 16$, the OLS scale estimate was, as expected, unbiased aside from random sampling variability as suggested by $\hat{E}[\hat{\sigma}^2] \approx 16.050$. However, $\hat{\sigma}_{LTS}^2$ exhibited a negative bias (an underestimate) since $\hat{E}[\hat{\sigma}^2] \approx 12.877$. In response units (i.e. taking square roots) it was a difference between $\hat{\sigma} = \sqrt{\hat{E}[\hat{\sigma}^2]} \approx 3.589$ and $\sigma_\varepsilon = 4$. For BI, it was observed

Table 8.3: Simulation results for Monte Carlo study #1A, the fixed regressor case.

	OLS	LTS	MIS	SIS	BI	CBI	
$\hat{E}[\hat{\beta}]$	-24.884	-24.990	-25.109	-25.279	-24.909	-25.272	
	9.981	9.967	10.037	10.048	9.986	10.044	
	0.999	1.008	0.996	1.001	0.999	1.002	
$\hat{E}[\hat{\sigma}^2]$	16.050	12.877			15.093	13.064	$\hat{E}[v_w^2]$
$\hat{E}[v^2]$					1.291	22.004	15.948
$\hat{E}[se[\hat{\beta}]]$	17.073		21.925	19.737	4.932	23.434	19.830
	3.355		4.270	3.969	0.969	4.605	3.897
	0.144		0.199	0.157	0.042	0.200	0.169
$\hat{\beta}$:	129.038	280.537	226.765	219.432	128.923	205.190	
	-91.118	-151.693	-123.136	-137.899	-90.783	-114.261	
	37.920	128.844	103.628	81.532	38.140	90.929	
	23.023	52.816	29.211	28.845	23.855	34.604	
	24.890	54.768	42.934	42.965	24.735	39.512	
	-2.313	-19.412	-12.423	-11.039	-2.372	-12.027	
	22.577	35.356	30.511	31.926	22.363	27.484	
	4.604	10.195	5.610	5.394	4.613	6.771	
	1.053	2.873	3.144	2.292	1.060	2.332	
	0.455	-0.265	-1.065	-0.268	0.451	-0.453	
1.508	2.608	2.079	2.024	1.511	1.879		
0.195	0.407	0.270	0.252	0.198	0.295		

$\hat{E}[\hat{\sigma}^2] \approx 15.093$, this being closer to the true value than that for the LTS scale estimate. However, from $\hat{E}[v^2] \approx 1.291$ it was clear that the robust v scale estimate for the BI estimator was dramatically small in magnitude. Meanwhile, with $\hat{E}[\hat{\sigma}^2] \approx 13.064$ the CBI scale estimator had a modest improvement over that from LTS. Further, while $\hat{E}[v^2] \approx 22.004$, it was also seen that $\hat{E}[v_w^2] \approx 15.948$. In response units the later becomes $\hat{v}_w = \sqrt{\hat{E}[v_w^2]} \approx 3.993$ compared to the true value of $\sigma_\epsilon = 4$. Thus, the CBI scale estimate v_w^2 was minimally biased, outperformed the other non-OLS scale estimates and was very competitive with the OLS scale estimate.

Standard Errors

The expected standard errors for the BI coefficients were by far the smallest, yet this resulted from a substantial underestimate of scale. Of more interest is that the standard errors

that (using v_w^2) for the CBI coefficients are comparable to those for the S1S coefficients (the standard errors for the M1S coefficients were somewhat higher), even though the M1S and S1S standard errors were based on a dramatic underestimate of scale! OLS, as expected, outperformed S1S and CBI by possessing smaller standard errors.

Coefficient Stability

Given that there was no contamination in the dataset construction, it would be expected that the BI coefficients be extremely competitive regarding coefficient stability, and this was indeed the case. The range for each of its three coefficients was by far the smallest across all non-OLS competitors, as was the IQR. There was near agreement between OLS and BI, the results were quite similar. Regarding the high breakdown estimators, M1S and S1S fared slightly better than CBI with respect to the IQR (LTS had substantially larger IQR values than any other method). However, the CBI coefficients were much more stable in the extremes than were the LTS, M1S or S1S coefficients as the CBI coefficients had smaller ranges (except for $\hat{\beta}_2$ when using S1S). In general, the LTS, M1S and S1S coefficients each had maximums and/or minimums that drifted substantially farther from the true parameter values than did the CBI observed extremes. Dramatic differences in estimator extremes are much more relevant to coefficient stability than are mild differences in the middle half of the corresponding sampling distributions, especially in the context of estimator breakdown in a practical, non-theoretic sense. Accordingly, there was an improvement in coefficient stability when using CBI over LTS, M1S or S1S under the constructs of this Monte Carlo study.

§8.1.2 Results for Monte Carlo Study #1B (Random Regressor Space)

The second case of the first Monte Carlo study enabled the regressor space to be generated anew for each simulation run. The results of this study, Monte Carlo study #1B, are provided in Table 8.4.

Table 8.4: Simulation results for Monte Carlo study #1B, the random regressor case.

	OLS	LTS	MIS	SIS	BI	CBI	
$\hat{E}[\hat{\beta}]$	-25.184	-25.029	-25.298	-25.226	-25.180	-24.867	
	10.027	9.989	10.042	10.030	10.025	9.965	
	1.002	1.004	1.004	1.003	1.003	1.002	
$\hat{E}[\hat{\sigma}^2]$	15.773	12.492			14.653	12.805	$\hat{E}[v_w^2]$
$\hat{E}[v^2]$					1.266	21.480	15.492
$\hat{E}[se[\hat{\beta}]]$	13.657		17.099	16.633	3.981	18.907	15.958
	2.605		3.250	3.160	0.759	3.608	3.045
	0.163		0.205	0.198	0.048	0.226	0.191
$\hat{\beta}$: Range Minimum Maximum IQR	110.258	252.680	295.839	300.659	122.372	183.479	
	-77.214	-157.513	-228.669	-228.669	-89.541	-111.939	
	33.045	95.167	67.170	71.990	32.831	71.541	
	18.638	41.977	22.947	22.567	18.859	28.013	
	19.765	52.406	58.675	58.170	23.532	37.230	
	-0.993	-14.527	-8.552	-8.047	-0.845	-10.743	
	18.772	37.879	50.123	50.123	22.686	26.487	
	3.616	7.872	4.417	4.284	3.642	5.493	
	1.229	2.749	4.266	4.197	1.245	2.296	
	0.406	-0.361	-1.625	-1.641	0.379	-0.275	
1.635	2.388	2.641	2.557	1.624	2.021		
0.220	0.498	0.255	0.255	0.220	0.327		

Unbiasedness of the Regression Estimator

All entries corresponding to $\hat{E}[\hat{\beta}]$ were, aside from random sampling variability, in agreement with the true underlying model parameters. As in the fixed regressor case, each of the six regression procedures demonstrated the ability to produce an unbiased regression estimator.

Scale Estimation

The true scale was $\sigma_\varepsilon^2 = 16$, with OLS having $\hat{E}[\hat{\sigma}^2] \approx 15.773$, a little farther off-target than what was witnessed during the fixed regressor case. As before, $\hat{\sigma}_{LTS}^2$ exhibited a negative bias (an underestimate) since $\hat{E}[\hat{\sigma}^2] \approx 12.492$. In response units it was a difference between $\hat{\sigma} = \sqrt{\hat{E}[\hat{\sigma}^2]} \approx 3.534$ and $\sigma_\varepsilon = 4$. With $\hat{E}[\hat{\sigma}^2] \approx 14.653$ it was observed that BI fared better than LTS regarding a scale estimate. Yet again the robust scale estimate for the BI estimator was

dramatically small in magnitude, as indicated by $\hat{E}[v^2] \approx 1.266$. The CBI scale estimate, with $\hat{E}[\hat{\sigma}^2] \approx 12.805$, saw a modest improvement over that for LTS. While $\hat{E}[v^2] \approx 21.480$ signaled a positive bias, it was also seen that $\hat{E}[v_w^2] \approx 15.492$. The later becomes $\hat{v}_w = \sqrt{\hat{E}[v_w^2]} \approx 3.936$ versus the true value of $\sigma_\varepsilon = 4$, a much improved scale estimate. Thus, the CBI procedure using v_w^2 again outperformed the other non-OLS scale estimates.

Standard Errors

The expected standard errors for the BI coefficients were by far the smallest, yet this resulted from a substantial underestimate of scale. As before, the standard errors (using v_w^2) for the CBI coefficients were slightly smaller than those for either the M1S coefficients or S1S coefficients, again with the later two being based on a dramatic underestimate of scale. OLS, as expected, had lower standard errors than those obtained for M1S, S1S or CBI.

Coefficient Stability

As with Monte Carlo study #1A, it was expected that the BI coefficients would be extremely competitive regarding coefficient stability, and this was indeed the case. The range for each of the three coefficients was by far the smallest across all non-OLS competitors, as was the IQR. Unlike the fixed regressor setting (study #1A), here the BI estimators had larger ranges than those obtained for OLS. Regarding the high breakdown estimators, M1S and S1S fared slightly better than LTS and CBI with respect to the IQR (CBI showed an improvement over LTS). However, the CBI coefficients were much more stable in the extremes than were the LTS, M1S or S1S coefficients; the CBI coefficients had much smaller ranges, whereas the LTS, M1S and S1S coefficients each had maximums and/or minimums that drifted substantially farther from the true parameter values than did the CBI observed extremes. In the context of estimator breakdown, there was an improvement in coefficient stability when using CBI over LTS, M1S or S1S.

Perhaps the most interesting comparison between the results of studies #1A and #1B relates to the observed estimator ranges for M1S and S1S. While the LTS estimator ranges were smaller for the random regressor case, both M1S and S1S demonstrated dramatic increases in the random regressor case from what each observed during the fixed regressor case. In other words, even though the initial estimate (i.e. LTS) was more stable, M1S and S1S became more extreme due to the randomness of the regressors. This is another illustration of the internal instability that can manifest within M1S and S1S analyses.

§8.2 Study #2: Simple Linear Regression with a High Influence Cluster

A simple linear regression example with a cluster of two high influence points has been useful in demonstrating the limitations of the low breakdown methods as well as the internal instability of the high breakdown methods. In this second Monte Carlo study, the competing regression methods are compared via the analysis of datasets each consisting of $n=10$ observations in one regressor variable.

A single regressor variable was randomly generated via

$$x_i \sim \begin{cases} U[3, 7], & i = 1 \text{ or } i = 2, \\ U[7, 11], & i = 3 \text{ or } i = 4, \\ U[15, 19], & i = 5 \text{ or } i = 6, \\ U[19, 23], & i = 7 \text{ or } i = 8, \\ 30, & i = 9, \\ 31, & i = 10. \end{cases}$$

The response variable was then generated according to the linear model

$$y_i = \begin{cases} 100 - 5x_i + \varepsilon_i, & i \leq 8, \\ 80, & i > 8, \end{cases}$$

with random errors following

$$\varepsilon_i \sim N[\mu_\varepsilon = 0, \sigma_\varepsilon^2 = 9], \text{ for } i \leq 8.$$

Under such a construction observations 9 and 10 were always tandem high influence points with the other eight observations located so that the general trend should be easily discernible.

For Monte Carlo study #2A the regressor space was generated just once, with a new response vector generated for each simulation run. This regressor space is presented in Appendix B. For Monte Carlo study #2B the regressor space was generated anew for each simulation run (with a new response vector generated for each simulation run as well).

§8.2.1 Results for Monte Carlo Study #2A (Fixed Regressor Space)

Table 8.5 provides the simulation summary for the fixed regressor case of the second Monte Carlo study. It is noted that MIS and SIS produced identical results for every simulation case. The parameters to be estimated are $\beta' = [100 \quad -5]$ and $\sigma_e^2 = 9$. A discussion of these results is now presented.

Unbiasedness of the Regression Estimator

From $\hat{E}[\hat{\beta}]$ it is clear that OLS was repeatedly pulled towards the tandem high influence cluster and exhibited a large bias in its estimation of β . By contrast, the LTS and CBI entries corresponding to $\hat{E}[\hat{\beta}]$ suggest no bias in the estimation of β . However, both MIS and SIS

Table 8.5: Simulation results for Monte Carlo study #2A, the fixed regressor case.

	OLS	LTS	MIS	SIS	BI	CBI	
$\hat{E}[\hat{\beta}]$	42.316	100.004	89.876	89.876	42.796	100.03	
	0.006	-5.005	-4.133	-4.133	-0.044	-5.003	
$\hat{E}[\hat{\sigma}^2]$	1250.37	20.529			2095.83	13.937	$\hat{E}[v_w^2]$
$\hat{E}[v^2]$					546.381	10.433	8.608
$\hat{E}[se[\hat{\beta}]]$	24.739		7.780	7.780	16.515	2.857	2.592
	1.313		0.650	0.650	0.882	0.196	0.178
$\hat{\beta}$:	13.132	34.549	79.404	79.404	13.183	31.803	
	35.377	82.450	25.591	25.591	3.580	79.23	
Range	48.509	117	104.995	104.995	48.984	111.04	
Minimum	2.557	5.314	6.456	6.456	2.582	3.716	
Maximum	0.523	2.180	4.937	4.937	0.547	1.878	
IQR	-0.261	-6.031	-5.284	-5.284	-0.330	-5.687	
	0.263	-3.851	-0.346	-0.346	0.217	-3.809	
	0.105	0.375	0.523	0.523	0.107	0.264	

exhibited a bias in both coefficients under this dataset construction. This surprising result indicates that the high influence cluster did indeed play a significant role in the final regression estimator (recall, the initial estimator, LTS, was unbiased). The bias, then, stems from the fact that the high influence cluster was always above the general trend (i.e. the shift was always positive) so that the influence did not average out in order to yield the unbiasedness property. Additionally, BI, like OLS, was completely overwhelmed by the high influence cluster and was severely biased as well.

Scale Estimation

While the true scale was $\sigma_\varepsilon^2 = 9$, the breakdown of OLS was clearly exhibited from $\hat{E}[\hat{\sigma}^2] \approx 1250.37$. For LTS, $\hat{E}[\hat{\sigma}^2] \approx 20.529$ signaled that $\hat{\sigma}_{LTS}^2$ was positively biased, an overestimate. Or, a difference between $\hat{\sigma} = \sqrt{\hat{E}[\hat{\sigma}^2]} \approx 4.531$ and $\sigma_\varepsilon = 3$. For BI, $\hat{E}[\hat{\sigma}^2] \approx 2095.83$ clearly indicated the breakdown of BI regression in the presence of a high influence cluster of size two. Even the robust scale estimate for the BI estimator was dramatically large in magnitude, as evident by $\hat{E}[\hat{\sigma}^2] \approx 546.381$. The CBI scale estimator also had a positive bias as $\hat{E}[\hat{\sigma}^2] \approx 13.937$, still an improvement over that from LTS. Improvements were seen by $\hat{E}[\hat{v}^2] \approx 10.433$ and by $\hat{E}[\hat{v}_w^2] \approx 8.608$, the later becoming $\hat{v}_w = \sqrt{\hat{E}[\hat{v}_w^2]} \approx 2.933$ (versus $\sigma_\varepsilon = 3$) in response units. The CBI procedure using \hat{v}_w^2 outperformed all other scale estimates.

Standard Errors

The expected standard errors for OLS and BI were by far the largest due to the breakdown of these procedures and their enormous scale estimates. The standard errors (using \hat{v}_w^2) for the CBI coefficients were dramatically smaller than those for either the M1S coefficients or S1S coefficients.

Coefficient Stability

The complete breakdowns of OLS and BI made the tight observed estimate ranges inconsequential. It actually reflects the complete dominance of the two high influence points in their respective analyses. The CBI coefficients were much more stable than either the LTS, M1S or S1S coefficients; each CBI coefficient's range was considerably the smallest of the group. In addition, each CBI coefficient's IQR was much smaller as well. Thus, there was a clear improvement in coefficient stability when using CBI over LTS, M1S or S1S.

§8.2.2 Results for Monte Carlo Study #2B (Random Regressor Space)

The second case of the second Monte Carlo study enabled the regressor space to be generated anew for each simulation run. The results of this study, Monte Carlo study #2B, are provided in Table 8.6.

Unbiasedness of the Regression Estimator

As was the case during the fixed regressor case of this study, the OLS entries corresponding to

Table 8.6: Simulation results for Monte Carlo study #2B, the random regressor case.

	OLS	LTS	M1S	S1S	BI	CBI	
$\hat{E}[\hat{\beta}]$	51.678	99.910	92.572	92.572	52.305	99.988	
	-0.463	-4.990	-4.304	-4.304	-0.524	-4.995	
$\hat{E}[\hat{\sigma}^2]$	1412.58	20.378			2283.20	13.812	$\hat{E}[v_w^2]$
$\hat{E}[v^2]$					587.435	10.207	8.430
$\hat{E}[se[\hat{\beta}]]$	24.678		5.834	5.834	16.043	2.496	2.266
	1.311		0.522	0.522	0.859	0.172	0.156
$\hat{\beta}$:	33.199	53.749	72.067	72.067	33.128	22.142	
	32.090	61.993	37.253	37.253	32.715	89.368	
	65.289	115.742	109.320	109.320	65.844	111.510	
	6.320	4.747	5.141	5.141	6.246	3.297	
	1.746	6.710	6.063	6.063	1.744	1.824	
Minimum	-1.176	-6.119	-5.570	-5.570	-1.231	-5.725	
Maximum	0.570	0.591	0.493	0.493	0.513	-3.901	
IQR	0.342	0.329	0.415	0.415	0.336	0.223	

$\hat{E}[\hat{\boldsymbol{\beta}}]$ indicate that the regression was repeatedly pulled towards the tandem high influence cluster and exhibited bias in its estimation of $\boldsymbol{\beta}$. Again, LTS and CBI exhibited unbiasedness. M1S and S1S, as in the previous section, each demonstrated a clear bias in both coefficients, with marginally better results observed when the regressor was randomly generated than in the fixed regressor case. BI, like OLS, was again completely overwhelmed by the high influence cluster and was severely biased.

Scale Estimation

The scale results mirror those obtained during the fixed regressor case. Although $\sigma_\varepsilon^2 = 9$, the breakdown of OLS was clearly exhibited from $\hat{E}[\hat{\sigma}^2] \approx 1412.58$. LTS was again positively biased since $\hat{E}[\hat{\sigma}^2] \approx 20.378$, this representing a difference between $\hat{\sigma} = \sqrt{\hat{E}[\hat{\sigma}^2]} \approx 4.514$ and $\sigma_\varepsilon = 3$. $\hat{E}[\hat{\sigma}^2] \approx 2283.20$ clearly indicated the breakdown of BI regression in the presence of a high influence cluster of size two, with the robust scale estimate also being dramatically large in magnitude since $\hat{E}[v^2] \approx 587.43$. With $\hat{E}[\hat{\sigma}^2] \approx 13.812$, the CBI scale estimator had a positive bias, but still an improvement over that from LTS. While $\hat{E}[v^2] \approx 10.207$, it was also seen that $\hat{E}[v_w^2] \approx 8.430$. In response units the later becomes $\hat{v}_w = \sqrt{\hat{E}[v_w^2]} \approx 2.903$ (versus $\sigma_\varepsilon = 3$). As such, the CBI procedure using v_w^2 outperformed the other scale estimates.

Standard Errors

As in the fixed regressor case, the expected standard errors for OLS and BI were by far the largest due to the breakdown of these procedures and their enormous scale estimates. Again, the standard errors (using v_w^2) for the CBI coefficients were dramatically smaller than those for either the M1S coefficients or S1S coefficients.

Coefficient Stability

An interesting result was that the CBI had an observed range for the intercept that was much lower than those obtained by OLS or BI, with the three ranges (OLS, BI and CBI) being quite similar for the slope. Further, the CBI had the smallest IQR values for both estimators as well. Of course, while the complete breakdown of OLS and BI makes their observed ranges inconsequential, in a random regressor setting these two methods were both biased and more variable than was the CBI method. As in the fixed regressor setting, the CBI coefficients were much more stable than either the LTS, M1S or S1S coefficients. Each CBI coefficient's range was considerably the smallest of the group, with each CBI coefficient's IQR being much smaller as well. Overall, there was a clear improvement in coefficient stability when using CBI over the five competitive methods.

§8.3 Study #3: Pendleton-Hocking Data

As shown previously in Chapter 7, the Pendleton-Hocking dataset illustrates the effectiveness of BI regression when data contamination is isolated; four of the twenty-six observations are problematic (due to their leverage and/or outlier nature), but on their own accord (no sizable joint influence to deal with). In this third Monte Carlo study each simulated dataset utilized the original regressor values of the Pendleton-Hocking dataset. Furthermore, the original values for y_{11} , y_{17} , y_{18} and y_{24} were also maintained throughout all simulation runs. Specifically, the $n = 26$ observations were generated as

$$y_i = \begin{cases} 20 + 3x_{1i} - 2x_{2i} + \varepsilon_i, & i \notin (11, 17, 18, 24), \\ 59.289, & i = 11, \\ 58.699, & i = 17, \\ 50.086, & i = 18, \\ 56.741, & i = 24, \end{cases},$$

where for $i \notin (11, 17, 18, 24)$ the random errors follow

$$\varepsilon_i \sim N[\mu_\varepsilon = 0, \sigma_\varepsilon^2 = 0.25].$$

The results of this Monte Carlo study are provided in Table 8.7. Here, the parameters to be estimated are $\beta' = [20 \ 3 \ -2 \ 0]$ and $\sigma_\epsilon^2 = 0.25$. A discussion of these results is now presented.

Unbiasedness of the Regression Estimator

The four observations y_{11} , y_{17} , y_{18} and y_{24} were not generated via a linear model with random errors symmetric about zero. One potential consequence of this fact would be an introduction of bias to the regression estimators. From $\hat{E}[\hat{\beta}]$ it was observed that all six regression estimators were biased, albeit to differing degrees. OLS demonstrated the most drift

Table 8.7: Simulation results for Monte Carlo study #3.

	OLS	LTS	MIS	SIS	BI	CBI	
$\hat{E}[\hat{\beta}]$	8.583	15.450	16.058	18.169	11.702	19.126	
	3.535	3.230	3.202	3.095	3.432	3.045	
	-1.639	-1.893	-1.908	-1.957	-1.815	-1.981	
	0.338	0.148	0.128	0.060	0.271	0.028	
$\hat{E}[\hat{\sigma}^2]$	3.256	0.227			0.360	0.274	$\hat{E}[v_w^2]$
$\hat{E}[v^2]$					0.108	0.273	0.221
$\hat{E}[se[\hat{\beta}]]$	6.317		14.436	18.182	8.588	10.729	9.613
	0.363		0.734	0.924	0.436	0.549	0.492
	0.160		0.341	0.429	0.202	0.255	0.228
	0.179		0.469	0.590	0.278	0.345	0.309
$\hat{\beta}$: Range Minimum Maximum IQR	6.206	184.498	136.294	400.787	106.519	150.366	
	5.235	-68.664	-63.428	-163.884	-46.675	-56.170	
	11.441	115.835	72.866	236.902	59.844	94.196	
	1.313	22.413	16.981	21.977	17.518	22.003	
	0.436	9.458	6.835	20.147	5.423	7.632	
	3.338	-1.897	0.363	-7.833	0.969	-0.737	
	3.775	7.560	7.199	12.314	6.392	6.894	
	0.089	1.163	0.859	1.116	0.887	1.109	
	0.193	4.323	3.308	9.650	2.506	3.477	
	-1.739	-4.326	-3.277	-7.270	-2.934	-3.777	
	-1.546	-0.003	0.031	2.380	-0.428	-0.299	
	0.036	0.533	0.400	0.521	0.414	0.519	
0.124	6.047	4.657	13.525	3.421	4.998		
0.279	-3.168	-1.777	-7.429	-1.248	-2.500		
0.403	2.879	2.880	6.096	2.173	2.497		
0.021	0.710	0.549	0.708	0.562	0.712		

away from β , with BI being only marginally better. While S1S outperformed LTS and M1S, CBI was clearly the least biased of the group.

Scale Estimation

While the true scale was $\sigma_\varepsilon^2 = 0.25$, the MSE for OLS exhibited a large, positive bias since $\hat{E}[\hat{\sigma}^2] \approx 3.256$. This dramatic increase in scale was another indication of the breakdown of OLS. With $\hat{E}[\hat{\sigma}^2] \approx 0.227$, it was observed that $\hat{\sigma}_{LTS}^2$ was a slight underestimate of scale, on average. In response units it was a difference between $\hat{\sigma} = \sqrt{\hat{E}[\hat{\sigma}^2]} \approx 0.476$ and $\sigma_\varepsilon = 0.5$. The scale estimate from BI was an overestimate, given $\hat{E}[\hat{\sigma}^2] \approx 0.360$, while $\hat{E}[v^2] \approx 0.108$ signifies that the robust BI scale estimate was a dramatic underestimate. Another overestimate was witnessed by the CBI scale estimator since $\hat{E}[\hat{\sigma}^2] \approx 0.274$, which was virtually the same as $\hat{E}[v^2] \approx 0.273$. However, $\hat{E}[v_w^2] \approx 0.221$ denotes a modest underestimate, where $\hat{v}_w = \sqrt{\hat{E}[v_w^2]} \approx 0.470$ (versus $\sigma_\varepsilon = 0.5$) still indicates a decent estimate of scale. The CBI scale estimate (using v_w^2) was competitive with the LTS scale estimate, with both of these scale estimates outperforming the BI scale estimate regarding bias.

Standard Errors

The breakdown of OLS and BI in regression estimation diminishes the importance of these two methods having obtained the smallest the expected standard errors. The standard errors (using v_w^2) for the CBI coefficients were dramatically smaller than those for either the M1S coefficients or S1S coefficients, even though the CBI scale estimate was only marginally smaller than the LTS scale estimate (in response units). This comparison was true even when considering the CBI standard errors based on v^2 , a positively biased estimator.

Coefficient Stability

The OLS coefficients were by far the most stable. OLS was also the most severely biased and that must be considered in concert with this evaluation. Here, the true value for any one of the four parameters was not contained in the corresponding observed ranges for the OLS coefficients! BI also demonstrated better coefficient stability than the collection of high breakdown estimators, but also fell prey to being severely biased. M1S was the most stable of the high breakdown estimators, clearly improving over LTS. Yet M1S also has bias issues with the first two coefficients, the same two coefficients that demonstrated the most pronounced improvement in stability. More interesting is the comparison of S1S to CBI, the two least biased estimators. While the respective four IQR values were nearly the same, the CBI coefficients were much more stable than S1S as each CBI range was considerably smaller. The four S1S ranges were over twice as large as the corresponding four CBI ranges! LTS, in addition to its bias, also had larger ranges than did CBI.

§8.4 Study #4: Hawkins-Bradru-Kass Data

As shown previously in Chapter 7, the HBK dataset illustrates the breakdown of BI regression when data contamination is moderate but clustered, while also demonstrating the effectiveness of the high breakdown procedures. Along with the ten high influence points, a cluster of four good leverage points were contained within this dataset. The difficulty with the original dataset was that these four good leverage points did not lie near the OLS fit of the 65 other good observations; they became, essentially, just a second cluster of high influence points. It is of interest to investigate the performance of the various regression methods if the good leverage cluster were generated under the same underlying model as the bulk of the data. This fourth Monte Carlo study utilized the original regressor values of the HBK dataset, but generated a new response vector while maintaining observations 1 through 10 as a high influence cluster. Specifically, the $n = 75$ observations were generated by the linear model

$$y_i = \begin{cases} \varepsilon_i, & i \leq 10, \\ 0.2 - 0.15 x_{1i} + 0.1 x_{3i} + \varepsilon_i, & i > 10, \end{cases}$$

with random errors following

$$\varepsilon_i \sim \begin{cases} N[\mu_\varepsilon = 0, \sigma_\varepsilon^2 = 0.25], & i \notin (1:10), \\ N[\mu_\varepsilon = 10, \sigma_\varepsilon^2 = 0.385^2], & i \in (1:10). \end{cases}$$

This generation scheme, and the numerical values assigned therein, was based on the structure of the original HBK dataset.

The results of this Monte Carlo are provided in Table 8.8. The parameters to be estimated are $\beta' = [0.20 \ -0.15 \ 0 \ 0.10]$ and $\sigma_\varepsilon^2 = 0.25$. A discussion of these results follows.

Table 8.8: Simulation results for Monte Carlo study #4.

	OLS	LTS	MIS	SIS	BI	CBI	
$\hat{E}[\hat{\beta}]$	0.029	-0.103	0.035	0.196	-0.423	0.190	
	-0.019	-0.121	-0.141	-0.153	-0.093	-0.150	
	-0.307	0.064	0.032	0.001	0.125	0.005	
	0.456	0.187	0.157	0.104	0.292	0.101	
$\hat{E}[\hat{\sigma}^2]$	3.478	0.300			0.485	0.341	$\hat{E}[v_w^2]$
$\hat{E}[v^2]$					0.020	0.436	0.268
$\hat{E}[se[\hat{\beta}]]$	0.345		0.153	0.211	0.027	0.189	0.148
	0.217		0.073	0.075	0.017	0.094	0.073
	0.128		0.067	0.074	0.015	0.076	0.059
	0.107		0.057	0.075	0.011	0.070	0.055
$\hat{\beta}$: Range Minimum Maximum IQR	0.689	2.556	1.487	2.426	0.671	2.034	
	-0.310	-1.176	-0.629	-0.669	-0.770	-0.835	
	0.379	1.381	0.858	1.757	-0.098	1.199	
	0.127	0.674	0.312	0.298	0.131	0.245	
	0.401	1.001	0.686	0.991	0.405	0.595	
	-0.223	-0.604	-0.572	-0.870	-0.302	-0.437	
	0.179	0.397	0.114	0.121	0.104	0.157	
	0.078	0.199	0.090	0.095	0.079	0.118	
	0.235	1.046	0.724	1.102	0.384	0.693	
	-0.423	-0.421	-0.256	-0.583	-0.068	-0.329	
	-0.188	0.624	0.468	0.519	0.316	0.365	
	0.046	0.181	0.097	0.103	0.074	0.094	
0.202	0.894	0.864	1.380	0.255	0.665		
0.364	-0.253	-0.343	-0.961	0.172	-0.212		
0.566	0.640	0.522	0.419	0.427	0.453		
0.039	0.194	0.111	0.109	0.052	0.094		

Unbiasedness of the Regression Estimator

Based on $\hat{E}[\hat{\boldsymbol{\beta}}]$, S1S and CBI were similar, with little exhibited bias (the CBI intercept was a little lower than S1S, but just by 0.006). Not surprisingly, OLS and BI were severely biased as, on average, their respective estimates were not at all near $\boldsymbol{\beta}$. It was surprising to have LTS demonstrate a bias. It is conjectured that perhaps the effects of the LTS subsampling algorithm may have led to this result. While S1S was able to overcome this bias, M1S was not. This was an interesting result since the main difference between M1S and S1S lay in their respective handling (via weighting) of high leverage points. The HBK dataset, with its high influence cluster, demonstrated the advantage of S1S over M1S under such a condition.

Scale Estimation

The true scale was $\sigma_\varepsilon^2 = 0.25$. The MSE from OLS was an extremely positively biased scale estimator, as indicated by $\hat{E}[\hat{\sigma}^2] \approx 3.478$. From $\hat{E}[\hat{\sigma}^2] \approx 0.300$ it was determined that $\hat{\sigma}_{LTS}^2$ was, on average, merely a mild overestimate of scale. In response units it was a difference between $\hat{\sigma} = \sqrt{\hat{E}[\hat{\sigma}^2]} \approx 0.548$ and $\sigma_\varepsilon = 0.5$. No decent estimate of scale emerged from BI as $\hat{E}[\hat{\sigma}^2] \approx 0.485$ was indicative of an overestimate of scale while $\hat{E}[v^2] \approx 0.020$ signified a gross underestimate of scale. Overestimation was also indicated for the CBI scale estimator from $\hat{E}[\hat{\sigma}^2] \approx 0.341$. While $\hat{E}[v^2] \approx 0.436$ was even more biased, $\hat{E}[v_w^2] \approx 0.268$ demonstrated the effectiveness of v_w^2 as a scale estimator. In response units the later becomes $\hat{v}_w = \sqrt{\hat{E}[v_w^2]} \approx 0.518$ (versus $\sigma_\varepsilon = 0.5$). Thus, the CBI procedure using v_w^2 outperformed the other scale estimates with respect to bias.

Standard Errors

The expected standard errors for the BI coefficients were by the smallest due to the gross underestimate of scale. The breakdown of BI and OLS distorted their standard errors. M1S and

CBI (using v_w^2) had very similar standard errors, but S1S was the other regression estimator (besides CBI) with no pronounced bias. The CBI coefficients had lower standard errors than those for the S1S coefficients.

Coefficient Stability

That both OLS and BI exhibited very tight distributions for each of the four coefficients was of little consequence given the extreme bias that was exhibited. In other words, their breakdown was consistent. Of the high breakdown estimator group, LTS was generally outperformed and MIS, aside from the smallest intercept range, was not overly impressive. Between S1S and CBI, the CBI coefficients were more stable than the S1S coefficients, both in terms of the observed range as well as with respect to the IQR.

§8.5 Study #5: Random 40% Contamination with 3 Regressors

Given that one of the performance characteristics of the CBI regression method is the high breakdown point, a comparison between methods under the dataset condition that the level of contamination is quite large was performed. In this fifth Monte Carlo study, 40% of the observations were generated as outliers, some with high leverage and some with moderate to low leverage. Furthermore, these outliers were generated according to a second linear model so that the ability of each of the competing regression methods to sift out a secondary trend and still correctly determine the general trend could be investigated. A consequence of this scheme was that the outliers would not be generated under a symmetric rule about the general trend, thereby potentially introducing bias into a regression estimate. Specifically, the $n = 80$ observations had three regressor variables that were randomly generated independently according to the following rule:

If $i \leq 32$, then

$$\begin{aligned} x_{1i} &\sim N\left[\mu_{x_1} = 0, \sigma_{x_1}^2 = 100\right], \\ x_{2i} &\sim N\left[\mu_{x_2} = 0, \sigma_{x_2}^2 = 400\right] \text{ and} \\ x_{3i} &\sim N\left[\mu_{x_3} = 0, \sigma_{x_3}^2 = 25\right]. \end{aligned}$$

Else, if $i > 32$ then

$$\begin{aligned}x_{1i} &\sim N[\mu_{x_1} = 0, \sigma_{x_1}^2 = 16], \\x_{2i} &\sim N[\mu_{x_2} = 0, \sigma_{x_2}^2 = 4] \text{ and} \\x_{3i} &\sim N[\mu_{x_3} = 0, \sigma_{x_3}^2 = 1].\end{aligned}$$

Under such a construction, the expected location for each of the eighty generated regressor locations was the origin, $(0, 0, 0)$. However, the space filling ellipsoids are of dramatically different volumes and have different dominant axes. The first thirty-two observations are much more likely to possess high leverage than would the remaining forty-eight observations.

The response variable was then generated according to the linear model

$$y_i = \begin{cases} -30 + 0 x_{1i} - 5 x_{2i} + 10 x_{3i} + \varepsilon_i, & i \leq 32, \\ 100 - 10 x_{1i} + 5 x_{2i} + 20 x_{3i} + \varepsilon_i, & i > 32, \end{cases}$$

with random errors following

$$\varepsilon_i \sim \begin{cases} N[\mu_\varepsilon = 0, \sigma_\varepsilon^2 = 4], & i \leq 32, \\ N[\mu_\varepsilon = 0, \sigma_\varepsilon^2 = 1], & i > 32. \end{cases}$$

For Monte Carlo study #5A the regressor space was generated just once, with a new response vector generated for each simulation run. This regressor space is presented in Appendix B. For Monte Carlo study #5B the regressor space was generated anew for each simulation run (with a new response vector generated for each simulation run as well).

§8.5.1 Results for Monte Carlo Study #5A (Fixed Regressor Space)

Table 8.9 provides the simulation summary for the fixed regressor case of the fifth Monte Carlo study. The parameters to be estimated are $\beta' = [100 \quad -10 \quad 5 \quad 20]$ and $\sigma_\varepsilon^2 = 1$. A discussion of these results is now presented.

Table 8.9: Simulation results for Monte Carlo study #5A, fixed regressor case.

	OLS	LTS	MIS	SIS	BI	CBI	
$\hat{E}[\hat{\beta}]$	49.211	100.000	99.739	99.743	51.213	100.000	
	-0.741	-9.998	-9.835	-9.836	-0.919	-10.000	
	-4.853	4.997	4.270	4.275	-4.747	5.000	
	4.477	19.996	19.746	19.733	3.666	19.996	
$\hat{E}[\hat{\sigma}^2]$	5317.411	4.483			8418.740	4.448	$\hat{E}[v_w^2]$
$\hat{E}[v^2]$					346.234	1.788	0.975
$\hat{E}[se[\hat{\beta}]]$	8.329		0.228	0.228	2.163	0.210	0.155
	1.304		0.070	0.070	0.347	0.052	0.039
	0.549		0.262	0.260	0.145	0.094	0.069
	3.292		0.327	0.323	0.868	0.206	0.152
$\hat{\beta}$: Range Minimum Maximum IQR	1.077	1.307	1.022	1.023	1.133	1.072	
	48.661	99.330	99.210	99.212	50.644	99.432	
	49.738	100.637	100.232	100.235	51.777	100.504	
	0.221	0.268	0.217	0.217	0.225	0.220	
	0.257	0.366	0.318	0.321	0.236	0.300	
	-0.863	-10.180	-9.976	-9.976	-1.027	-10.135	
	-0.606	-9.814	-9.658	-9.656	-0.791	-9.835	
	0.045	0.078	0.059	0.058	0.044	0.055	
	0.099	0.734	0.949	0.942	0.103	0.515	
	-4.898	4.615	3.728	3.737	-4.794	4.748	
	-4.799	5.349	4.677	4.679	-4.691	5.263	
	0.020	0.137	0.158	0.157	0.021	0.100	
0.645	1.534	1.154	1.388	0.645	1.153		
4.138	19.128	19.212	19.195	3.311	19.369		
4.783	20.662	20.366	20.583	3.956	20.521		
0.120	0.293	0.198	0.200	0.127	0.209		

Unbiasedness of the Regression Estimator

From $\hat{E}[\hat{\beta}]$ it was apparent that only LTS and CBI demonstrated unbiasedness. OLS and BI were highly biased regression estimators, indicative of their respective breakdowns due to the secondary trend. Even with an unbiased LTS starting point, both MIS and SIS exhibited a bias in each of the four parameters, but especially with respect to the estimation of β_2 . An interesting result since x_2 was the dominant direction for high leverage. The asymmetric generation of the outliers posed difficulty for MIS and SIS.

Scale Estimation

Although the true scale was $\sigma_\varepsilon^2 = 1$, both OLS and BI yielded enormous expected estimates of scale due to their breakdown in identifying the general trend. The estimated expected value for MSE (OLS) was $\hat{E}[\hat{\sigma}^2] \approx 5317.411$. For BI, $\hat{E}[\hat{\sigma}^2] \approx 8418.740$ also signified a gross overestimate, with $\hat{E}[v^2] \approx 346.234$ indicating that v^2 is also severely positively biased. For LTS, $\hat{E}[\hat{\sigma}^2] \approx 4.483$ demonstrated that the LTS scale estimator was positively biased as well, where in response units it was a difference between $\hat{\sigma} = \sqrt{\hat{E}[\hat{\sigma}^2]} \approx 2.117$ and $\sigma_\varepsilon = 1$. With $\hat{E}[\hat{\sigma}^2] \approx 4.448$, the CBI scale estimator was very similar to that obtained via LTS. An improvement was demonstrated with $\hat{E}[v^2] \approx 1.788$, but $\hat{E}[v_w^2] \approx 0.975$ illustrated the clear advantage of v_w^2 as a scale estimator. In response units the later becomes $\hat{v}_w = \sqrt{\hat{E}[v_w^2]} \approx 0.987$, quite close to $\sigma_\varepsilon = 1$.

Standard Errors

The larger expected standard errors for the coefficients from both OLS and BI were indicative of their breakdown in detecting the general trend. The expected standard errors (using v_w^2) for the CBI coefficients were lower than those witnessed for either the M1S or S1S coefficients.

Coefficient Stability

The breakdown of OLS and BI offsets their ability to provide stable coefficients with tight sampling distributions. It was observed that the CBI coefficients had smaller ranges and, aside from the intercept, smaller IQR's than those obtained for LTS. M1S and S1S had bias issues, but still demonstrated less stability with respect to the estimation of β_2 .

§8.5.2 Results for Monte Carlo Study #5B (Random Regressor Space)

The second case of the fifth Monte Carlo study enabled the regressor space to be generated anew for each simulation run. The results of this study, Monte Carlo study #5B, are provided in Table 8.10.

Unbiasedness of the Regression Estimator

As in the fixed regressor case, by viewing $\hat{E}[\hat{\beta}]$ only LTS and CBI demonstrated unbiasedness. OLS and BI were highly biased regression estimators, again indicative of their respective breakdowns due to the secondary trend. Even with an unbiased LTS starting point,

Table 8.10: Simulation results for Monte Carlo study #5B, random regressor case.

	OLS	LTS	MIS	SIS	BI	CBI	
$\hat{E}[\hat{\beta}]$	50.577	99.999	99.640	99.639	55.034	100.000	
	-2.072	-9.999	-9.912	-9.912	-2.768	-9.999	
	-4.830	5.002	4.520	4.525	-4.710	5.001	
	10.638	19.995	19.769	19.770	10.981	19.997	
$\hat{E}[\hat{\sigma}^2]$	5124.729	4.386			7237.757	4.444	$\hat{E}[v_w^2]$
$\hat{E}[v^2]$					338.363	1.779	0.974
$\hat{E}[se[\hat{\beta}]]$	8.143		0.222	0.222	2.144	0.212	0.157
	1.187		0.078	0.078	0.324	0.054	0.040
	0.666		0.225	0.224	0.185	0.103	0.076
	2.597		0.395	0.393	0.715	0.211	0.156
$\hat{\beta}$: Range Minimum Maximum IQR	34.837	1.341	2.788	2.783	51.256	1.276	
	35.498	99.363	98.423	98.427	38.071	99.322	
	70.335	100.705	101.211	101.211	89.327	100.599	
	5.649	0.291	0.349	0.349	7.997	0.215	
	9.614	0.440	0.959	0.959	14.285	0.357	
	-7.697	-10.212	-10.311	-10.311	-12.415	-10.183	
	1.916	-9.772	-9.352	-9.352	1.870	-9.826	
	1.778	0.079	0.111	0.111	2.309	0.057	
	6.435	1.123	4.040	4.038	10.133	0.693	
	-7.926	4.498	1.272	1.272	-9.281	4.695	
	-1.491	5.621	5.313	5.310	0.852	5.388	
	1.056	0.155	0.335	0.343	1.368	0.115	
22.433	1.788	3.426	3.464	30.538	1.243		
-1.450	19.067	18.021	17.983	-2.793	19.437		
20.983	20.855	21.447	21.447	27.746	20.680		
4.035	0.336	0.575	0.580	5.186	0.230		

both M1S and S1S exhibited a bias in each of the four parameters, but especially with respect to the estimation of β_2 .

Scale Estimation

While the true scale was $\sigma_\varepsilon^2 = 1$, both OLS and BI yielded enormous expected estimates of scale due to their breakdown in identifying the general trend. The estimated expected value for MSE (OLS) was $\hat{E}[\hat{\sigma}^2] \approx 5124.729$. For BI, $\hat{E}[\hat{\sigma}^2] \approx 7237.757$ also signified a gross overestimate, with $\hat{E}[v^2] \approx 338.363$ indicating that v^2 was also severely positively biased. For LTS, $\hat{E}[\hat{\sigma}^2] \approx 4.386$ demonstrated that the LTS scale estimator was positively biased as well, where in response units it was a difference between $\hat{\sigma} = \sqrt{\hat{E}[\hat{\sigma}^2]} \approx 2.094$ and $\sigma_\varepsilon = 1$. Meanwhile, the CBI scale estimator had $\hat{E}[\hat{\sigma}^2] \approx 4.444$, very similar to that for LTS. An improvement was demonstrated with $\hat{E}[v^2] \approx 1.779$, but $\hat{E}[v_w^2] \approx 0.974$ again illustrated the clear advantage of v_w^2 as a scale estimator. In response units the later becomes $\hat{v}_w = \sqrt{\hat{E}[v_w^2]} \approx 0.987$, again quite close to $\sigma_\varepsilon = 1$.

Standard Errors

The larger expected standard errors for the coefficients from both OLS and BI were indicative of their breakdown in detecting the general trend. The expected standard errors (using v_w^2) for the CBI coefficients were lower than those witnessed for either the M1S or S1S coefficients.

Coefficient Stability

The breakdown of OLS and BI offsets their ability to provide stable coefficients with tight sampling distributions. It was observed that the CBI coefficients had smaller ranges and,

aside from the intercept, smaller IQR's than those obtained for LTS. MIS and SIS had bias issues, but also demonstrated less stability with respect to the estimation of β .

§8.6 Study #6: Clustered 40% Contamination (Random Sign) with 3 Regressors

The previous Monte Carlo study was supplemented with this next study in which the regressor structure for the 40% contamination was placed in a common high influence region of the regressor space instead of following a secondary linear trend. The $n = 80$ observations had three regressor variables that were randomly generated independently according to the following rule:

If $i \leq 32$, then

$$\begin{aligned} x_{1i} &\sim N[\mu_{X_1} = 15, \sigma_{X_1}^2 = 100], \\ x_{2i} &\sim N[\mu_{X_2} = 10, \sigma_{X_2}^2 = 400] \text{ and} \\ x_{3i} &\sim N[\mu_{X_3} = 60, \sigma_{X_3}^2 = 25]. \end{aligned}$$

Else, if $i > 32$ then

$$\begin{aligned} x_{1i} &\sim N[\mu_{X_1} = 0, \sigma_{X_1}^2 = 16], \\ x_{2i} &\sim N[\mu_{X_2} = 0, \sigma_{X_2}^2 = 4] \text{ and} \\ x_{3i} &\sim N[\mu_{X_3} = 0, \sigma_{X_3}^2 = 10]. \end{aligned}$$

The response variable was then generated according to the linear model

$$y_i = 50 + 5 x_{1i} - 10 x_{2i} + x_{3i} + \varepsilon_i,$$

with random errors following

$$\varepsilon_i \sim \begin{cases} U[250, 275], & i \leq 32, \\ N[\mu_\varepsilon = 0, \sigma_\varepsilon^2 = 1], & i > 32. \end{cases}$$

For Monte Carlo study #5A the regressor space was generated just once, with a new response vector generated for each simulation run. This regressor space is presented in

Appendix B. For Monte Carlo study #5B the regressor space was generated anew for each simulation run (with a new response vector generated for each simulation run as well).

§8.6.1 Results for Monte Carlo Study #6A (Fixed Regressor Space)

Table 8.11 provides the simulation summary for the fixed regressor case of the sixth Monte Carlo study. The parameters to be estimated are $\beta' = [50 \ 5 \ -10 \ 1]$ and $\sigma_\varepsilon^2 = 1$. A discussion of these results is now presented.

Unbiasedness of the Regression Estimator

LTS and CBI were the only two regression estimators that demonstrated unbiasedness, as

Table 8.11: Simulation results for Monte Carlo study #6A, fixed regressor case.

	OLS	LTS	MIS	SIS	BI	CBI	
$\hat{E}[\hat{\beta}]$	51.084	50.009	50.757	50.989	50.170	50.006	
	6.927	4.999	6.798	6.479	7.307	4.999	
	-1.956	-10.000	-2.696	-3.579	-1.831	-9.999	
	3.466	1.000	2.653	2.584	3.380	1.000	
$\hat{E}[\hat{\sigma}^2]$	525.746	4.492			280.430	7.697	$\hat{E}[v_w^2]$
$\hat{E}[v^2]$					26.141	1.677	0.963
$\hat{E}[se[\hat{\beta}]]$	3.343		0.212	0.261	0.813	0.196	0.149
	0.778		0.322	0.265	0.194	0.048	0.037
	1.209		1.298	1.141	0.313	0.097	0.073
	0.272		0.293	0.281	0.072	0.023	0.018
$\hat{\beta}$: Range Minimum Maximum IQR	0.984	1.311	5.409	7.023	4.398	1.018	
	50.576	49.348	47.592	47.079	47.468	49.494	
	51.560	50.658	53.000	54.102	51.865	50.511	
	0.194	0.277	0.340	1.223	1.489	0.200	
	0.439	0.374	1.909	3.182	1.418	0.255	
	6.689	4.810	6.015	5.598	6.692	4.877	
	7.128	5.184	7.924	8.781	8.110	5.132	
	0.081	0.075	0.341	0.568	0.469	0.050	
	1.030	0.883	6.980	8.271	2.979	0.545	
	-2.458	-10.396	-5.761	-6.640	-3.533	-10.311	
	-1.429	-9.513	1.219	1.631	-0.555	-9.766	
	0.185	0.155	1.327	1.497	0.993	0.103	
	0.170	0.185	1.667	1.704	0.574	0.124	
3.375	0.905	1.925	1.882	3.098	0.932		
3.546	1.090	3.592	3.586	3.672	1.056		
0.034	0.037	0.298	0.317	0.294	0.023		

evident from $\hat{E}[\hat{\beta}]$. OLS and BI were highly biased regression estimators, particularly with respect to the non-intercept coefficients. Although based on an unbiased LTS starting point, both MIS and SIS also exhibited a bias in each of the four parameters. Here, the high influence cluster always possessed a positive residual by its generation rule, thereby introducing an asymmetric outlier structure (with respect to the general trend) that biased MIS and SIS.

Scale Estimation

While the true scale was $\sigma_\varepsilon^2 = 1$, both OLS and BI yielded enormous expected estimates of scale due to their breakdown in identifying the general trend. The estimated expected value for MSE (OLS) was $\hat{E}[\hat{\sigma}^2] \approx 525.746$. For BI, $\hat{E}[\hat{\sigma}^2] \approx 280.430$ also signified a gross overestimate, with $\hat{E}[v^2] \approx 26.141$ indicating that v^2 was also severely positively biased. For LTS, $\hat{E}[\hat{\sigma}^2] \approx 4.492$ demonstrated that the LTS scale estimator was positively biased as well, where in response units it was a difference between $\hat{\sigma} = \sqrt{\hat{E}[\hat{\sigma}^2]} \approx 2.119$ and $\sigma_\varepsilon = 1$. With $\hat{E}[\hat{\sigma}^2] \approx 7.697$, the CBI scale estimator was very similar to that obtained via LTS. An improvement was demonstrated with $\hat{E}[v^2] \approx 1.677$, but $\hat{E}[v_w^2] \approx 0.963$ illustrated the clear advantage of using v_w^2 as a scale estimator. In response units the later becomes $\hat{v}_w = \sqrt{\hat{E}[v_w^2]} \approx 0.981$, modestly smaller than $\sigma_\varepsilon = 1$.

Standard Errors

The larger expected standard errors for the coefficients from both OLS and BI were indicative of their breakdown in detecting the general trend. The expected standard errors (using v_w^2) for the CBI coefficients were much lower than those witnessed for either the MIS or SIS coefficients.

Coefficient Stability

The breakdown of OLS and BI offsets their ability to provide stable coefficients with tight sampling distributions. It was observed that the CBI coefficients had smaller ranges and smaller IQR's than those obtained for LTS, MIS or SIS. The differences were more dramatic with respect to MIS and SIS versus CBI.

§8.6.2 Results for Monte Carlo Study #6B (Random Regressor Space)

The second case of the sixth Monte Carlo study enabled the regressor space to be generated anew for each simulation run. The results of this study, Monte Carlo study #6B, are provided in Table 8.12.

Table 8.12: Simulation results for Monte Carlo study #6B, random regressor case.

	OLS	LTS	MIS	SIS	BI	CBI	
$\hat{E}[\hat{\beta}]$	53.444	49.991	52.144	52.112	53.001	49.992	
	8.248	4.999	7.458	7.202	8.280	5.000	
	-2.327	-10.006	-3.459	-4.155	-2.566	-10.002	
	3.141	1.000	2.565	2.393	3.196	1.000	
$\hat{E}[\hat{\sigma}^2]$	583.227	4.524			370.221	7.780	$\hat{E}[v_w^2]$
$\hat{E}[v^2]$					28.124	1.688	0.971
$\hat{E}[se[\hat{\beta}]]$	3.522		0.639	0.585	0.813	0.201	0.153
	0.796		0.440	0.394	0.188	0.051	0.039
	1.394		1.168	1.043	0.327	0.102	0.077
	0.256		0.277	0.247	0.060	0.020	0.015
$\hat{\beta}$: Range Minimum Maximum IQR	29.301	1.365	35.522	29.138	30.008	1.053	
	37.969	49.352	35.904	37.308	37.464	49.463	
	67.270	50.717	71.426	66.446	67.472	50.515	
	5.851	0.286	4.583	4.054	6.003	0.204	
	7.127	0.443	10.801	9.691	7.638	0.294	
	4.562	4.793	3.864	3.593	4.277	4.865	
	11.689	5.236	14.665	13.284	11.915	5.158	
	1.265	0.078	1.508	1.305	1.316	0.051	
	11.885	0.882	26.660	23.071	12.556	0.552	
	-8.580	-10.459	-9.623	-9.397	-9.106	-10.271	
3.305	-9.577	17.037	13.674	3.450	-9.719		
2.155	0.161	3.165	2.671	2.192	0.102		
2.270	0.192	4.939	4.071	2.264	0.122		
2.014	0.914	1.439	1.451	2.124	0.942		
4.284	1.105	6.379	5.523	4.388	1.064		
0.411	0.032	0.650	0.572	0.441	0.020		

Unbiasedness of the Regression Estimator

As in the fixed regressor case, it was apparent from $\hat{E}[\hat{\beta}]$ that only LTS and CBI demonstrated unbiasedness. OLS and BI were highly biased regression estimators, again indicative of their respective breakdowns due to the high influence cluster. Even with an unbiased LTS initial value, both M1S and S1S exhibited a substantial bias in each of the four parameters.

Scale Estimation

Although the true scale was $\sigma_\varepsilon^2 = 1$, both OLS and BI yielded enormous expected estimates of scale due to their breakdown in identifying the general trend. The estimated expected value for MSE (OLS) was $\hat{E}[\hat{\sigma}^2] \approx 583.227$. For BI, $\hat{E}[\hat{\sigma}^2] \approx 370.221$ also signified a gross overestimate, with $\hat{E}[v^2] \approx 28.124$ indicating that v^2 was also severely positively biased. For LTS, $\hat{E}[\hat{\sigma}^2] \approx 4.524$ demonstrated that the LTS scale estimator was positively biased as well, where in response units it was a difference between $\hat{\sigma} = \sqrt{\hat{E}[\hat{\sigma}^2]} \approx 2.127$ and $\sigma_\varepsilon = 1$. The CBI scale estimator was very similar to that of LTS since $\hat{E}[\hat{\sigma}^2] \approx 7.780$. An improvement was demonstrated with $\hat{E}[v^2] \approx 1.688$, but $\hat{E}[v_w^2] \approx 0.971$ illustrated the clear advantage of using v_w^2 as a scale estimator. In response units the later becomes $\hat{v}_w = \sqrt{\hat{E}[v_w^2]} \approx 0.985$, modestly lower than $\sigma_\varepsilon = 1$.

Standard Errors

The larger expected standard errors for the coefficients from both OLS and BI were indicative of their breakdown in detecting the general trend. The expected standard errors (using v_w^2) for the CBI coefficients were lower than those witnessed for either the M1S or S1S coefficients.

Coefficient Stability

Again, the breakdown of OLS and BI overrides the stability of their respective coefficients. It was observed that the CBI coefficients had smaller ranges and smaller IQR's than those obtained for LTS. MIS and SIS each demonstrated highly variable sampling distributions for each of the four coefficients, well beyond the magnitude exhibited by LTS and CBI.

§8.7 Study #7: Clustered 40% Contamination (Single Cluster) with 3 Regressors

The previous Monte Carlo study was supplemented with this next study in which the regressor structure for the 40% contamination was still clustered together to form a single large high influence cluster, but that this cluster could, across replicates, fall either above or below the general trend. This action allows for an investigation of the bias exhibited during a process that has outliers symmetrically distributed about the general trend. Here, the $n = 80$ observations had three regressor variables that were randomly generated independently according to the following rule:

If $i \leq 32$, then

$$\begin{aligned}x_{1i} &\sim N[\mu_{X_1} = 20, \sigma_{X_1}^2 = 1], \\x_{2i} &\sim N[\mu_{X_2} = 5, \sigma_{X_2}^2 = 1] \text{ and} \\x_{3i} &\sim N[\mu_{X_3} = 20, \sigma_{X_3}^2 = 1].\end{aligned}$$

Else, if $i > 32$ then

$$\begin{aligned}x_{1i} &\sim N[\mu_{X_1} = 0, \sigma_{X_1}^2 = 16], \\x_{2i} &\sim N[\mu_{X_2} = 0, \sigma_{X_2}^2 = 4] \text{ and} \\x_{3i} &\sim N[\mu_{X_3} = 0, \sigma_{X_3}^2 = 1].\end{aligned}$$

The response variable was then generated according to the linear model

$$y_i = 50 + 5 x_{1i} - 10 x_{2i} + \varepsilon_i,$$

with random errors following

$$\varepsilon_i \sim \begin{cases} N[\mu_\varepsilon = 0, \sigma_\varepsilon^2 = 25], & i \notin (1:32), \\ \text{sign}(U[0,1]-0.5) \cdot U[1000, 1250], & i \in (1:32). \end{cases}$$

Here, the expression $\text{sign}(U[0,1]-0.5)$ represents a random sign, where half the time it is equal to -1 and half the time it is equal to $+1$, on average. This random sign is generated once per simulation replicate; in other words it represents a common sign for each and every outlier. It is noted that this study had a larger error variance and more extreme outliers than in the previous study.

For Monte Carlo study #7A the regressor space was generated just once, with a new response vector generated for each simulation run. This regressor space is presented in Appendix B. For Monte Carlo study #7B the regressor space was generated anew for each simulation run (with a new response vector generated for each simulation run as well).

§8.7.1 Results for Monte Carlo Study #7A (Fixed Regressor Space)

Table 8.13 provides the simulation summary for the fixed regressor case of the seventh Monte Carlo study. Again, the parameters to be estimated are $\boldsymbol{\beta}' = [50 \ 5 \ -10 \ 0]$ and $\sigma_\varepsilon^2 = 25$. A discussion of these results is now presented in four subsections.

Unbiasedness of the Regression Estimator

LTS and CBI, by viewing $\hat{E}[\hat{\boldsymbol{\beta}}]$, were on average in agreement with the true underlying model parameters. However, OLS and BI exhibited minor bias in estimating β_3 , a parameter that was zero. Actually, OLS and BI were bi-modal; a biased central tendency for each location (above or below the general trend) of the high influence cluster that essentially canceled out to near unbiasedness. Furthermore, M1S and S1S demonstrated an even larger bias in estimating β_3 , with a smaller bias in estimating β_0 . In other words, some of the noise was being modeled as being attributable to x_3 , with this requiring an offset on the intercept.

Table 8.13: Simulation results for Monte Carlo study #7A, fixed regressor case.

	OLS	LTS	MIS	SIS	BI	CBI	
$\hat{E}[\hat{\beta}]$	50.110	49.993	50.556	50.517	50.116	49.983	
	5.035	4.998	5.089	5.109	5.035	5.001	
	-9.929	-10.000	-9.984	-10.061	-9.931	-10.000	
	0.444	-0.020	3.451	2.890	0.444	-0.014	
$\hat{E}[\hat{\sigma}^2]$	4678.875	113.013			4003.166	193.079	$\hat{E}[v_w^2]$
$\hat{E}[v^2]$					217.926	41.833	24.050
$\hat{E}[se[\hat{\beta}]]$	0.235		10.447	7.582	2.281	1.038	0.787
	0.055		1.951	1.496	0.540	0.249	0.188
	0.086		0.906	0.674	0.835	0.395	0.300
	0.066		53.220	48.427	0.645	1.150	0.872
$\hat{\beta}$: Range Minimum Maximum IQR	33.941	7.753	275.864	222.485	35.015	5.170	
	33.151	46.252	-100.931	-54.547	32.606	47.340	
	67.092	54.005	174.933	167.938	67.622	52.511	
	24.214	1.493	116.546	80.735	25.784	1.097	
	10.554	2.017	44.313	54.153	10.436	1.380	
	-0.224	4.024	-17.713	-20.662	-0.225	4.276	
	10.329	6.042	26.600	33.491	10.211	5.656	
	6.723	0.391	21.699	15.670	6.776	0.257	
	12.550	3.046	31.549	31.720	12.564	2.025	
	-16.477	-11.551	-25.609	-25.617	-16.555	-11.063	
	-3.927	-8.505	5.941	6.102	-3.991	-9.038	
	4.179	0.620	8.033	3.595	3.788	0.419	
110.007	8.995	920.195	836.035	110.175	6.078		
-55.421	-4.706	-455.702	-412.925	-55.404	-3.082		
54.586	4.289	464.492	423.110	54.772	2.996		
102.863	1.831	596.828	542.255	102.821	1.177		

Scale Estimation

As stated earlier, the true scale was $\sigma_\varepsilon^2 = 25$. While, on average, OLS and BI were not too deviant from being unbiased in terms of regression estimation, the scale estimates for each method were clearly affected by a breakdown in fit. With $\hat{E}[\hat{\sigma}^2] \approx 4678.875$, the MSE (OLS) was completely overwhelmed and highly biased. Likewise, with $\hat{E}[\hat{\sigma}^2] \approx 4003.166$ and $\hat{E}[v^2] = 217.926$, either BI scale estimator became excessively huge as well. Furthermore, $\hat{E}[\hat{\sigma}^2] \approx 113.013$ indicated that $\hat{\sigma}_{LTS}^2$ was also substantially positively biased. In response units

it was a difference between $\hat{\sigma} = \sqrt{\hat{E}[\hat{\sigma}^2]} \approx 10.631$ and $\sigma_\varepsilon = 5$. The CBI scale estimator witnessed $\hat{E}[\hat{\sigma}^2] \approx 193.079$, yet another gross overestimate. Yet while $\hat{E}[v^2] \approx 41.833$ was an improvement, $\hat{E}[v_w^2] \approx 24.050$ was even better. In response units the later becomes $\hat{v}_w = \sqrt{\hat{E}[v_w^2]} \approx 4.904$ (versus $\sigma_\varepsilon = 5$). Thus, v_w^2 was the only scale estimator to be anywhere near the correct parametric value.

Standard Errors

The expected standard errors for the OLS coefficients were the smallest of the group, although this fact alone masks the breakdown that occurred in one of two directions. The standard errors (using v_w^2) for the CBI coefficients were much lower than those obtained for either the MIS coefficients or S1S coefficients, the later two methods struggling to handle x_3 . The standard errors for $\hat{\beta}_3$ were enormous.

Coefficient Stability

The MIS and S1S coefficients were by far the least stable of the group, with substantially larger ranges and IQR's across four coefficients. Even though their initial estimator, LTS, performed adequately, the one-step improvements were prone to wild departures from the general trend. Meanwhile, CBI was more stable than LTS, with smaller ranges and smaller IQR's across the four coefficients. BI, like OLS, followed the high influence cluster, whether it was above or below the general trend. Thus, BI, like OLS, had a bi-modal sampling distribution for its coefficients, neither being unbiased. The symmetric nature of the outlier generation allowed OLS and BI to appear better than they really performed.

§8.7.2 Results for Monte Carlo Study #7B (Random Regressor Space)

The second case of the seventh Monte Carlo study enabled the regressor space to be generated anew for each simulation run. The results of this study, Monte Carlo study #7B, are provided in Table 8.14.

Table 8.14: Simulation results for Monte Carlo study #7B, random regressor case.

	OLS	LTS	MIS	SIS	BI	CBI	
$\hat{E}[\hat{\beta}]$	50.093	50.023	50.810	50.905	50.080	49.998	
	4.930	4.994	4.889	4.904	4.933	5.000	
	-9.902	-10.002	-9.591	-9.615	-9.891	-9.999	
	-0.578	-0.004	1.113	1.084	-0.583	0.005	
$\hat{E}[\hat{\sigma}^2]$	4874.563	113.384			4354.882	192.635	$\hat{E}[v_w^2]$
$\hat{E}[v^2]$					233.476	41.919	24.091
$\hat{E}[se[\hat{\beta}]]$	0.229		1.702	1.635	2.274	1.005	0.762
	0.054		0.693	0.652	0.546	0.255	0.193
	0.106		1.016	0.959	1.069	0.510	0.386
	0.061		8.869	8.226	0.617	1.018	0.772
$\hat{\beta}$: Range Minimum Maximum IQR	75.699	7.555	250.246	197.427	74.452	5.701	
	16.379	46.116	-84.596	-47.024	16.731	47.271	
	92.078	53.671	165.650	150.403	91.183	52.972	
	12.791	1.437	1.779	1.770	12.731	1.039	
	25.536	2.067	121.742	124.666	25.228	1.595	
	-6.810	3.932	-58.405	-62.221	-6.948	4.259	
	18.726	5.999	63.337	62.445	18.280	5.854	
	9.858	0.403	0.620	0.609	9.811	0.272	
	38.435	4.424	177.934	132.227	38.254	2.897	
	-29.152	-12.173	-105.701	-73.702	-29.114	-11.307	
	9.283	-7.749	72.233	58.526	9.139	-8.409	
	9.432	0.805	0.941	0.934	9.344	0.524	
115.197	8.952	847.288	701.404	115.503	5.880		
-58.036	-4.689	-450.312	-345.583	-57.972	-3.154		
57.161	4.262	396.976	355.821	57.530	2.727		
99.290	1.540	7.926	7.732	99.535	1.045		

Unbiasedness of the Regression Estimator

As in the fixed regressor case, LTS and CBI, by viewing $\hat{E}[\hat{\beta}]$, were on average in agreement with the true underlying model parameters. However, OLS and BI still exhibited minor bias in estimating β_3 , a parameter that was zero. OLS and BI were again bi-modal with a biased central tendency for each location (above or below the general trend) of the high influence cluster that essentially canceled out to near unbiasedness. MIS and SIS demonstrated a bias in estimating β_3 , with a smaller bias in estimating the other three parameters.

Scale Estimation

The true scale was $\sigma_\varepsilon^2 = 25$ with OLS and BI still yielding enormous scale estimates attributable to a breakdown in fit. With $\hat{E}[\hat{\sigma}^2] \approx 4874.563$, the MSE (OLS) was completely overwhelmed and highly biased. Likewise, with $\hat{E}[\hat{\sigma}^2] \approx 4354.882$ and $\hat{E}[v^2] = 233.476$, either BI scale estimator became excessively huge as well. Furthermore, $\hat{E}[\hat{\sigma}^2] \approx 113.384$ indicated that $\hat{\sigma}_{LTS}^2$ was also substantially positively biased. In response units it was a difference between $\hat{\sigma} = \sqrt{\hat{E}[\hat{\sigma}^2]} \approx 10.648$ and $\sigma_\varepsilon = 5$. The CBI scale estimator witnessed $\hat{E}[\hat{\sigma}^2] \approx 192.635$, yet another gross overestimate. Yet while $\hat{E}[v^2] \approx 41.919$ was an improvement, $\hat{E}[v_w^2] \approx 24.091$ was even better. In response units the later becomes $\hat{v}_w = \sqrt{\hat{E}[v_w^2]} \approx 4.908$ (versus $\sigma_\varepsilon = 5$). Again, v_w^2 was the only scale estimator to be anywhere near the correct parametric value.

Standard Errors

As before, the expected standard errors for the OLS coefficients were the smallest of the group, although this fact alone masks the breakdown that occurred in one of two directions. The standard errors (using v_w^2) for the CBI coefficients were much lower than those obtained for either the MIS coefficients or SIS coefficients, the later two methods still struggling to handle x_3 . The standard errors for $\hat{\beta}_3$ remained quite large in this random regressor case.

Coefficient Stability

The MIS and SIS coefficients were by far the least stable of the group, with substantially larger ranges and larger IQR's across the four coefficients. Again, this occurred even though their initial estimator, LTS, performed adequately. Meanwhile, CBI was more stable than LTS, with smaller ranges and smaller IQR's across the four coefficients.

§8.8 Study #8: Mixed 40% Contamination with 5 Regressors

A final simulation study was constructed so that the dataset has a complex contamination structure and five regressor variables. The $n=100$ observations were classified into five categories. There were 55 observations that comprised the general trend. An additional 5 observations formed a good high leverage cluster. Another 5 observations became randomly dispersed low leverage outliers. Yet another 5 observations were randomly dispersed high influence outliers. The remaining 30 observations formed a high influence outlier cluster.

For every observation, $1 \leq i \leq 100$,

$$x_{4i} \sim U[0, 10] \text{ and}$$

$$x_{5i} \sim U[-5, 5].$$

The first three regressor variables were generated according to specific rules in order to produce the desired data structure. Namely, if $i \leq 55$ or $66 \leq i \leq 70$, then

$$x_{1i} \sim N[\mu_{x_1} = 0, \sigma_{x_1}^2 = 4],$$

$$x_{2i} \sim N[\mu_{x_2} = 0, \sigma_{x_2}^2 = 2] \text{ and}$$

$$x_{3i} \sim N[\mu_{x_3} = 0, \sigma_{x_3}^2 = 1].$$

Else, if $56 \leq i \leq 60$, then

$$x_{1i} \sim N[\mu_{x_1} = -15, \sigma_{x_1}^2 = 4],$$

$$x_{2i} \sim N[\mu_{x_2} = 5, \sigma_{x_2}^2 = 2] \text{ and}$$

$$x_{3i} \sim N[\mu_{x_3} = -10, \sigma_{x_3}^2 = 1].$$

Else, if $61 \leq i \leq 65$, then first compute three high leverage directional indicators u_{1i} , u_{2i} and u_{3i} (each observation generates its own indicators) such that

$$(u_{1i}, u_{2i}, u_{3i}) \in \{(1, 0, 0), (0, 1, 0), (0, 0, 1), (1, 1, 0), (1, 0, 1), (0, 1, 1)\},$$

with each of these six realizations being equally likely. Then,

$$x_{1i} \sim N[\mu_{x_1} = \text{sign}(U[0, 1] - 0.5)_i \cdot U[12, 22]_i \cdot u_{1i}, \sigma_{x_1}^2 = 4],$$

$$x_{2i} \sim N\left[\mu_{x_2} = \text{sign}(U[0,1]-0.5)_i \cdot U[6,12]_i \cdot u_{2i}, \sigma_{x_2}^2 = 2\right] \text{ and}$$

$$x_{3i} \sim N\left[\mu_{x_3} = \text{sign}(U[0,1]-0.5)_i \cdot U[3,8]_i \cdot u_{3i}, \sigma_{x_3}^2 = 1\right].$$

Finally, if $i > 70$ then

$$x_{1i} \sim N\left[\mu_{x_1} = 20, \sigma_{x_1}^2 = 1\right],$$

$$x_{2i} \sim N\left[\mu_{x_2} = 5, \sigma_{x_2}^2 = 1\right] \text{ and}$$

$$x_{3i} \sim N\left[\mu_{x_3} = 20, \sigma_{x_3}^2 = 1\right].$$

The response variable was then generated according to the linear model

$$y_i = \begin{cases} 50 + 5x_{1i} - 10x_{2i} + x_{3i} + \varepsilon_i, & i \leq 70, \\ -200 + 5x_{2i} - 20x_{3i} + \varepsilon_i, & i > 70, \end{cases}$$

with random errors following

$$\varepsilon_i \sim \begin{cases} N\left[\mu_\varepsilon = 0, \sigma_\varepsilon^2 = 1\right], & i \leq 60, \\ \text{sign}(U[0,1]-0.5)_i \cdot U[15, 25], & 61 \leq i \leq 70, \\ N\left[\mu_\varepsilon = 0, \sigma_\varepsilon^2 = 9\right], & i > 70. \end{cases}$$

Here, the expression $\text{sign}(U[0,1]-0.5)_i$ represents a random sign that is generated independently for each of the ten observations from sixty-one to seventy (inclusive). These signs are generated anew for each simulation replicate.

§8.8.1 Results for Monte Carlo Study #8A (Fixed Regressor Space)

Table 8.15 provides the simulation summary for the fixed regressor case of the eighth Monte Carlo study. Again, the parameters to be estimated are $\beta' = [50 \ 5 \ -10 \ 1 \ 0 \ 0]$ and $\sigma_\varepsilon^2 = 1$. A discussion of these results is now presented.

Unbiasedness of the Regression Estimator

LTS and CBI, by viewing $\hat{E}[\hat{\boldsymbol{\beta}}]$, demonstrated unbiasedness. OLS and BI, however, still substantial bias in estimating $\boldsymbol{\beta}$. These two methods clearly had difficulties detecting the general trend given this level and structure of contamination. M1S and S1S demonstrated smaller amounts of bias in each of the six coefficients. As with other studies, there appears to be some interplay between insignificant terms and the intercept.

Scale Estimation

The true scale was $\sigma_\varepsilon^2 = 1$, with OLS and BI yielding enormous scale estimates attributable to a breakdown in fit. With $\hat{E}[\hat{\sigma}^2] \approx 5106.134$, the MSE (OLS) was completely overwhelmed and highly biased. Likewise, with $\hat{E}[\hat{\sigma}^2] \approx 2185.454$ and $\hat{E}[v^2] = 89.827$, either BI scale estimator became excessively huge as well. Furthermore, $\hat{E}[\hat{\sigma}^2] \approx 2.350$ indicated that $\hat{\sigma}_{LTS}^2$ was also moderately positively biased. In response units it was a difference between $\hat{\sigma} = \sqrt{\hat{E}[\hat{\sigma}^2]} \approx 1.533$ and $\sigma_\varepsilon = 1$. The CBI scale estimator witnessed $\hat{E}[\hat{\sigma}^2] \approx 5.677$, yet another overestimate. Yet while $\hat{E}[v^2] \approx 1.673$ was an improvement, $\hat{E}[v_w^2] \approx 0.966$ was even better. In response units the later becomes $\hat{v}_w = \sqrt{\hat{E}[v_w^2]} \approx 0.983$ (versus $\sigma_\varepsilon = 1$). Thus, v_w^2 was the least biased estimator of scale.

Standard Errors

The expected standard errors (using v_w^2) for the CBI coefficients were the smallest of the group, while those for the OLS coefficients were the largest overall. M1S and S1S were similar.

Coefficient Stability

The CBI coefficients were by far the most stable of the group, with substantially smaller ranges and comparable IQR's across all six coefficients. M1S and S1S marginally improved

stability versus LTS for a few coefficients, but this gain was offset by the introduction of bias that occurred with these one-step improvements.

Table 8.15: Simulation results for Monte Carlo study #8A, fixed regressor case.

	OLS	LTS	MIS	SIS	BI	CBI	
$\hat{E}[\hat{\beta}]$	54.177	50.026	50.879	50.853	73.589	49.993	
	0.464	4.999	4.927	4.931	0.180	5.000	
	-25.246	-9.999	-10.160	-10.162	-11.398	-9.998	
	-22.515	0.999	0.428	0.464	-28.311	1.000	
	4.559	0.001	0.050	0.051	-0.927	-0.001	
	6.991	0.004	0.142	0.143	1.670	-0.001	
$\hat{E}[\hat{\sigma}^2]$	5106.134	2.350			2185.454	5.677	$\hat{E}[v_w^2]$
$\hat{E}[v^2]$					89.827	1.673	0.966
$\hat{E}[se[\hat{\beta}]]$	39.292		1.235	1.239	5.403	0.955	0.725
	1.421		0.080	0.081	0.250	0.040	0.031
	2.516		0.107	0.107	0.508	0.079	0.060
	1.625		0.192	0.197	0.292	0.079	0.060
	2.637		0.098	0.098	0.374	0.070	0.053
	5.165		0.165	0.165	0.725	0.133	0.101
$\hat{\beta}$: Range Minimum Maximum IQR	20.024	10.210	7.811	8.788	23.531	4.893	
	43.638	45.021	47.116	46.791	62.274	47.722	
	63.662	55.231	54.927	55.579	85.805	52.616	
	4.764	2.178	1.688	1.756	5.888	1.011	
	1.356	0.466	0.530	0.667	3.436	0.237	
	-0.202	4.750	4.632	4.542	-1.287	4.885	
	1.154	5.216	5.161	5.209	2.149	5.122	
	0.406	0.096	0.109	0.132	1.224	0.044	
	1.833	0.980	0.679	0.775	7.068	0.413	
	-26.167	-10.466	-10.524	-10.604	-15.928	-10.207	
	-24.334	-9.486	-9.845	-9.829	-8.860	-9.794	
	0.611	0.176	0.138	0.140	1.508	0.083	
	1.465	1.205	1.683	2.392	3.818	0.432	
	-23.261	0.406	-0.742	-1.323	-30.111	0.796	
	-21.796	1.611	0.942	1.069	-26.293	1.228	
	0.431	0.180	0.237	0.426	1.130	0.081	
2.083	0.780	0.679	0.702	2.805	0.352		
3.545	-0.361	-0.292	-0.290	-2.277	-0.196		
5.629	0.418	0.388	0.411	0.528	0.155		
0.477	0.153	0.134	0.136	0.597	0.073		
2.692	1.355	1.045	1.156	3.356	0.723		
5.565	-0.669	-0.372	-0.362	0.052	-0.340		
8.257	0.686	0.673	0.795	3.409	0.382		
0.602	0.305	0.224	0.228	0.758	0.146		

§8.8.2 Results for Monte Carlo Study #8B (Random Regressor Space)

The second case of the eighth Monte Carlo study enabled the regressor space to be generated anew for each simulation run. The results of this study, Monte Carlo study #8B, are provided in Table 8.16.

Table 8.16: Simulation results for Monte Carlo study #8B, random regressor case.

	OLS	LTS	MIS	SIS	BI	CBI	
$\hat{E}[\hat{\beta}]$	26.488	49.982	47.611	47.680	47.556	50.027	
	5.126	5.000	5.058	5.116	2.705	5.000	
	-25.560	-9.999	-12.886	-12.892	-10.990	-9.997	
	-27.651	1.001	-4.647	-3.912	-30.846	1.001	
	0.059	0.004	0.019	0.012	0.032	0.000	
	0.184	0.000	0.016	0.025	0.033	0.004	
$\hat{E}[\hat{\sigma}^2]$	4914.801	2.342			1728.356	5.637	$\hat{E}[v_w^2]$
$\hat{E}[v^2]$					79.396	1.662	0.959
$\hat{E}[se[\hat{\beta}]]$	40.208		2.351	2.348	5.373	1.015	0.771
	1.587		0.101	0.108	0.245	0.044	0.033
	2.659		0.598	0.600	0.487	0.090	0.068
	1.888		1.088	0.964	0.285	0.092	0.070
	2.489		0.143	0.144	0.332	0.063	0.048
	4.983		0.290	0.290	0.665	0.126	0.096
$\hat{\beta}$: Range Minimum Maximum IQR	271.704	11.790	126.395	157.261	156.986	6.415	
	-101.799	44.302	-21.285	-49.112	-29.057	47.122	
	169.905	56.092	105.110	108.149	127.929	53.537	
	51.576	2.272	6.226	5.753	26.282	1.060	
	17.021	0.537	11.181	10.574	11.731	0.217	
	-2.658	4.743	-3.583	-1.953	-2.237	4.888	
	14.363	5.279	7.598	8.621	9.494	5.104	
	3.050	0.099	0.263	0.359	1.771	0.046	
	26.934	1.134	14.671	19.407	30.460	0.499	
	-40.752	-10.667	-24.283	-25.879	-34.408	-10.255	
	-13.818	-9.533	-9.612	-6.472	-3.948	-9.756	
	4.739	0.206	5.292	5.317	3.152	0.091	
	17.834	1.308	53.819	78.187	10.789	0.788	
	-37.362	0.349	-52.867	-76.994	-35.214	0.574	
	-19.528	1.657	0.952	1.194	-24.425	1.362	
	3.573	0.209	9.432	7.683	2.002	0.095	
16.748	0.670	7.762	7.690	11.806	0.323		
-7.897	-0.336	-3.947	-3.769	-6.362	-0.177		
8.851	0.334	3.815	3.921	5.444	0.146		
3.446	0.138	0.270	0.267	1.649	0.066		
30.709	1.416	14.512	21.136	18.069	0.817		
-15.062	-0.681	-7.974	-11.885	-9.132	-0.401		
15.647	0.735	6.538	9.251	8.937	0.416		
6.574	0.278	0.539	0.545	3.301	0.126		

Unbiasedness of the Regression Estimator

As with the fixed regressor case, by viewing $\hat{E}[\hat{\boldsymbol{\beta}}]$ it was observed that both LTS and CBI demonstrated unbiasedness. OLS and BI, however, still had substantial bias in estimating $\boldsymbol{\beta}$. Again, these two methods clearly had difficulties detecting the general trend given this level and structure of contamination. M1S and S1S demonstrated larger amounts of bias in each of the six coefficients during this random regressor case than they did during the fixed regressor case.

Scale Estimation

The true scale was $\sigma_\varepsilon^2 = 1$, with OLS and BI yielding enormous scale estimates attributable to a breakdown in fit. With $\hat{E}[\hat{\sigma}^2] \approx 4914.801$, the MSE (OLS) was completely overwhelmed and highly biased. Likewise, with $\hat{E}[\hat{\sigma}^2] \approx 1728.356$ and $\hat{E}[v^2] = 79.396$, either BI scale estimator became excessively huge as well. Furthermore, $\hat{E}[\hat{\sigma}^2] \approx 2.342$ indicated that $\hat{\sigma}_{LTS}^2$ was also moderately positively biased. In response units it was a difference between $\hat{\sigma} = \sqrt{\hat{E}[\hat{\sigma}^2]} \approx 1.530$ and $\sigma_\varepsilon = 1$. The CBI scale estimator witnessed $\hat{E}[\hat{\sigma}^2] \approx 5.637$, yet another overestimate. Yet while $\hat{E}[v^2] \approx 1.662$ was an improvement, $\hat{E}[v_w^2] \approx 0.959$ was even better. In response units the later becomes $\hat{v}_w = \sqrt{\hat{E}[v_w^2]} \approx 0.979$ (versus $\sigma_\varepsilon = 1$). Thus, v_w^2 was again the least biased estimator of scale.

Standard Errors

The expected standard errors (using v_w^2) for the CBI coefficients were the smallest of the group, while those for the OLS coefficients were the largest overall. M1S and S1S were again very similar to one another.

Coefficient Stability

The CBI coefficients were by far the most stable of the group, with substantially smaller ranges and comparable IQR's across all six coefficients. M1S and S1S exhibited large variability in the intercept (as did OLS), with very extreme observations for the other five coefficients as well ($\hat{\beta}_3$ was particularly unstable for M1S and S1S).

§8.9 Chapter Summary

This chapter focused on comparing six regression methods under a variety of data conditions. The results were discussed relative to issues of unbiasedness of the coefficient estimates, unbiasedness of the scale estimate, expected standard errors for the coefficient estimates, high breakdown capability and coefficient stability. The following are the major findings drawn from the Monte Carlo studies.

- Regarding the CBI estimation procedure, there was a clear, marked improvement in the estimate of scale when v_w^2 was used instead of v^2 or $\hat{\sigma}_{CBI}$. The benefit over competing estimators of scale was typically most evident and dramatic when the level of contamination was large and the effective sample size was much smaller than the actual sample size.
- There was no surprise that BI regression can demonstrate a bias in the coefficients under certain dataset structures (Monte Carlo study #2 is one case in point). Its low breakdown point has been well documented in the literature. However, BI regression did offer a baseline for comparison when the data has no contamination present.
- It was surprising to observe the frequent bias in the M1S and S1S coefficients. While LTS was unbiased, each one-step method often showed evidence of coefficient drift away from the true parameter values.
- In general, the scale estimate, v_w^2 , for the CBI method was considerably better than the LTS scale estimate with respect to expected bias and had an adequate performance under the no contamination scenario. Related to the performance of the scale estimation, the CBI procedure produced standard errors for the coefficients that were very competitive

with the other high breakdown procedures. In fact, The CBI standard errors were generally smaller, often rather dramatically smaller than the standard errors for LTS, M1S or S1S.

- In theory, high breakdown capability relates to the ability to drive a coefficient to infinity by the movement of certain observations. However, it was also interesting to view the sampling distributions for each method's coefficient estimates and determine if there was evidence that more stability was offered by a particular method over another. Throughout these Monte Carlo studies, the CBI estimates dominated the stability analysis by consistently possessing small ranges when contamination was present. In the few instances where a CBI coefficient did not possess the smallest range outright, it was nonetheless competitive. Often, the CBI ranges were considerably smaller than those observed by LTS, M1S and S1S. In addition, it was also witnessed that LTS may demonstrate smaller ranges than either of the one-step improvements based on LTS. Regarding the central bulk of the sampling distribution for a particular coefficient, the IQR was viewed. The CBI method was quite competitive versus the other methods with respect to this quantity as well; often possessing the smallest IQR's in contaminated data settings.

Whether analyzing a contamination-free dataset or one with large quantities of contamination in an assortment of configurations, the proposed CBI method demonstrated an ability to be very competitive with the current state-of-the-art method, S1S, and the other high breakdown methods, LTS and M1S. The only study where CBI was clearly inferior to another method was versus OLS and BI regression in a contamination-free setting, which was to be expected. Yet in this study it should be noted that CBI outperformed LTS, M1S and S1S regarding scale estimation, coefficient stability and standard errors. CBI (and LTS, for that matter) was always unbiased, unlike M1S and S1S. Overall, these Monte Carlo studies demonstrated that the CBI regression procedure has significant merit and, when combined with its cluster summary analysis, has taken high breakdown regression to a new level.