

Solutions to Voleon Problem

Yi Gao*

(Dated: December 28, 2016)

I. PATTERN OF DATA

Since the model is linear of the form

$$y_i = ax_i + b + e_i, \quad (1)$$

in which b is the intercept, and a is the linear coefficient. Meanwhile, the noise has property

$$\mathbb{E}(e_i | x_i) = 0, \quad (2)$$

and also the conditional variance of e_i is depending on x_i , i.e.,

$$\mathbb{V}(e_i | x_i) = \sigma^2(x_i). \quad (3)$$

Therefore, the linear model is heteroskedastic, and the ordinary linear least square (OLS) is no longer the best linear unbiased estimator (BLUE).

To see the pattern of the noise terms, we first use OLS to fit the data in each data set. We then use the residuals of OLS, \hat{e}_i , as an approximation of noise terms, and see the relation between noise and predictors x_i 's. We fit each data set using OLS, and regress the square of residuals on to the terms $|x_i|$ and x_i^2 ,

$$\hat{e}_i^2 = \alpha + \beta |x_i| + \gamma x_i^2. \quad (4)$$

The results are summarized in Table I.

	Value	t -stat	p -value
α	-88.6353	-1.469	0.143
β	54.4057	3.024	0.003
γ	-1.1313	-1.240	0.216

TABLE I: Regression results of OLS residuals on $|x_i|$ and x_i^2 .

According to the above regression result, it can be concluded that the model is indeed heteroskedastic, i.e., the noise term is depending on x_i . Moreover, the regression result strongly suggests the relation as $\mathbb{E}(e_i^2 | x_i) \propto |x_i|$. We now plot $\hat{e}_i/|x_i|^{1/2}$ as a function of x_i 's in Fig. 1. We can reasonably hypothesize that $e_i/|x_i|^{1/2}$ is white noise. Note that here e_i is the true noise, instead of OLS residuals.

*Electronic address: yi_gao@mfe.berkeley.edu

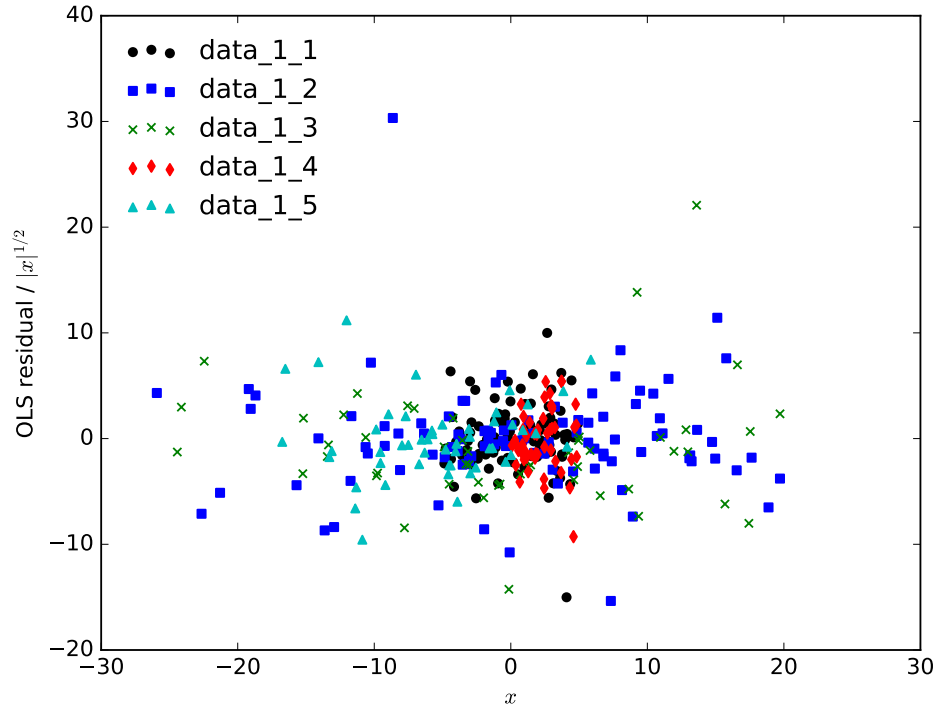


FIG. 1: $\hat{e}_i / |x_i|^{1/2}$ as a function of predictors x_i 's.