

**KU LEUVEN**

**G0B16**  
**ROBUST STATISTICS**

---

**Project 2. Mallow's type M-estimator of  
(Cantoni and Ronchetti, 2001)**

---

**Dzmitry Sei (r0913246)**

**Supervised by:  
Prof. Mia Hubert  
Prof. Stefan Van Aelst**

**June 6, 2023**



**FACULTY OF SCIENCE**

# 1 Introduction

This report presents the development of a robust classifier capable of accurately predicting class labels. To achieve this, we employ a logistic regression model with Mallow's Type M-Estimator, as introduced by Cantoni and Ronchetti in 2001. By leveraging this classifier, we address the challenges posed by outliers and data discrepancies, making it suitable for datasets containing noisy or inconsistent samples.

## 2 Limitations of the Maximum Likelihood Approach

The influence function of logistic regression with the standard Maximum Likelihood Estimator (MLE) for the coefficients  $\beta$  is determined by the following equation:

$$IF(y, x, \hat{\beta}) = J(\hat{\beta})^{-1}(y - \mu(\hat{\beta}^t x))x^t$$

where  $J(\hat{\beta})$  is Fisher information matrix  $E[(y - \mu)^2 x^t x]$ . The influence function is unbounded in the design space (x) and bounded in the response space (y). While extreme values in the design space can have a significant impact on the estimator, the effect of errors in the response variable (misclassification errors) is not as clear. Studies conducted by Copas (1988) and Pregibon (1982) have demonstrated that misclassification errors can introduce bias into the MLE. Consequently, it is crucial to employ robust estimators that are less affected by deviations in both the predictor variables (x) and the responses (y)[4].

## 3 Mallow's estimator GLM approach

### 3.1 Mallow estimator

Logistic regression is a GLM for binary classification. The estimation of  $\beta$  in a generalized linear model (GLM) with non-ordinary variance is achieved through the utilization of the quasi-likelihood method.

$$\sum_{i=1}^n \left( \frac{y_i - \mu_i}{\phi v_{\mu_i}} \right) \mu_i' = 0, \quad (1)$$

where  $var(y_i) = \phi v_{\mu_i}$ . **In the context of Bernoulli logistic distribution (in our case)  $\phi = 1$ .** Non-robustness of this approach for logistic regression was shown in the previous section.

The Mallow's estimator is a robust version of 1, Mallow's estimator modifies the quasi-likelihood method by down-weighting the influence of outliers or influential observations, making it more resistant to deviations from the assumed model structure.

Given the Pearson residuals  $r_i = \frac{y_i - \mu_i}{\sqrt{\phi v_{\mu_i}}}$ , the M-estimating equations for  $\beta$  of a GLM from exponential family is given by the with the following formula

$$\sum_{i=1}^n \left[ \psi(r_i; \beta, \phi, c) \omega(x_i) r_i \frac{1}{\sqrt{\phi v_{\mu_i}}} \mu'_i - \alpha(\beta) \right] = 0, \quad (2)$$

where  $\mu'_i = \frac{\partial \mu_i}{\partial \beta}$  and  $\alpha(\beta)$  is a correction term to ensure Fisher consistency. It is computed explicitly in [2].

The function  $\psi(r_i; \beta, \phi, c)$  and the weights  $\omega(x_i)$  introduce new ingredients in comparison to the classical GLM estimators obtained through maximum quasi-likelihood, where the estimating equations 1 is derived using  $\psi(r_i; \beta, \phi, c) = r_i$  and  $\omega(x_i) = 1$  for all  $i$ .  $\psi(r_i; \beta, \phi, c)$  is added to control deviations in the y-space and the weights  $\omega(x_i)$  to down weight leverage points.

2 can be rewritten as

$$\sum_{i=1}^n \left[ \tilde{\omega}(r_i; \beta, \phi, c) \omega(x_i) r_i \frac{1}{\sqrt{\phi v_{\mu_i}}} \mu'_i - \alpha(\beta) \right], \quad (3)$$

where  $\tilde{\omega}(r_i; \beta, \phi, c) = \psi(r_i; \beta, \phi, c)/r_i$ . Weights  $\tilde{\omega}$  and  $\omega$  can be used for diagnostics of outliers and leverage points.

## 3.2 Choice of $\psi$ and $\omega(x)$

### 3.2.1 Choice $\psi$

The function  $\psi$  plays a crucial role in controlling the impact of large residuals in the estimation process. It is essential for  $\psi$  to be bounded to effectively manage these residuals. Common choices for  $\psi$  include functions that level off, such as the Huber function, or functions that exhibit a redescending. The selection of  $\psi$  is often accompanied by tuning it with a constant  $c$ . This constant is typically chosen to ensure a desired level of asymptotic efficiency in the estimation procedure.

Determining the value of  $c$  for a desired efficiency level in GLM models is more complex than in linear regression due to the influence of design and lack of general results. The estimated efficiency can be inspected afterward, allowing for model refitting with different  $c$  values if needed.

In this project Huber  $\psi$  function is used. **The value of  $c$  equal to 1.8 was empirically obtained as the optimal choice.**

### 3.2.2 Choice of $\omega(x_i)$

Robust estimators in linear models suggest a specific approach for determining the weights  $\omega(x_i)$ . One simple method is to set  $\omega(x_i) = \sqrt{(1 - h_{ii})}$ , where  $h_{ii}$  represents the leverage of observation  $x_i$ . In the current implementation of the *robustbase* package, besides the default equal weights  $\omega(x_i) = 1$  for all  $i$ , an alternative option is available. This option lets users choose weights based on the Mahalanobis distances  $d_i$  [3].

$$\omega(x_i) = \frac{1}{\sqrt{1 + 8\max(0, (d_i^2 - q)) / \sqrt{2q}}}$$

On the test data, the  $h_{ii}$  based approach shows better performance. That is why we use it during the model-building process.

### 3.3 $\beta$ distribution and significance test

According to [2]  $\hat{\beta} \sim N(\beta, \widehat{Var})$ , where  $\widehat{Var}$  is estimated using  $\hat{\beta}$ . This allows testing the significance of a coefficient estimation, using z-statistics, namely

$$z - statistic = \hat{\beta}_j / \sqrt{\widehat{Var}}. \quad (4)$$

Using 4 p-value for a coefficient can be calculated.

## 4 Robust Estimation of the Full Model

In this section, we construct a robust classifier that utilizes all the covariates. Next, we compare the performance of this estimator with logistic regression. Then, we employ the weights of residuals and data points obtained during the training of the robust classifier to diagnose robust outliers and leverage points. Only the train dataset is used for model building.

### 4.1 Model building

On the 1 the result of training is shown. As we can see models disagree about the significance of *MCoreJ* and *AreaJ*.

	Non-robust_Estimate	Non-robust_std	Non-robust_z_value	Non-robust_p	Non-robust_is_significant	Robust_Estimate	Robust_std	Robust_z_value	Robust_p	Robust_is_significant
(Intercept)	-55.87347	"0.63705"	-8.43852	"0"	"T"	-49.36688	"0.2398"	-3.2743	"0"	"T"
paperf	-0.00203	"6e-04"	-3.3885	"7e-04"	"T"	-0.00231	"0.00066"	-3.51509	"0.00044"	"T"
mpof	-0.00095	"0.00081"	-1.18554	"0.24337"	"F"	-0.00099	"0.00089"	-1.10737	"0.2684"	"F"
ncoref	"0.00387"	"0.00084"	"4.63331"	"T"	"T"	"0.00462"	"0.00093"	"4.98971"	"0"	"T"
ncore	"4e-05"	"4e-05"	"1.92147"	"0e-05"	"T"	"4e-05"	"4e-05"	"4.68839"	"0"	"T"
l2p	-0.00082	"0.00055"	-1.47984	"0.13892"	"F"	-0.00114	"0.00061"	-1.87961	"0.06016"	"F"
cffp	"0.00128"	"0.00039"	"3.84145"	"0"	"T"	"0.00121"	"0.00031"	"3.76337"	"0"	"T"
l11p	"0.00115"	"4.34108"	"1e-05"	"T"	"T"	"0.00488"	"0.00127"	"3.84779"	"0.00012"	"T"
paperj	-0.00028	"0.00041"	-0.67283	"0.50104"	"F"	-0.00026	"0.00046"	-0.60613	"0.42383"	"F"
ncorej	-0.00047	"0.00075"	-0.62849	"0.52903"	"F"	-0.00028	"0.00081"	-0.34836	"0.72757"	"F"
area	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"
area2	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"
area3	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"
area4	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"
area5	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"
area6	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"
area7	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"
area8	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"
area9	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"
area10	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"
area11	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"
area12	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"
area13	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"
area14	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"
area15	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"
area16	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"
area17	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"
area18	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"
area19	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"
area20	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"
area21	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"
area22	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"
area23	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"
area24	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"
area25	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"
area26	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"
area27	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"
area28	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"
area29	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"
area30	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"
area31	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"
area32	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"
area33	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"
area34	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"
area35	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"
area36	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"
area37	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"
area38	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"
area39	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"
area40	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"
area41	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"
area42	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"
area43	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"
area44	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"
area45	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"
area46	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"
area47	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"
area48	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"
area49	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"
area50	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"
area51	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"
area52	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"
area53	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"
area54	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"
area55	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"
area56	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"
area57	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"
area58	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"
area59	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"
area60	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"
area61	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"
area62	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"
area63	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"
area64	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"
area65	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"
area66	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"
area67	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"
area68	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"
area69	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"
area70	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"
area71	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"
area72	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"
area73	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"
area74	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"
area75	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"
area76	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"
area77	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"
area78	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"
area79	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"
area80	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"
area81	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"
area82	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"
area83	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"
area84	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"	"0.00047"	"0.00047"	"1.00000"	"0.31831"	"F"

Figure 2 depicts the fractional relationship between the non-robust and robust estimators. As we can non-robust method underestimates standard error. The observed phenomenon can be attributed to the inherent characteristics of the M-estimator and the non-robust estimator. The M-estimator, specifically designed to exhibit robustness against outliers, employs a strategy that reduces the influence of extreme observations. Consequently, this down weighting of extreme observations results in estimated coefficients with larger standard errors. The M-estimator achieves this robustness by being less reliant on individual observations and their potential impact on the estimated coefficients. In contrast, the non-robust estimator treats all observations equally, including outliers. This approach assigns equal weight to every data point, regardless of its magnitude or potential to distort the estimation. Consequently, if outliers exert a substantial impact on the estimated coefficients, the non-robust estimator can yield smaller standard errors due to the lack of down weighting.

The presence of outliers in the data significantly also impacts the observed effect of a  $AreaN$  coefficient, leading to pronounced disparities.

	Non-rob. Estimate/rob	Estimate Non-rob.	Std. Error/rob.	Std. Error
(Intercept)	0.81725631		0.8035271	
WageF	0.87126832		0.0041326	
MTGCF	0.95624617		0.0077376	
MCORF	0.83810500		0.0025718	
AreaF	0.68475559		0.8186686	
IR2F	0.71901639		0.9132517	
cstF	0.89703895		0.8891861	
e111pf	1.02566227		0.0091137	
WageF3	0.77980232		0.9270831	
MTGCF3	1.66042454		0.9188808	
MCORF3	1.17951369		0.9284968	
Area3	0.34151095		0.9078398	
IR23	0.82815775		0.8942091	
cst3	0.75044473		0.9150634	
e111p3	-3.00395984		0.9218334	
WageFN	4.55255634		0.9377782	
MTGCFN	-1.20155574		0.9214555	
MCORFN	1.03900002		0.9391429	
AreaN	1.05598202		0.9465593	
IR2N	0.97228597		0.9122235	
cstN	0.31638574		0.9299707	
e111pN				

Figure 2: Non-robust estimations/Robust estimation

## 4.2 Outlier detection

As mentioned before outlier detection is based on examining  $\tilde{\omega}(w_r)$  and  $\omega(w_x)$ . But it should be noted that there is no single approach (especially considering that  $\tilde{\omega}$  depends on tuning parameter  $c$ ).

For instance, in the context of *glmrob*, an observation may be detected as an outlier if  $|w_r w_x| < 0.0001$ . On the other hand [1] proposed an outlier detection method where an observation is flagged as an outlier if its weight, represented by  $w_r$ , falls below a specific threshold, such as 0.5.

The approach based on [1] finds 16 outliers (see the code for the exact outliers array) in the train data.3 highlights the outliers. All the dots under the 0.5 threshold are considered extreme.

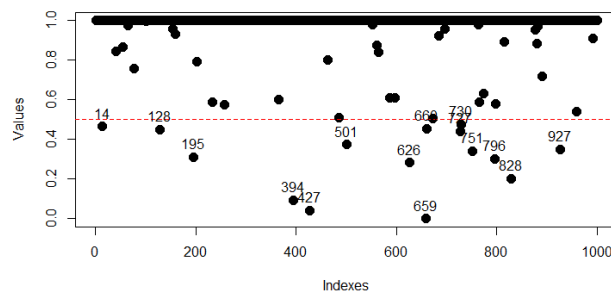


Figure 3: Full model outliers

## 5 Robust variable selection

According to [2], 2 can be seen as the derivative with respect to  $\beta$  of the robust quasi-likelihood function  $\sum Q_M(y_i, \mu_i)$ . This allows us to use the robust version of the deviance statistics.

$$\Lambda_{Q_M} = 2 \left[ \sum Q_M(y_i, \hat{\mu}_i) - \sum Q_M(y_i, \dot{\mu}_i) \right],$$

where  $\hat{\mu}_i$  and  $\dot{\mu}_i$  are the estimators under model  $M_p$  and nested model  $M_{p-q}$  respectively.

The quasi-deviance statistic can be used in an ANOVA test to compare two nested *glmrob* models.

Utilizing ANOVA we reduced the full model to the following:

$$\text{logit}(p_i) = \beta_0 + \beta_1 \text{MCoreF} + \beta_2 \text{AreaF} + \beta_3 \text{IR2F} + \beta_4 \text{csfF} + \beta_5 \text{EllipF} + \beta_6 \text{AreaJ} + \beta_7 \text{IR2J} + \beta_8 \text{csfJ} \quad (5)$$

5 (in robust and in non-robust context) is referred to as a truncated model In the text below.

## 6 Performance on Test Data

The algorithm, utilized for detecting the best-performing model and selecting the optimal cutoff point for separating stars from galaxies using predicted probabilities, is as follows:

1. Build a Receiver Operating Characteristic (ROC) curve for each available model.
2. Determine the cutoff point for each model based on the constructed ROC curve.
3. The best model is determined by selecting the one with the highest Area Under the Curve (AUC) value.

We obtained the following results.

Model	AUC	Cutoff Point
Full Robust	0.9951	0.354
Full Non-robust	0.9949	0.27
Truncated Robust	0.9948	0.284
Truncated Non-robust	0.9946	0.336

Based on the analysis, it can be observed that the Robust full model performs the best among the models, despite having similar AUC values. Furthermore, it is noteworthy that the Full robust model has the largest cutoff point compared to the other models.

The predicted probabilities for the test set are stored in the variable *rob\_preds* in the provided code.

## References

- [1] Eva Cantoni and Elvezio Ronchetti. “A robust approach for skewed and heavy-tailed outcomes in the analysis of health care expenditures”. In: *Journal of Health Economics* 25.2 (2006), pp. 198–213.
- [2] Eva Cantoni and Elvezio Ronchetti. “Robust inference for generalized linear models”. In: *Journal of the American Statistical Association* 96.455 (2001), pp. 1022–1030.
- [3] Stephane Heritier et al. *Robust methods in biostatistics*. John Wiley & Sons, 2009.
- [4] Maria-Pia Victoria-Feser. “Robust logistic regression for binomial responses”. In: *Available at SSRN 1763301* (2000).