

***Supplementary material to: Robust estimators for additive models
using backfitting***

Graciela Boente^{a*} Alejandra Martínez^a and Matías Salibián-Barrera^b

^a *Departamento de Matemática, Facultad de Ciencias Exactas y Naturales, Universidad de Buenos Aires and IMAS, CONICET, Ciudad Universitaria, Pabellón 1, 1428, Buenos Aires, Argentina;* ^b *Department of Statistics, University of British Columbia, 3182 Earth Sciences Building, 2207 Main Mall, Vancouver, BC, V6T 1Z4, Canada*

(Received 00 Month 20XX; accepted 00 Month 20XX)

To illustrate the sensitivity of the robust backfitting to the presence of single outliers, Section S.1 reports a numerical study of its empirical influence function. Some additional figures for the real data set are provided in Section S.2. In this supplement, Figures and Tables are numbered S.1, S.2, References to Sections on the main body of the paper are indicated without the capital “S”.

Keywords: Fisher-consistency; Kernel weights; Robust estimation; Smoothing

AMS Subject Classification: MSC 62G35; 62G08

*Corresponding author. Email: gboente@dm.uba.ar

S.1. Empirical Influence

A well-known measure of robustness of an estimator is given by its influence function (see Hampel *et al.* 1986). The influence function measures resistance of an estimator against infinitesimal proportions of outliers and helps study the local robustness and asymptotic efficiency of an estimator. The finite-sample version of the influence function, called the empirical influence function (Tukey, 1977), is a useful measure of sensitivity quantifying the effect of a single outlier on the estimator computed on a given sample. Although influence functions have been widely studied for many parametric models, much less attention has been paid to nonparametric estimators. To measure the influence of a contaminating point on the estimators, we follow the approach of Manchester (1996), who proposed a graphical method to display the sensitivity of a scatter plot smoother that is related to the finite-sample influence function introduced by Tukey (1977).

Given a data set $\{(\mathbf{X}_i^T, Y_i)^T\}_{1 \leq i \leq n}$ satisfying the additive model $Y = \mu_0 + \sum_{j=1}^d g_{0,j}(X_j) + \sigma_0 \varepsilon$, let $\hat{g}_{n,j}(\tau)$ be the estimator of the j -th component based on this data set evaluated at the point $\tau \in \mathbb{R}$. Assume that $\mathbf{z}_0 = (\mathbf{x}_0^T, y_0)^T$ represents a contaminating point and let $\hat{g}_{n,j}^{(\mathbf{z}_0)}(\tau)$ be the estimator based on the augmented data set $\{(\mathbf{X}_1^T, Y_1)^T, \dots, (\mathbf{X}_n^T, Y_n)^T, \mathbf{z}_0\}$ evaluated at the point τ . For a fixed value of τ , we define the empirical influence function of $\hat{g}_{n,j}(\tau)$ at \mathbf{z}_0 as the surface

$$\text{EIF}_{j,\tau}(\mathbf{z}_0) = (n+1) \left[\hat{g}_{n,j}^{(\mathbf{z}_0)}(\tau) - \hat{g}_{n,j}(\tau) \right], \quad (\text{S.1})$$

as \mathbf{z}_0 varies in $\mathbb{R}^d \times \mathbb{R}$. To explore the sensitivity of the backfitting estimators to the presence of outliers using the empirical influence function (S.1), we generated a data set of size $n = 500$ following an additive model with location $\mu_0 = 0$, additive components $g_{0,1}(x_1) = 24(x_1 - 0.5)^2 - 2$ and $g_{0,2}(x_2) = 2\pi \sin(\pi x_2) - 4$ and covariates $\mathbf{X}_i = (X_{i,1}, X_{i,2})^T \sim U([0, 1] \times [0, 1])$. The data and the regression function are shown in Figure S.1.

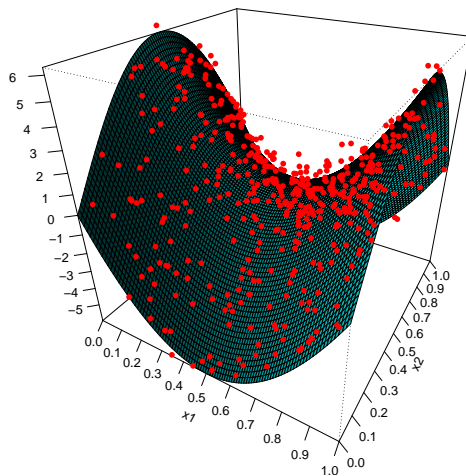


Figure S.1. Data used for the influence function study, and the corresponding regression function g_0 .

We used an Epanechnikov kernel with bandwidths $h_1 = h_2 = 0.10$, local constant smoothers ($q = 0$) and the same tuning constants as in our simulation study. We com-

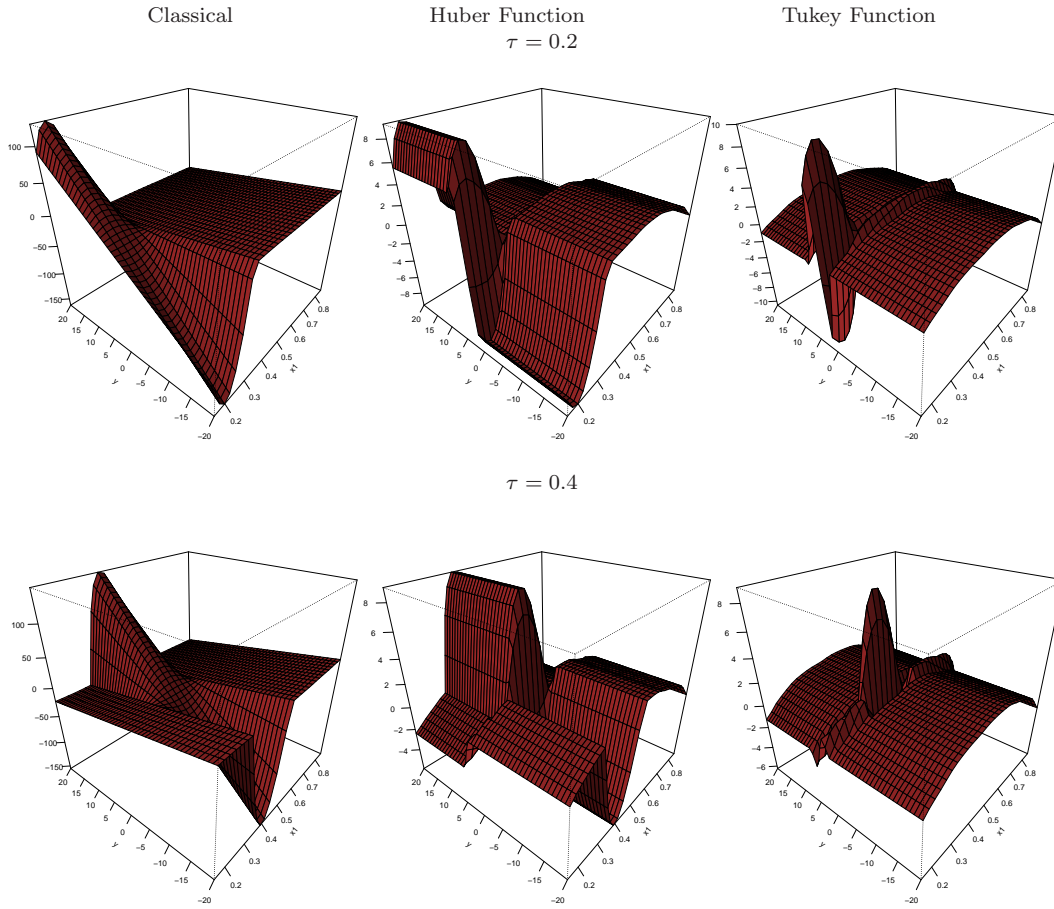


Figure S.2. Empirical influence for the classical and robust estimators, $\text{EIF}_{1,\tau}(\mathbf{x}, y)$ when $\tau = 0.2$ and 0.4 and $\mathbf{x} = (x_1, 0.5)$.

puted $\text{EIF}_{j,\tau}(\mathbf{z}_0)$ for $\tau = 0.20, 0.40, 0.60$ and 0.80 and a grid of points $\mathbf{z}_0 = ((x_1, 0.5)^T, y)^T$, where x_1 ranges over 30 equidistant points in the interval $[0.15, 0.85]$ and y takes 50 equally spaced points in $[-20, 20]$.

The results for each estimator and for $\tau = 0.2$ and 0.4 are displayed in Figure S.2, while the results for $\tau = 0.6$ and 0.8 are given in Figure S.3.

These plots illustrate the expected lack of robustness of the classical backfitting estimator, for which the empirical influence function takes very large values. Note the EIF attain the largest absolute value when x_1 is close to τ , and estimators based on Tukey's bisquare loss function have a slightly larger $|\text{EIF}|$ than those based on Huber's loss. The redescending structure of the score function can also be observed in the plot, showing that very large values of the responses have less effect on the estimator based on the Tukey loss function than in that based on the Huber loss, as noted also in the simulation study. It is important to note that, when the nonparametric regression model does not take into account an additive structure and when using a kernel with compact support to compute a kernel regression estimator only outliers near the value at which the regression function estimator is evaluated may impact the regression estimator. However, the situation is different for the backfitting method, which involves the estimation of the location parameter and an iterative algorithm involving all the residuals.

Since the absolute value of $\text{EIF}_{1,\tau}(\mathbf{x}, y)$ attains its maximum value near τ , Figure

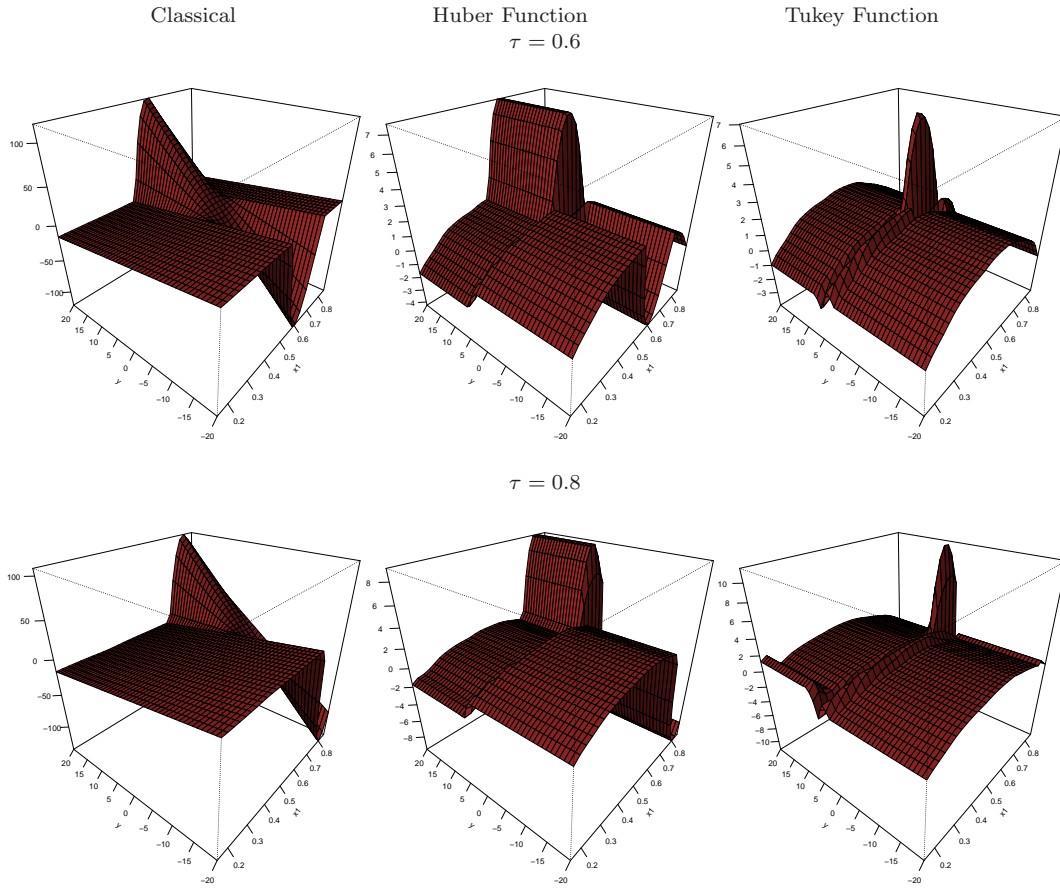


Figure S.3. Empirical influence for the classical and robust estimators, $\text{EIF}_{1,\tau}(\mathbf{x}, y)$ when $\tau = 0.6$ and 0.8 and $\mathbf{x} = (x_1, 0.5)$

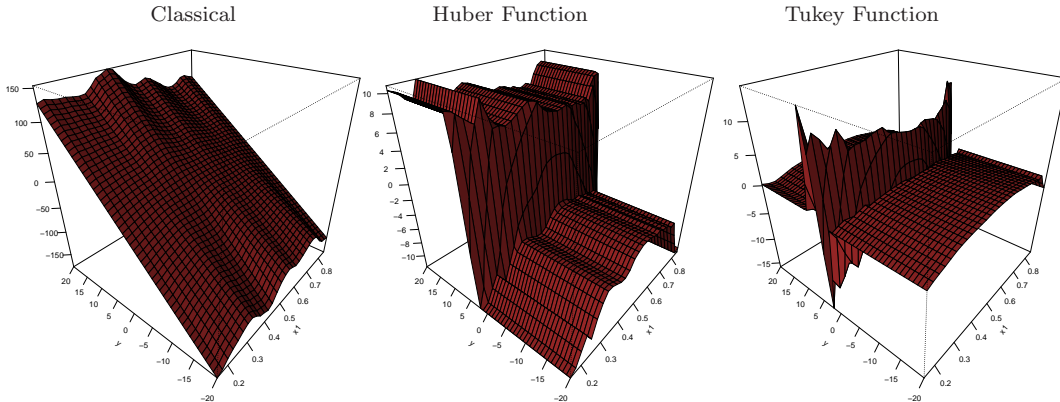


Figure S.4. Empirical influence $\text{EIF}_{1,x_1}((x_1, 0.5), y)$ for the classical and robust estimators.

S.4 shows the surfaces $\text{EIF}_{1,x_1}((x_1, 0.5), y)$, which represent the worst possible bias of these estimators in this setting. The plots of $|\text{EIF}_{1,x_1}((x_1, 0.5), y)|$ are given in Figure S.5. As expected, the bias of the classical estimators follows the size of the contaminated responses. On the other hand, the empirical functions of the robust estimators are bounded, and the most influential points correspond to x_1 near 0.2 and 0.8 , which

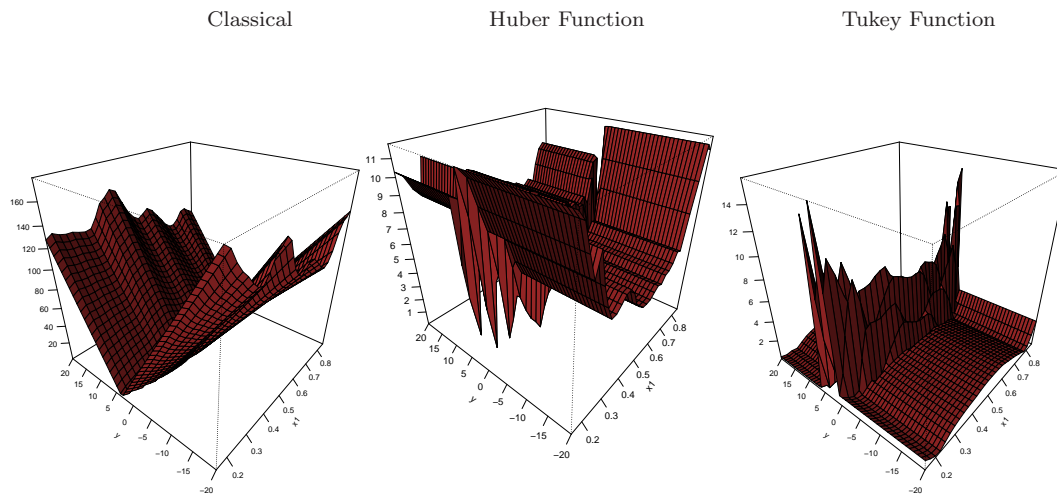


Figure S.5. Absolute value of the empirical influence, $|\text{EIF}_{1,x_1}((x_1, 0.5), y)|$ for the classical and robust estimators.

reflects the expected boundary effect. Due to the re-descending nature of the Tukey score function, the absolute value of the empirical function for larger values of y ($|y| > 5$, say) remains very low, near its minimum absolute value of 0.019.

S.2. Real data example

In this section, we complement the results obtained in Section 5, where we considered the `airquality` data set available in `R`. Figures S.6 gives the plots for the partial residuals obtained using the classical and robust estimators with all the data. On the other hand, Figure S.7 provides similar plots when using the classical estimators on the data set without the 5 detected atypical observations.

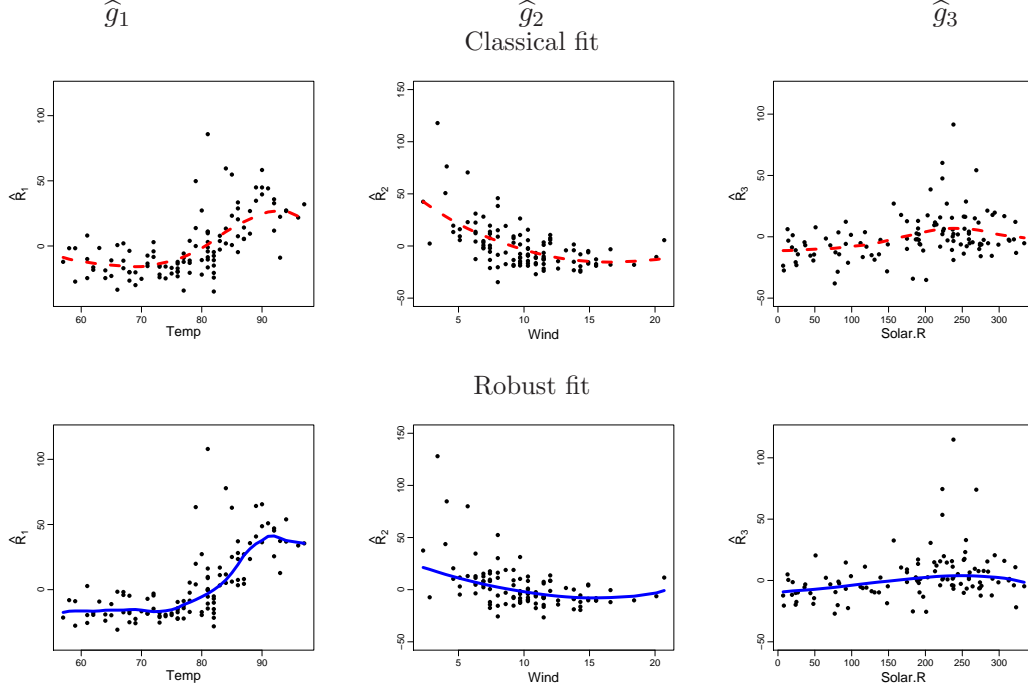


Figure S.6. Partial residuals, \hat{R}_j for $1 \leq j \leq 3$, and estimated curves for the classical (in red dashed lines) and robust (in blue solid lines) backfitting estimators with data-driven bandwidths \mathbf{h}_{LS} and \mathbf{h}_R , respectively.

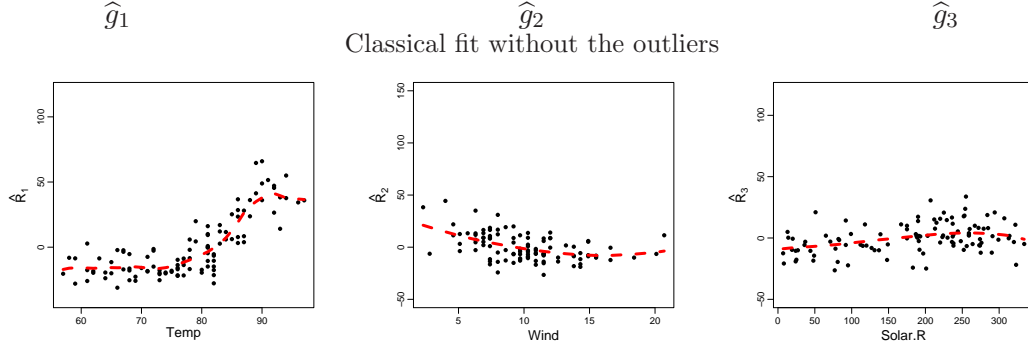


Figure S.7. Partial residuals and estimated curves for the classical backfitting estimator, $\hat{g}_j^{(-5)}$, (in red dashed lines) with data-driven bandwidth $\mathbf{h}_{LS}^{(-5)}$.

Acknowledgements. This research was partially supported by Grants PIP 112-201101-00339 from CONICET, PICT 2014-0351 from ANPCYT and 20020130100279BA from the Universidad de Buenos Aires at Buenos Aires, Argentina (G. Boente and A. Martínez) and Discovery Grant 250048-11 of the Natural Sciences and Engineering Research Council of Canada (M. Salibián Barrera). This research was begun while Alejandra Martínez was visiting the University of British

References

- Hampel, F.R., Ronchetti, E.M., Rousseeuw, P.J. and Stahel, W.A. (1986) *Robust Statistics: The Approach Based on Influence Functions*, Wiley, New York.
- Manchester, L. (1996). Empirical Influence for robust smoothing. *Australian Journal of Statistics*, **38**, 275-296.
- Tukey, J. (1977). *Exploratory Data Analysis*. Reading, MA: Addison–Wesley.