

Robust VIF Regression

D.J. Dupuis
M.-P. Victoria-Feser

G-2011-58

October 2011

Robust VIF Regression

Debbie J. Dupuis

GERAD & HEC Montréal
Montréal (Québec) Canada, H3T 2A7
debbie.dupuis@hec.ca

Maria-Pia Victoria-Feser

Research Center for Statistics
HEC Genève
Genève, Switzerland
maria-pia.victoriafeser@unige.ch

October 2011

Les Cahiers du GERAD

G-2011-58

Copyright © 2011 GERAD

Abstract

The sophisticated and automated means of data collection used by an increasing number of institutions and companies leads to extremely large datasets. Subset selection in regression is essential when a huge number of covariates can potentially explain a response variable of interest. The recent statistical literature has seen an emergence of new selection methods that provide some type of compromise between implementation (computational speed) and statistical optimality (e.g. prediction error minimization). Global methods such as Mallows' C_p have been supplanted by sequential methods such as stepwise regression. More recently, streamwise regression, faster than the former, has emerged. A recently proposed streamwise regression approach based on the variance inflation factor (VIF) is promising but its least-squares based implementation makes it susceptible to the outliers inevitable in such large data sets. This lack of robustness can lead to poor and suboptimal feature selection. This article proposes a robust VIF regression, based on fast robust estimators, that inherits all the good properties of classical VIF in the absence of outliers, but also continues to perform well in their presence where the classical approach fails. The analysis of two real data sets shows the necessity of a robust approach for policy makers.

Key Words: Variable selection; Linear regression; Multicollinearity, M -estimator.

Résumé

Un nombre croissant d'établissements et de compagnies utilisent des moyens sophistiqués et automatisés de collection de données qui mènent à des ensembles de données extrêmement grands. Lorsqu'il existe un nombre énorme de variables explicatives pour une variable réponse, le choix de sous-ensembles dans la régression est essentiel. La littérature statistique récente a vu l'apparition de nouvelles méthodes de sélection qui fournissent un certain type de compromis entre l'exécution (la vitesse de calcul) et l'optimalité statistique (par exemple, la minimisation de l'erreur de prévision). Des méthodes globales telles que le C_p de Mallows ont été supplantées par des méthodes séquentielles telles que la régression stepwise. Plus récemment, la régression streamwise, plus rapide que cette dernière, a émergé. Une approche récemment proposée de régression streamwise basée sur le variance inflation factor (VIF) est prometteuse, mais son implémentation basée sur les moindres carrés la rend susceptible aux valeurs aberrantes inévitables dans de grands ensembles de données. Ce manque de robustesse peut mener à une sélection mauvaise et sous-optimale. Cet article propose une régression VIF robuste, basée sur des estimateurs robustes rapides, qui hérite de toutes les bonnes propriétés de VIF classique dans l'absence de valeurs aberrantes, mais continue également de bien performer dans leur présence, où l'approche classique échoue. L'analyse de deux vrais jeux de données montre la nécessité d'une approche robuste pour les responsables de politiques économique, sociale et juridique.

Acknowledgments: The first author acknowledges the support of the Natural Sciences and Engineering Research Council of Canada. The second author acknowledges the support of the Swiss National Science Foundation (grant no 100014-131906).

1 Introduction

Datasets with millions of observations and a huge number of variables are now quite common, especially in business- and finance-related fields, as well as computer sciences, health sciences, etc. An important challenge is to provide statistical tools and algorithms that can be used with such datasets. In particular, for regression models, a first data analysis requires that the number of potential explanatory variables be reduced to a reasonable and tractable amount. Consider p potential explanatory variables $[1 \ x_1 \dots x_p]^T = \mathbf{x}$ and a response variable y observed on n subjects. The classical normal linear model supposes $y|\mathbf{x} \sim N(\mathbf{x}^T\boldsymbol{\beta}; \sigma^2)$ with slope parameters $\boldsymbol{\beta} = [\beta_0, \beta_1, \dots, \beta_p]^T$. Let also $\mathbf{X} = [\mathbf{1} \ \mathbf{x}_j]_{j=1, \dots, p}$ be the $n \times (p+1)$ design matrix with $\mathbf{1}$ a vector of ones. The aim is to find a subset of explanatory variables that satisfies a given criterion and such that the regression model holds.

The selection criteria are numerous and can be based on prediction, fit, etc. The available selection procedures can be broadly classified into three classes according to their general strategy and, as a result, their computational speed. A first class considers all the possible combinations of covariates as potential models, evaluates each according to a fixed criterion, and chooses the model which best suits the selected criterion. A second class is formed of sequential selection procedures in which a covariate at a time is entered in (or removed from) the model, based on a criterion that can change from one step to the next and that is computed for all potential variables to enter (or to exit) until another criterion is reached. Finally, the third class of selection procedures is also sequential in nature, but each covariate is only considered once as a potential covariate. For the first class, we find criteria such as the AIC (Akaike 1973), BIC (Schwarz 1978), Mallows' C_p (Mallows 1973), cross-validation, etc (see also Efron 2004). These methods are not adapted to large datasets since the number of potential models becomes too large and the computations are no longer feasible. For the second class, we find for example the classical stepwise regression which can be considered as a simple algorithm to compute the estimator of regression coefficients $\boldsymbol{\beta}$ that minimizes an l_q penalized sum of squared errors $\|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|_2^2 + \lambda_q \|\boldsymbol{\beta}\|_{l_q}$, with $q = 0$, i.e. $\|\boldsymbol{\beta}\|_{l_0} = \sum_{j=1}^p \eta(\beta_j \neq 0)$ (see Lin, Foster, and Ungar 2011). Fast algorithms for stepwise regressions are available, e.g. Foster and Stine (2004). Procedures for the l_1 problem are also available, e.g. Lasso/LARS (Efron, Hastie, Johnstone, and Tibshirani 2004), the Dantzig Selector (Candes and Tao 2007), or coordinate descent (Friedman, Hastie, and Tibshirani 2010). But these algorithms may also become very slow for large data sets, not only because all remaining variables are evaluated at each stage, but also because the penalty λ_q needs to be computed, and often via cross-validation. The last class is a variation of stepwise regression in which covariates are tested sequentially but only *once* for addition to the model. An example is the *streamwise regression* of Zhou, Foster, Stine, and Ungar (2006) which uses the α -investing rule (Foster and Stine 2008), is very fast, and guards against overfitting. An improved streamwise regression approach was recently proposed in Lin, Foster, and Ungar (2011) where a very fast to compute test statistic based on the variance inflation factor (VIF) of the candidate variable given the currently selected model is proposed. The approach takes into account possible multicollinearity, seeking to find the best predictive model, even if it is not the most parsimonious. Comparisons in Lin, Foster, and Ungar (2011) establish that the method performs well and is the fastest available.

Our concern in this paper is to provide model selection tools for the regression model that are robust to small model deviations. As argued in Dupuis and Victoria-Feser (2011) (see also Ronchetti and Staudte 1994), spurious model deviations such as outliers can lead to a completely different, and suboptimal, selected model when a non robust criterion, like Mallows' C_p or the VIF, is used. This happens because under slight data contamination, the estimated model parameters, using for example the least squares estimator (LS), and consequently the model choice criterion can be seriously biased. The consequence is that when the estimated criteria are compared to an absolute level (like a quantile of the χ^2 distribution), the decisions are taken at the wrong level. For the first class of selection procedures, robust criteria have been proposed such as the robust AIC of Ronchetti (1982), the robust BIC of Machado (1993), the robust Mallows' C_p of Ronchetti and Staudte (1994) and a robust criterion based on cross-validation (CV) in Ronchetti, Field, and Blanchard (1997). Since standard robust estimators are impossible to compute when the number of covariates is too large, Dupuis and Victoria-Feser (2011) proposed the use of a forward search procedure together with adjusted robust estimators when there is a large number of potential covariates. Their selection procedure, called Fast Robust Forward Selection (FRFS), falls in the second class of selection procedures. FRFS outperforms classical approaches

such as Lasso/LARS when data contamination is present and outperforms, in all studied instances, a robust version of the LARS algorithm proposed by Khan, Van Aelst, and Zamar (2007).

However, although FRFS is indeed very fast and robust, it too can become quite slow when the number of potential covariates is very large as *all* covariates are reconsidered after one is selected for entry in the model. It is therefore important to have a robust selection procedure in the streamwise regression class so that very large datasets can be analyzed in a robust fashion. In this paper, we develop a robust VIF approach that is fast, very efficient, and clearly outperforms non-robust VIF in the presence of outliers.

The remainder of the paper is organized as follows. In Section 2, we review the classical VIF approach and present our robust VIF approach. A simulation study in Section 3 shows the good performance of the new approach. In Section 4, we consider two real data sets and show how policy makers are better served by robust VIF regression than by classical VIF or Lasso. Section 5 contains a few closing remarks.

2 Robust VIF Regression

2.1 The classical approach

Lin, Foster, and Ungar (2011) propose a procedure that allows one to sweep through all available covariates and to enter those that can reduce a statistically sufficient part of the variance in the predictive model. Let \mathbf{X}_S be the design matrix that includes the selected variables at a given stage, and $\tilde{\mathbf{X}}_S = [\mathbf{X}_S \ \mathbf{z}_j]$ with \mathbf{z}_j the new potential covariate to be considered for inclusion. Without loss of generality, we suppose all variables have been standardized. Consider the following two models

$$\mathbf{y} = \mathbf{X}_S \boldsymbol{\beta}_S + \mathbf{z}_j \beta_j + \boldsymbol{\varepsilon}_{step}, \quad \boldsymbol{\varepsilon}_{step} \sim N(\mathbf{0}, \sigma_{step}^2 \mathbf{I}) \quad (1)$$

$$\mathbf{r}_S = \mathbf{z}_j \gamma_j + \boldsymbol{\varepsilon}_{stage}, \quad \boldsymbol{\varepsilon}_{stage} \sim N(\mathbf{0}, \sigma_{stage}^2 \mathbf{I}). \quad (2)$$

where $\mathbf{r}_S = (\mathbf{I} - \mathbf{X}_S (\mathbf{X}_S^T \mathbf{X}_S)^{-1} \mathbf{X}_S^T) \mathbf{y}$ are the residuals of the projection of \mathbf{y} on \mathbf{X}_S . All known estimators of the parameters β_j , σ_{step}^2 and γ_j , σ_{stage}^2 will provide different estimates when the covariates present some degree of multicollinearity, and consequently, significance tests based on estimates of β_j or γ_j do not necessarily lead to the same conclusions. While in stepwise regression the significance of β_j in model (1) is at the core of the selection procedure, in streamwise regression, one estimates more conveniently γ_j . Lin, Foster, and Ungar (2011) show that, when LS are used to estimate, $\hat{\gamma}_j = \rho \hat{\beta}_j$ where $\rho = \mathbf{z}_j^T (\mathbf{I} - \mathbf{X}_S (\mathbf{X}_S^T \mathbf{X}_S)^{-1} \mathbf{X}_S^T) \mathbf{z}_j$. They then compare $T_\gamma = \hat{\gamma}_j / (\rho^{1/2} \sigma)$, with suitable estimates for ρ and σ , to the standard normal distribution to decide whether or not \mathbf{z}_j should be added to the current model. The procedure is called VIF regression since Marquardt (1970) called $1/\rho$ the VIF for \mathbf{z}_j .

2.2 A robust weighted slope estimator

Since the test statistic T_γ is based on : (1) the LS estimator $\hat{\gamma}_j$, (2) ρ , in-turn based on the design matrix \mathbf{X}_S and \mathbf{z}_j , and (3) the classical estimator of σ , it is obviously very sensitive to outliers, a form of model deviation. An extreme response or a very badly placed design point can have a drastic effect on T_γ . The latter is then compared to the null distribution : the correct asymptotic distribution under the hypothesis that the regression model holds. With model deviations, the null distribution is not valid, and hence selection decisions (to add the covariate or not) are taken rather arbitrarily. We propose here to limit the influence of extreme observations by considering weighted LS estimators of the form

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^{wT} \mathbf{X}^w)^{-1} \mathbf{X}^{wT} \mathbf{y}^w \quad (3)$$

with $\mathbf{X}^w = \text{diag}(\sqrt{w_i^0}) \mathbf{X}$ and $\mathbf{y}^w = \text{diag}(\sqrt{w_i^0}) \mathbf{y}$. The weights w_i^0 depend on the data and are such that extreme observations in the response and/or in the design have a nil or limited effect on the value of $\hat{\boldsymbol{\beta}}$. Dupuis and Victoria-Feser (2011) propose Tukey's redescending biweight weights

$$w_i(r_i; c) = \begin{cases} \left(\left(\frac{r_i}{c} \right)^2 - 1 \right)^2 & \text{if } |r_i| \leq c, \\ 0 & \text{if } |r_i| > c, \end{cases} \quad (4)$$

where $r_i = (y_i - \mathbf{x}_i^T \boldsymbol{\beta})/\sigma$ are standardized residuals that are computed in practice for chosen estimators of $\boldsymbol{\beta}$ and σ (see below). The constant c controls the efficiency and the robustness of the estimator. Indeed, the most efficient estimator is the LS estimator, i.e. (3) with all weights equal to one (i.e. $c \rightarrow \infty$), but it is very sensitive to (small) model deviations, while a less efficient but more robust estimator is obtained by downweighting observations that have a large influence on the estimator, i.e. by setting $c < \infty$ in (4). The value $c = 4.685$ corresponds to an efficiency level of 95% for the robust estimator compared to the LS estimator at the normal model and is the value used throughout the paper.

We follow Dupuis and Victoria-Feser (2011) and use for the weights $w_i^0 = w_i(r_i^0; c)$ in (3), where the residuals $r_i^0 = (y_i - \mathbf{x}_i^T \hat{\boldsymbol{\beta}}^0)/\hat{\sigma}^0$ and $\hat{\sigma}^0 = 1.483 \text{med}|\tilde{r}_i^0 - \text{med}(\tilde{r}_i^0)|$, the median absolute deviation (MAD) of the residuals $\tilde{r}_i^0 = y_i - \mathbf{x}_i^T \hat{\boldsymbol{\beta}}^0$. The slope estimates are $\hat{\boldsymbol{\beta}}^0 = [(\mathbf{X}_0^w)^T \mathbf{X}_0^w]^{-1} (\mathbf{X}_0^{w2})^T \mathbf{y}$, with $\mathbf{X}_0^w = [1 \ \sqrt{w_{i1}}x_{i1} \ \dots \ \sqrt{w_{ip}}x_{ip}]$ and $\mathbf{X}_0^{w2} = [1 \ w_{i1}x_{i1} \ \dots \ w_{ip}x_{ip}]$, $i = 1, \dots, n$, with weights $w_{ij}, \forall j = 1, \dots, p$, computed using (4) at the residuals $r_{ij} = (y_i - \hat{\beta}_{0j} - x_{ij}\hat{\beta}_j)/\hat{\sigma}_j$, with $\hat{\sigma}_j = \text{MAD}(y_i - \hat{\beta}_{0j} - x_{ij}\hat{\beta}_j)$. The slope estimators $\hat{\beta}_1, \dots, \hat{\beta}_p$ and the intercept estimators $\hat{\beta}_{01}, \dots, \hat{\beta}_{0p}$ are computed on the p marginal models $y = \beta_{01} + x_1\beta_1 + \varepsilon_1, \dots, y = \beta_{0p} + x_p\beta_p + \varepsilon_p$ using a robust weighted estimator defined implicitly through

$$\sum_{i=1}^n w_i(r_i; c) r_i \mathbf{x}_i = 0. \quad (5)$$

Here we actually propose to consider Huber's weights given for the regression model by

$$w_i(r_i; c) = \min \left\{ 1; \frac{c}{\|\mathbf{r}_i \mathbf{x}_i\|} \right\} \quad (6)$$

with $c = 1.345$. Estimators in (5) belong to the class of M -estimators (Huber 1964, 1967). With (6) in (5), the marginal intercepts and slope estimators are simpler (and faster) to compute than the ones based on Tukey's biweight weights as originally proposed in Dupuis and Victoria-Feser (2011). For the scale in the weights in (5), we propose to use the MAD of the residuals.

The estimator in (3) is a one-step estimator that is actually biased when there is multicollinearity in the covariates. Dupuis and Victoria-Feser (2011) show that the bias can be made smaller and even nil if $\hat{\boldsymbol{\beta}} = \hat{\boldsymbol{\beta}}^1$ is iterated further to get say $\hat{\boldsymbol{\beta}}^k$ computed at the updated weights w_i^1, \dots, w_i^{k-1} based on the residuals $r_i^{(1)} = (y_i - \mathbf{x}_i^T \hat{\boldsymbol{\beta}}^{(1)})/\hat{\sigma}^{(1)}, \dots, r_i^{(k-1)} = (y_i - \mathbf{x}_i^T \hat{\boldsymbol{\beta}}^{(k-1)})/\hat{\sigma}^{(k-1)}$. In the simulation study in Section 3 however, we find that the bias is very small even with relatively large multicollinearity, so that in practice there is often no need to proceed with this iterative correction.

2.3 Robust VIF selection criterion

Let $\mathbf{X}_S^w = \text{diag}(\sqrt{w_{iS}^0}) \mathbf{X}_S$ be the weighted design matrix at stage S with say q columns (hence $q - 1$ covariates) and $\mathbf{z}_j^w = \text{diag}(\sqrt{w_{ij}^w}) \mathbf{z}_j$ the new candidate covariate that is evaluated at the current stage $S + 1$. One could use the weights w_{iS}^0 for \mathbf{z}_j^w instead of the weights w_{ij} computed at the marginal models with only \mathbf{z}_j as covariate, but this would require more computational time. The simulation results in Section 3 show that one gets very satisfactory results with w_{ij} . Let also $\tilde{\mathbf{X}}_S^w = [\mathbf{X}_S^w | \mathbf{z}_j^w]$ and define $\hat{\beta}_j^w$ as the last element of the vector $[\tilde{\mathbf{X}}_S^{wT} \tilde{\mathbf{X}}_S^w]^{-1} \tilde{\mathbf{X}}_S^{wT} \mathbf{y}^w$ with $\mathbf{y}^w = \text{diag}(\sqrt{w_{iS}^0}) \mathbf{y}$. $\hat{\beta}_j^w$ is actually a robust estimator of β_j in (1). Let $\mathbf{H}_S^w = \mathbf{X}_S^w (\mathbf{X}_S^{wT} \mathbf{X}_S^w)^{-1} \mathbf{X}_S^{wT}$ and $\hat{\boldsymbol{\beta}}_S = (\mathbf{X}_S^{wT} \mathbf{X}_S^w)^{-1} \mathbf{X}_S^{wT} \mathbf{y}^w$, then

$$\begin{aligned} \hat{\beta}_j^w &= -(\mathbf{z}_j^{wT} \mathbf{z}_j^w - \mathbf{z}_j^{wT} \mathbf{H}_S^w \mathbf{z}_j^w)^{-1} \mathbf{z}_j^{wT} \mathbf{X}_S^w (\mathbf{X}_S^{wT} \mathbf{X}_S^w)^{-1} \mathbf{X}_S^{wT} \mathbf{y}^w \\ &\quad + (\mathbf{z}_j^{wT} \mathbf{z}_j^w - \mathbf{z}_j^{wT} \mathbf{H}_S^w \mathbf{z}_j^w)^{-1} \mathbf{z}_j^{wT} \mathbf{y}^w \\ &= (\mathbf{z}_j^{wT} \mathbf{z}_j^w - \mathbf{z}_j^{wT} \mathbf{H}_S^w \mathbf{z}_j^w)^{-1} \mathbf{z}_j^{wT} (\mathbf{y}^w - \mathbf{X}_S^w \hat{\boldsymbol{\beta}}_S) \\ &= (\mathbf{z}_j^{wT} \mathbf{z}_j^w - \mathbf{z}_j^{wT} \mathbf{H}_S^w \mathbf{z}_j^w)^{-1} \mathbf{z}_j^{wT} \mathbf{r}_S^w \\ &= (\mathbf{z}_j^{wT} \mathbf{z}_j^w - \mathbf{z}_j^{wT} \mathbf{H}_S^w \mathbf{z}_j^w)^{-1} (\mathbf{z}_j^{wT} \mathbf{z}_j^w) (\mathbf{z}_j^{wT} \mathbf{z}_j^w)^{-1} \mathbf{z}_j^{wT} \mathbf{r}_S^w \end{aligned}$$

where \mathbf{r}_S^w are the residuals of the weighted fit of \mathbf{y}^w on \mathbf{X}_S^w . Let

$$\rho^w = (\mathbf{z}_j^{wT} \mathbf{z}_j^w - \mathbf{z}_j^{wT} \mathbf{H}_S^w \mathbf{z}_j^w) (\mathbf{z}_j^{wT} \mathbf{z}_j^w)^{-1}$$

then

$$\widehat{\beta}_j^w = (\rho^w)^{-1} \widehat{\gamma}_j^w$$

with $\widehat{\gamma}_j^w = (\mathbf{z}_j^{wT} \mathbf{z}_j^w)^{-1} \mathbf{z}_j^{wT} \mathbf{r}_S^w$, i.e. the weighted estimator of the fit of \mathbf{z}_j^w on the weighted residuals \mathbf{r}_S^w , i.e. model (2). Note however that $\widehat{\beta}_j^w$ is not equal to the last element of $\widehat{\beta}_{S+1}^w$ unless the weights w_{iS}^0 are used for \mathbf{z}_j^w . Note also that we can write

$$\rho^w = 1 - R_{jS}^{w2}$$

with

$$R_{jS}^{w2} = \mathbf{z}_j^{wT} \mathbf{H}_S^w \mathbf{z}_j^w (\mathbf{z}_j^{wT} \mathbf{z}_j^w)^{-1} \quad (7)$$

a robust estimate of the coefficient of determination R^2 . Renaud and Victoria-Feser (2010) propose a robust R^2 based on weighted responses and covariates and (7) is equivalent to their proposal (with $a = 1$, see their Theorem 1) but with other sets of weights. Moreover, ρ^w is the partial variance of \mathbf{z}_j^w given \mathbf{X}_S^w (see Dupuis and Victoria-Feser 2011). Lin, Foster, and Ungar (2011) note that using all the data to compute ρ (in the classical setting) is quite computationally expensive and they propose a subsampling approach. For the same reason, we also propose to actually estimate ρ^w by computing (7) on a randomly chosen subset of size $m = 200$.

To derive the t -statistic based on $\widehat{\gamma}_j^w$, we follow Lin, Foster, and Ungar (2011) who base their comparison on the expected value of the estimated variance of respectively $\widehat{\beta}_j^w$ and $\widehat{\gamma}_j^w$. Let $\widehat{\sigma}_{step}^2$ and $\widehat{\sigma}_{stage}^2$ be respectively robust residual variance estimates for models (1) and (2). Let also $\mathbf{A}_{(i)(j)}$ denote the element (i, j) of matrix \mathbf{A} . For $\widehat{\beta}_j^w$, supposing that $w_{ij}/w_i^0 \approx 1$, we can use

$$\begin{aligned} \widehat{\text{Var}}(\widehat{\beta}_j^w) &\approx \widehat{\sigma}_{step}^2 \left[\widetilde{\mathbf{X}}_S^{wT} \widetilde{\mathbf{X}}_S^w \right]_{(q+1)(q+1)}^{-1} e_c^{-1} \\ &= \widehat{\sigma}_{step}^2 (\mathbf{z}_j^{wT} \mathbf{z}_j^w - \mathbf{z}_j^{wT} \mathbf{H}_S^w \mathbf{z}_j^w)^{-1} e_c^{-1} \\ &= \widehat{\sigma}_{step}^2 (\rho^w)^{-1} (\mathbf{z}_j^{wT} \mathbf{z}_j^w)^{-1} e_c^{-1} \\ &= \frac{\widehat{\sigma}_{step}^2}{n} (\rho^w)^{-1} \left(\frac{1}{n} \sum_i (z_{ij}^w)^2 \right)^{-1} e_c^{-1} \end{aligned}$$

with e_c given by equation (3.20) in Heritier et al. (2009). For $\widehat{\gamma}_j^w$, based on the model with \mathbf{r}_S^w as the response and \mathbf{z}_j^w as the explanatory variable (without intercept) we have

$$\begin{aligned} \widehat{\text{Var}}(\widehat{\gamma}_j^w) &\approx \widehat{\sigma}_{stage}^2 (\mathbf{z}_j^{wT} \mathbf{z}_j^w)^{-1} \tilde{e}_c^{-1} \\ &= \frac{\widehat{\sigma}_{stage}^2}{n} \left(\frac{1}{n} \sum_i (z_{ij}^w)^2 \right)^{-1} \tilde{e}_c^{-1} \end{aligned}$$

with \tilde{e}_c^{-1} the efficiency of a robust slope estimator computed using Huber's weights relative to the LS, which is not equal to e_c^{-1} , the efficiency of a robust slope estimator computed using Tukey's weights relative to the LS. We will see below that the computation of the former is not needed. Hence, approximating $\widehat{\sigma}_{step}^2 \approx \widehat{\sigma}_{stage}^2 = \widehat{\sigma}^2$, we have

$$\widehat{\text{Var}}(\widehat{\beta}_j^w) \approx (\rho^w)^{-1} \widehat{\text{Var}}(\widehat{\gamma}_j^w) (e_c / \tilde{e}_c)^{-1}.$$

An honest approximate robust test statistic T_w is then given

$$\frac{\widehat{\beta}_j^w}{\sqrt{\widehat{\text{Var}}(\widehat{\beta}_j^w)}} \approx \frac{(\rho^w)^{-1} \widehat{\gamma}_j^w}{\sqrt{(\rho^w)^{-1} \widehat{\text{Var}}(\widehat{\gamma}_j^w) (e_c / \tilde{e}_c)^{-1}}}$$

i.e.

$$T_w = (\rho^w)^{-1/2} \frac{\hat{\gamma}_j^w}{\sqrt{\frac{\hat{\sigma}^2}{n} \left(\frac{1}{n} \sum_i z_{ij}^w\right)^{-1} e_c^{-1}}} \quad (8)$$

with $\hat{\sigma}^2$ a robust mean squared error for the model with \mathbf{r}_S^w as response and \mathbf{z}_j^w as explanatory variable (i.e. model (2)). We use $\hat{\sigma} = \text{MAD}(\mathbf{r}_S^w - \mathbf{z}_j^w (\mathbf{z}_j^{wT} \mathbf{z}_j^w)^{-1} \mathbf{z}_j^{wT} \mathbf{r}_S^w)$.

Our *fast robust evaluation procedure* is summarized by the following five steps. Suppose that we are at stage S and a set of $q - 1$ covariates has been chosen in the model. We are considering covariate \mathbf{z}_j for possible entry. We are working with $c = 4.685$ and have computed e_c and the weights w_{ij} and w_{iS}^0 .

1. Obtain the residuals $\mathbf{r}_S^w = \mathbf{y}^w - \mathbf{X}_S^w (\mathbf{X}_S^{wT} \mathbf{X}_S^w)^{-1} \mathbf{X}_S^{wT} \mathbf{y}^w$.
2. Set $\mathbf{z}_j^w = \text{diag}(\sqrt{w_{ij}}) \mathbf{z}_j$. Compute $\hat{\gamma}_j^w = (\mathbf{z}_j^{wT} \mathbf{z}_j^w)^{-1} \mathbf{z}_j^{wT} \mathbf{r}_S^w$ and $\hat{\sigma} = \text{MAD}(\mathbf{r}_S^w - \mathbf{z}_j^w (\mathbf{z}_j^{wT} \mathbf{z}_j^w)^{-1} \mathbf{z}_j^{wT} \mathbf{r}_S^w)$.
3. Sample a small subset $\mathcal{I} = \{i_1, \dots, i_m\} \in \{1, \dots, n\}$ of the observations and let $\mathcal{I}\mathbf{x}$ denote the corresponding subsample from the regressor \mathbf{x} .
4. Let $\mathcal{I}\mathbf{H}_S^w = \mathcal{I}\mathbf{X}_S^w (\mathcal{I}\mathbf{X}_S^{wT} \mathcal{I}\mathbf{X}_S^w)^{-1} \mathcal{I}\mathbf{X}_S^{wT}$, compute $R_{jS}^{w2} = \mathcal{I}\mathbf{z}_j^{wT} \mathcal{I}\mathbf{H}_S^w \mathcal{I}\mathbf{z}_j^w (\mathcal{I}\mathbf{z}_j^{wT} \mathcal{I}\mathbf{z}_j^w)^{-1}$, and find $\rho^w = 1 - R_{jS}^{w2}$.
5. Compute the approximate t -ratio $T_w = (\rho^w)^{-1/2} \hat{\gamma}_j^w / \sqrt{\hat{\sigma}^2 \left(\sum_i z_{ij}^w\right)^{-1} e_c^{-1}}$ and compare it to an adapted quantile to decide to add or not \mathbf{z}_j to the current set.

A more detailed algorithm in which the decision rule (add or not the new variable) is also specified is given in the Appendix.

2.4 Comparison with the robust t -statistic of FRFS

The t -statistic proposed by Dupuis and Victoria-Feser (2011) (equation (5)) and used to test whether a candidate covariate is entered in the current model, can be written as

$$T^2 = \frac{1}{\sigma^2 \rho^w} \frac{n}{\sum w_{ij}} e_c \mathbf{y}_j^{wT} \mathbf{z}_j^w (\mathbf{z}_j^{wT} \mathbf{z}_j^w)^{-1} \mathbf{z}_j^{wT} (\mathbf{I} - \mathbf{H}_S^w) \mathbf{y}_j^w$$

with $\mathbf{y}_j^w = \text{diag}(\sqrt{w_{ij}}) \mathbf{y}$. Supposing that $\mathbf{y}_j^w \approx \mathbf{y}^w$ and $n / \sum w_{ij} \approx 1$, then

$$\begin{aligned} T^2 &\approx \frac{1}{\sigma^2 \rho^w} e_c \mathbf{y}^{wT} \mathbf{z}_j^w \hat{\gamma}_j^w \\ &= \frac{(\hat{\gamma}_j^w)^2}{\sigma^2 \rho^w (\mathbf{z}_j^{wT} \mathbf{z}_j^w)^{-1}} e_c \frac{1}{\hat{\gamma}_j^w} \mathbf{y}^{wT} \mathbf{z}_j^w (\mathbf{z}_j^{wT} \mathbf{z}_j^w)^{-1} \\ &= \frac{(\hat{\gamma}_j^w)^2}{\sigma^2 \rho^w (\mathbf{z}_j^{wT} \mathbf{z}_j^w)^{-1}} e_c \frac{\mathbf{y}^{wT} \mathbf{z}_j^w}{\mathbf{z}_j^{wT} (\mathbf{I} - \mathbf{H}_S^w) \mathbf{y}^w} \end{aligned} \quad (9)$$

Hence, T_w in (8) and T in (9) differ by a multiplicative factor of

$$\kappa = \sqrt{\frac{\mathbf{y}_j^{wT} \mathbf{z}_j^w}{\mathbf{z}_j^{wT} (\mathbf{I} - \mathbf{H}_S^w) \mathbf{y}_j^w}}$$

which is the square root of the ratio of the robustly estimated covariance between \mathbf{z}_j and \mathbf{y} , and the robustly estimated partial covariance between \mathbf{z}_j and \mathbf{y} given \mathbf{X}_S . One can notice that in the orthogonal case (and standardized covariates), we have $\mathbf{z}_j^{wT} \mathbf{H}_S^w \approx 0$ so that $\kappa \approx 1$.

3 Simulation Study

We carry out a simulation study to assess the effectiveness of the model selection approaches outlined above. First, we create a linear model

$$\mathbf{y} = X_1 + X_2 + \dots + X_k + \sigma \varepsilon \quad (10)$$

where X_1, X_2, \dots, X_k are multivariate normal (MVN) with $E(X_i) = 0$, $\text{Var}(X_i) = 1$, and $\text{corr}(X_i, X_j) = \theta$, $i \neq j, i, j = 1, \dots, k$, and ε an independent standard normal variable. We choose θ to produce a range of theoretical $R^2 = (\text{Var}(y) - \sigma^2)/\text{Var}(y)$ values for (10) and σ to give t values for our target regressors of about 6 under normality as in Ronchetti et al. (1997). The covariates X_1, \dots, X_k are our k target covariates. Let e_{k+1}, \dots, e_p be independent standard normal variables and use the first $2k$ to give the $2k$ covariates

$$\begin{aligned} X_{k+1} &= X_1 + \lambda e_{k+1}, & X_{k+2} &= X_1 + \lambda e_{k+2}, \\ X_{k+3} &= X_2 + \lambda e_{k+3}, & X_{k+4} &= X_2 + \lambda e_{k+4}, \\ &\vdots & & \\ X_{3k-1} &= X_k + \lambda e_{3k-1}, & X_{3k} &= X_k + \lambda e_{3k}; \end{aligned}$$

and the final $p - 3k$ to give the $p - 3k$ covariates

$$X_i = e_i, \quad i = 3k + 1, \dots, p.$$

Variables X_{k+1}, \dots, X_{3k} are noise covariates that are correlated with our target covariates, and variables X_{3k+1}, \dots, X_p are independent noise covariates. Note that the covariates X_1, \dots, X_p are then relabeled with a random permutation of $1 : p$ so that the target covariates do not appear in position $1 : k$, but rather in arbitrary positions. This is necessary to test the effectiveness of the streamwise variable selection as covariates considered early on are favored when many covariates are correlated.

We consider samples without and with contamination. Samples with no contamination are generated using $\varepsilon \sim N(0, 1)$. To allow for 5% outliers, we generate using $\varepsilon \sim 95\%N(0, 1) + 5\%N(30, 1)$. These contaminated cases also have high leverage X -values: $X_1, \dots, X_k \sim \text{MVN}$ as before, except $\text{Var}(X_i) = 5$, $i = 1, \dots, k$. This represents the most difficult contamination scheme: large residuals at high leverage points. We choose $\lambda = 3.18$ so that $\text{corr}(X_1, X_{k+1}) = \text{corr}(X_1, X_{k+2}) = \text{corr}(X_2, X_{k+3}) = \dots = \text{corr}(X_k, X_{3k}) = 0.3$.

In all simulations, we simulated n independent samples, with or without contamination, to use for variable selection. Then, another n independent samples without contamination were simulated for out-of-sample performance testing. The out-of-sample performance was evaluated using the mean sum of squared errors (MSE), $\sum_{i=n+1}^{2n} (y_i - \mathbf{x}_i^T \hat{\boldsymbol{\beta}})^2 / n$ where $\hat{\boldsymbol{\beta}}$ is the estimated coefficient determined by the classical and robust VIF regression selection procedures applied to the training set. Because the true predictors are known, we also compute the out-of-sample performance measure using the true $\boldsymbol{\beta}$. Classical VIF selection was carried out using the VIF package for R and default argument settings. Robust VIF was also implemented in R.

Simulation results for $n = 1000$, $k = 5$, and $p = 100$ and $p = 1000$, are presented in Table 1 and Figures 1 and 2, respectively. Entries in the top panel of the table give the percentage of runs falling into each category. The category ‘‘Correct’’ means that the correct model was chosen. ‘‘Extra’’ means that a model was chosen for which the true model is a proper subset. ‘‘Missing 1’’ means that the model chosen differed from the true model only in that it was missing one of the target covariates; ‘‘Missing 2’’ and ‘‘Missing 3’’ are defined analogously. The Monte Carlo standard deviation of entries is bounded by 3.5%. We also report the empirical marginal false discovery rate (mFDR) $\widehat{\text{mFDR}} = \widehat{E(V)} / (\widehat{E(V)} + \widehat{E(S)} + \eta)$ where $\widehat{E(S)}$ is the average number of true discoveries, $\widehat{E(V)}$ is the average number of false discoveries and $\eta = 10$ is selected following Lin, Foster, and Ungar (2011). We also report the required computation time. Note the particularly frugal robust approach : the cost of robustness is no more than a doubling of the computation time.

Both algorithms do not perform well in terms of the proportion of correctly selected models. Both algorithms do however choose a model for which the true model is a subset when there are no outliers. The classical VIF approach fails miserably in the presence of outliers, while the robust VIF approach is only slightly affected by the presence of outliers.

As the simulated data sets have noise covariates that are correlated with target covariates, the poor performance in terms of %Correct is expected given the streamwise approach of VIF regressions. The FRFS approach in Dupuis and Victoria-Feser (2011) chose the correct model in at least 80% of the cases in similar data sets. But as pointed out by Lin, Foster, and Ungar (2011), the goal here is different : good fast out-of-sample prediction. The streamwise approach is fast and the main purpose of an α -investing control is to avoid

Table 1: Model selection results. Simulated data, as described in §3, have $n = 1000$ observations with $p = 100$ and $p = 1000$ potential regressors, including $k = 5$ target regressors. Correlation among target regressors is $\theta = 0.1$ ($R^2 = 0.20$), $\theta = 0.53$ ($R^2 = 0.50$) and $\theta = 0.85$ ($R^2 = 0.80$). Correlation among each target regressor and two other regressors is 0.3 in all cases. Remaining regressors are uncorrelated. Methods are classical (C) and robust (R) VIF regression. Table entries are % of cases in categories listed in first column. Empirical mFDR appears in the second to last row. Mean execution times (in seconds) appear in the last row. Results are based on 200 simulations for each configuration.

	$p = 100$											
	$R^2 = 0.20$				$R^2 = 0.50$				$R^2 = 0.80$			
	No outliers		5% outliers		No outliers		5% outliers		No outliers		5% outliers	
	C	R	C	R	C	R	C	R	C	R	C	R
%Correct	13.5	33	0	20	6	12	0	10.5	11.5	18.5	0	15
%Extra	83.5	58	0	65.0	92	83	0	80	86.5	76.5	0	73.5
%Missing 1	1.5	3.5	0	6.5	1	1	0	2.0	0.5	1	0	3
%Missing 2	0	0.5	1	0.5	0	0	1.5	0	0	0	1.5	0
%Missing 3	0	0	2	0	0	0	9.5	0	0	0	11	0
%Other	1.5	5	97	8	1	4	89	7.5	1.5	4	87.5	8.5
%mFDR	11.0	6.3	6.1	9.3	17.9	14.2	7.5	14.7	16.1	13.2	10.7	13.8
Time	0.63	1.11	0.54	1.09	0.65	1.12	0.59	1.18	0.69	1.20	0.59	1.20

	$p = 1000$											
	$R^2 = 0.20$				$R^2 = 0.50$				$R^2 = 0.80$			
	No outliers		5% outliers		No outliers		5% outliers		No outliers		5% outliers	
	C	R	C	R	C	R	C	R	C	R	C	R
%Correct	30	32	0	25	14.5	19.5	0	17	14.5	16	0	10
%Extra	53	27	0	26	75.5	50.5	0	43.5	77	54	0	43
%Missing 1	5.5	17	0	20	1.5	7.5	0	7	3	7.5	0	6
%Missing 2	1	5.5	0	6	0	0	0.5	1	0	0	1.5	1.5
%Missing 3	0	0.5	0	1	0	0	2.5	0	0	0	13	0
%Other	10.5	18	100	22	8.5	22.5	97	31.5	5.5	22.5	85.5	39.5
%mFDR	7.0	4.4	4.6	6.0	14.8	11.0	4.8	11.1	15.5	13.6	6.2	13.5
Time	5.8	10.8	6.1	11.7	5.7	10.6	5.8	11.7	5.86	10.9	5.47	11.3

model overfitting. We assess the latter through out-of-sample performance. Figure 1 shows out-of-sample MSE for the case $p = 100$. Robust VIF is as efficient as classical VIF when there are no outliers (top panel) and clearly outperforms classical VIF when there are 5% outliers (bottom panel). Much of the same can be seen in Figure 2 where results for the case $p = 1000$ are shown.

4 Real Data

In this section, we analyze two real data sets. In each case, the variables have been centered. The analysis of these two data sets will show how classical, i.e. non robust, VIF regression can be inadequate for the policy maker : in the first example, failing to keep important features, and in the second, failing to give a usable result. The selected models are compared using the median absolute prediction error (MAPE), as measured by 10-fold CV. That is, we split the data into 10 roughly equal-sized parts. For the k^{th} part, we carry out model selection using the other nine parts of the data and calculate the MAPE of the chosen model when predicting the k^{th} part of the data. We do this for $k = 1, \dots, 10$ and show boxplots of the 10 estimates of the MAPE. For all methods, the data were split in the same way. For the first data set, we randomly generated the folds, whereas in the second data set, we used the provided fold variable. Note here that we look at MAPE instead of mean squared prediction error as these real data can contain outliers (as opposed to the simulated data which were clean) and the MAPE is a better choice.

For completeness, we compare the models selected by classical and robust VIF approaches with that of the popular least angle regression (LARS) of Efron, Hastie, Johnstone, and Tibshirani (2004), an extremely efficient algorithm for computing the entire Lasso (Tibshirani 1996) path. We use the R package `lars` to do the computations.

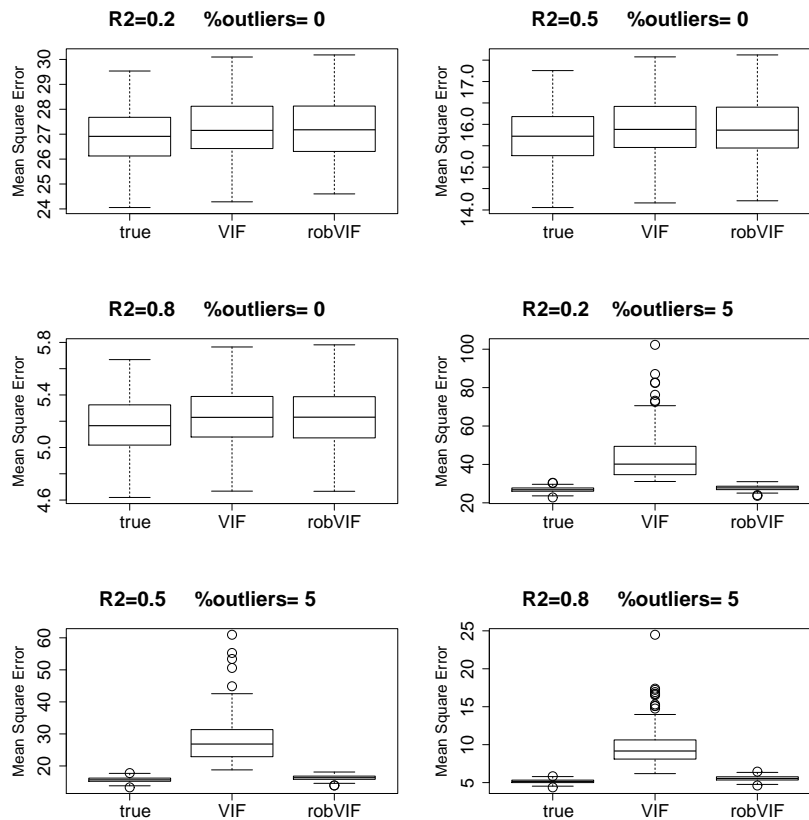


Figure 1: Out-of-sample mean square errors of the models chosen by classical and robust VIF regression. Simulated data with $n = 1000$ observations with $p = 100$ potential regressors, including $k = 5$ target regressors. Simulation scenarios as described in Table 1 where other details of the selected models can be found.

4.1 College Data

The data are in the R package *AER* and are a subset of the data previously analyzed in Rouse (1995). There are 4739 observations on 14 variables. The variables are listed in Table 2. We seek to predict the number of years of education using 13 economic and demographic variables. There are continuous and binary variables along with one categorical variable with three categories which is converted to two dummy variables. When considering only first-order variables we thus have $n = 4739$ and $p = 14$, when we include second-order interaction terms p rises to 104 (some interaction terms are constant and are removed).

Tables 3 and 4 list the VIF and robust VIF regression selected features, along with estimated slopes, for the $p = 14$ and $p = 104$ scenarios, respectively. For both scenarios, the robust VIF regression approach selects slightly more, and/or slightly different, features. When considering only first-order terms, we see that the classical and robust estimates of commonly selected features are almost the same. This serves as a good form of validation for the relative importance of these features. However, the presence of outliers in the data has lead classical VIF regression to completely miss two important features which are identified by robust VIF regression : **unemp** and **wage**. Even LS estimates (not shown) of the robust VIF regression selected model find these two features important with t-values of 3.15 and -2.70, but the classical VIF regression selection procedure could not detect this importance for the reasons outlined in Section 1. VIF regression also misses the two important features in the $p = 104$ scenario, see Table 4. As both the county unemployment rate and the state hourly wage in manufacturing are directly impacted by economic policy, policy makers must be equipped with the best feature selection tools to have an effective strategy to reach sought after goals : in this case, increasing the level of education among its constituents. These tools, we argue, must include a

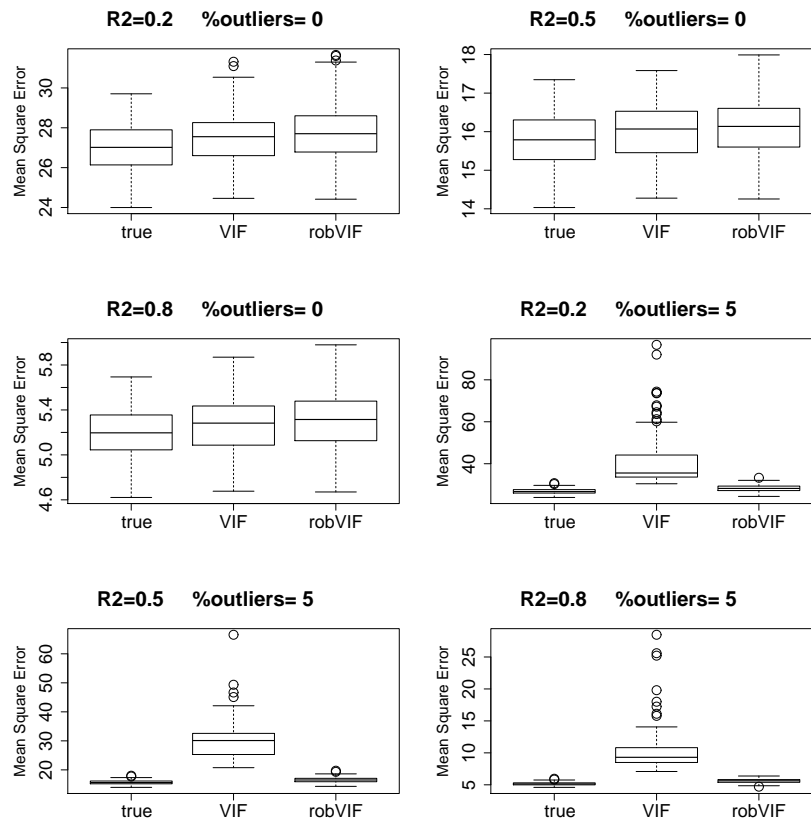


Figure 2: Out-of-sample mean square errors of the models chosen by classical and robust VIF regression. Simulated data with $n = 1000$ observations with $p = 1000$ potential regressors, including $k = 5$ target regressors. Simulation scenarios as described in Table 1 where other details of the selected models can be found.

Table 2: Original 14 Variables in College Data.

Variable	Description
<code>gender</code>	Factor indicating gender.
<code>ethnicity</code>	Factor indicating ethnicity (African-American, Hispanic or other).
<code>score</code>	Base year composite test score. These are achievement tests given to high school seniors in the sample.
<code>fcollege</code>	Factor. Is the father a college graduate?
<code>mcollege</code>	Factor. Is the mother a college graduate?
<code>home</code>	Factor. Does the family own their home?
<code>urban</code>	Factor. Is the school in an urban area?
<code>unemp</code>	County unemployment rate in 1980.
<code>wage</code>	State hourly wage in manufacturing in 1980.
<code>distance</code>	Distance from 4-year college (in 10 miles).
<code>tuition</code>	Average state 4-year college tuition (in 1000 USD).
<code>income</code>	Factor. Is the family income above USD 25,000 per year?
<code>region</code>	Factor indicating region (West or other).
<code>education</code>	Number of years of education.

robust selection procedure as shown effectively by this example. Further evidence is given in Figure 3 where MAPE for VIF, robust VIF and Lasso are shown for both scenarios. Robust VIF outperforms both of its competitors.

Table 3: VIF and robust VIF selected variables and estimated slope parameters (t-values) when only considering first-order terms. Significance : *0.05, **0.01, ***0.001.

Variable	$\hat{\beta}_{LS}$	$\hat{\beta}_{rob}$
ethnicityafam	0.349 (5.28)***	0.345 (4.90)***
ethnicityhispanic	0.360 (5.97)***	0.316 (4.92)***
score	0.088 (31.29)***	0.094 (31.77)***
fcollegeyes	0.540 (8.40)***	0.573 (8.51)***
mcollegeyes	0.380 (5.25)***	0.425 (5.60)***
homeyes	0.141 (2.39)*	0.148 (2.38)**
urbanyes	-	0.057 (0.96)
unemp	-	0.028 (3.00)**
wage	-	-0.047 (-2.56)**
distance	-0.027 (-2.81)**	-0.036 (-3.22)***
incomehigh	0.359 (6.70)***	0.398 (7.07)***

Table 4: VIF and robust VIF selected variables and estimated slope parameters (t-values) when including second-order interactions. Significance : *0.05, **0.01, ***0.001.

Variable	$\hat{\beta}_{LS}$	$\hat{\beta}_{rob}$
ethnicityafam	0.346 (5.25)***	0.342 (4.83)***
ethnicityhispanic	-0.028 (-0.20)	0.311 (4.83)***
score	0.088 (30.98)***	0.093 (27.26)***
fcollegeyes	0.543 (8.46)***	-0.080 (-0.16)
mcollegeyes	0.055 (0.28)	0.132 (0.24)
homeyes	0.136 (2.19)*	0.108 (1.60)
urbanyes	-	0.066 (1.12)
unemp	-	0.021 (2.09)*
wage	-	-0.050 (-2.35)**
distance	-0.036 (-3.53)***	-0.034 (-3.00)**
incomehigh	0.368 (6.87)***	0.089 (0.27)
genderfemale:score	0.001 (1.17)	-
genderfemale:fcollegeyes	-	0.008 (0.06)
genderfemale:mcollegeyes	0.407 (3.11)**	0.509 (3.43)***
ethnicityhispanic:unemp	0.048 (2.96)**	-
mcollegeyes:homeyes	0.120 (0.62)	-
score:incomehigh	-	0.006 (0.98)
fcollegeyes:homeyes	-	0.297 (1.74)*
fcollegeyes:unemp	-	0.028 (1.23)
fcollegeyes:wage	-	0.000 (0.00)
fcollegeyes:tuition	-	0.227 (1.44)
mcollegeyes:wage	-	0.000 (0.00)

4.2 Crime Data

We analyze recently made available crime data. These data are from the UCI Machine Learning Repository (Frank and Asuncion 2010) and are available at <http://archive.ics.uci.edu/ml/datasets/Communities+and+Crime>. We seek to predict the per capita violent crimes rate using economic, demographic, community and law enforcement related variables. After removing variables with missing data, we are left with $n = 1994$ observations on $p = 97$ first-order covariates. If we include second-order interactions (removing those that are constant) we have $p = 4753$. VIF regression selects 33 and 1993 variables, in the respective scenarios, while robust VIF regression selects 20 variables in both cases. Classical VIF experiences problems with the larger data set, which contains outliers in a highly multicollinear setting, and chooses too many covariates. This shows how the guarantee of no overfitting only holds at the model, i.e. without any outliers in the data. For these data, robust VIF regression provides the only viable option for policy makers as the 1993 features returned by classical VIF regression do not provide useful information. As can be seen in Figure 4, robust VIF is clearly the best performer for both scenarios. VIF regression chooses too many features for many of the folds and this leads to catastrophic results out-of-sample.

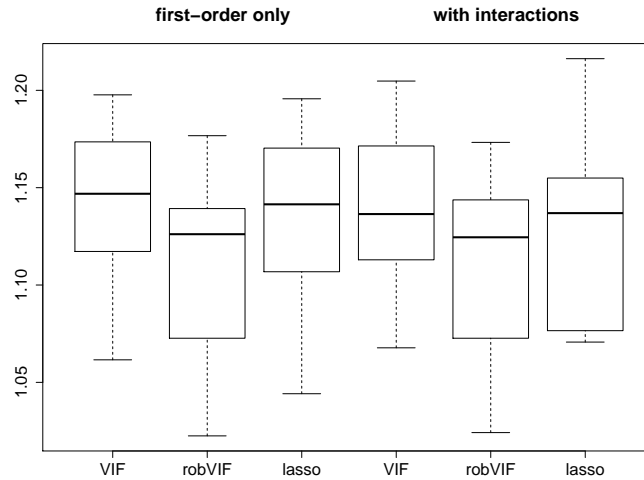


Figure 3: College data: Out-of-sample median absolute prediction errors of the models chosen by classical and robust VIF regression, and the Lasso, in 10-fold cross-validation.

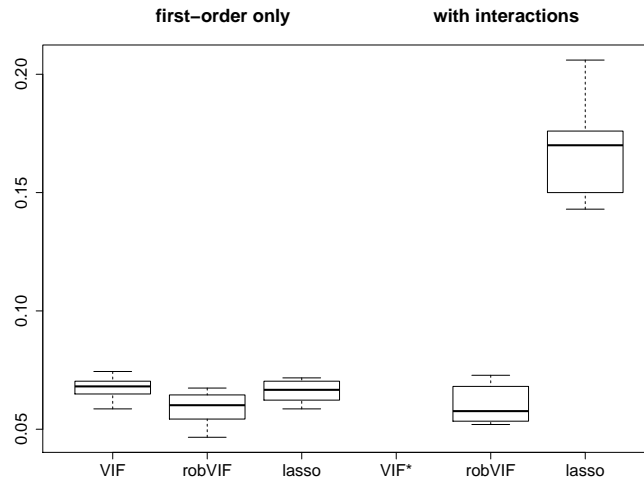


Figure 4: Crime and communities data: Out-of-sample median absolute prediction errors of the models chosen by classical and robust VIF regression, and the Lasso, in 10-fold cross-validation. * Results are not shown as VIF collapses in 4 folds, yielding MAPE of 5.62, 6.55, 6.82, 9.4, and 15.1, respectively. Results for other folds were good, 0.0652, 0.0676, 0.0686, 0.0694, 0.0744, but are excluded from the boxplot to allow for a better comparisons of all methods.

5 Concluding remarks

In Lin, Foster, and Ungar (2011), it was also shown that classical VIF regression equates or outperforms *stepwise regression*, *Lasso*, *FoBa*, an adaptive forward-backward greedy algorithm focusing on linear models (Zhang 2009), and *GPS*, the generalized path-seeking algorithm of Friedman (2008). In this paper, we present a very efficient robust VIF approach that clearly outperforms classical VIF in the case of contaminated data sets. This robust implementation comes with a very small cost in speed, computation time is less than doubled, and provides much needed robust model selection for large data sets.

Appendix - Algorithm Robust VIF regression

The robust VIF regression procedure, based on a streamwise regression approach and α -investing, can then be summarized by the following algorithm :

Input data $\mathbf{y}, \mathbf{x}_1, \mathbf{x}_2, \dots$ (standardized)

Set initial wealth $a_0 = 0.50$, pay-out $\Delta a = 0.05$, subsample size m , and robustness constant c .

Compute efficiency e_c^{-1} where

$$e_c = \left[\int_{-c}^c (5r^4/c^4 - 6r^2/c^2 + 1) d\Phi(r) \right]^2 / \int_{-c}^c r^2 (r^2/c^2 - 1)^4 d\Phi(r).$$

Get All marginal weights w_{ij} by fitting p marginal models $y = \beta_{01} + x_1\beta_1 + \varepsilon_1, \dots, y = \beta_{0p} + x_p\beta_p + \varepsilon_p$ using (5) and (6).

Initialize $j = 1, S = \{0\}, \mathbf{X}_S = \mathbf{1}, \mathbf{X}_S^w = \text{diag}(\sqrt{w_{iS}^0})\mathbf{X}_S$ and $\mathbf{y}^w = \text{diag}(\sqrt{w_{iS}^0})\mathbf{y}$ where w_{iS}^0 is computed using (4) where $\mathbf{r}^0 = (\mathbf{y} - \mathbf{1}\hat{\beta}^0)/\hat{\sigma}^0$ using $\mathbf{X}_0^w = \mathbf{X}_0^{w2} = \mathbf{1}, \hat{\beta}^0 = [(\mathbf{X}_0^w)^T \mathbf{X}_0^w]^{-1} (\mathbf{X}_0^{w2})^T \mathbf{y}$, where $\hat{\sigma}^0 = 1.483\text{med}|\hat{\mathbf{r}}^0 - \text{med}(\hat{\mathbf{r}}^0)|$ and $\hat{\mathbf{r}}^0 = \mathbf{y} - \mathbf{1}\hat{\beta}^0$.

repeat

set $\alpha_j = a_j/(1 + j - f)$

get T_w from the five-step *Fast Robust Evaluation Procedure* in §2.3.

if $2(1 - \Phi(|T_w|)) < \alpha_j$ **then**

$$S = S \cup \{j\}, \quad \mathbf{X}_S = [\mathbf{1} \quad \mathbf{x}_j], \quad \mathbf{X}_S^w = \text{diag}(\sqrt{w_{iS}^0})\mathbf{X}_S, \quad \mathbf{y}^w = \text{diag}(\sqrt{w_{iS}^0})\mathbf{y},$$

where w_{iS}^0 is computed using (4) where $\mathbf{r}^0 = (\mathbf{y} - \mathbf{X}_S\hat{\beta}^0)/\hat{\sigma}^0$ using $\mathbf{X}_0^w = [1 \quad \sqrt{w_{ij}x_{ij}}], \mathbf{X}_0^{w2} = [1 \quad w_{ij}x_{ij}], i = 1, \dots, n, \hat{\beta}^0 = [(\mathbf{X}_0^w)^T \mathbf{X}_0^w]^{-1} (\mathbf{X}_0^{w2})^T \mathbf{y}$, where $\hat{\sigma}^0 = 1.483\text{med}|\hat{\mathbf{r}}^0 - \text{med}(\hat{\mathbf{r}}^0)|$ and $\hat{\mathbf{r}}^0 = \mathbf{y} - \mathbf{X}_S\hat{\beta}^0$.

$$a_{j+1} = a_j + \Delta a$$

$$f = j$$

else $a_{j+1} = a_j - \alpha_j/(1 - \alpha_j)$

end if

$$j = j + 1$$

until all p covariates have been considered

References

- Akaike, H. (1973). Information theory and an extension of the maximum likelihood principle. In B. N. Petrov and F. Csaki (Eds.), *Proceedings of the Second International Symposium on Information Theory*, Budapest, pp. 267–281. Akademiai Kiado.
- Candes, E. J. and T. Tao (2007). The dantzig selector: Statistical estimation when p is much larger than n . *The Annals of Statistics* 35, 2313–2351.
- Dupuis, D. J. and M.-P. Victoria-Feser (2011). Fast robust model selection in large datasets. *Journal of the American Statistical Association*. 106, 203–212.
- Efron, B. (2004). The estimation of prediction error: Covariance penalties and cross-validation. *Journal of the American Statistical Association* 99, 619–632.
- Efron, B., T. Hastie, I. Johnstone, and R. Tibshirani (2004). Least angle regression (with discussion). *Annals of Statistics* 32, 407–499.
- Foster, D. and R. Stine (2004). Variable selection in data mining: Building a predictive model for bankruptcy. *Journal of the American Statistical Association* 99, 303–313.
- Frank, A. and Asuncion, A. (2010). UCI Machine Learning Repository [<http://archive.ics.uci.edu/ml>]. Irvine, CA: University of California, School of Information and Computer Science.
- Foster, D. P. and R. A. Stine (2008). Alpha-investing: A procedure for sequential control of expected false discoveries. *Journal of the Royal Statistical Society, Ser. B* 70, 429–444.
- Friedman, J.H. (2008). Fast sparse regression and classification. Technical report, Stanford University.
- Huber, P. J. (1964). Robust estimation of a location parameter. *Annals of Mathematical Statistics* 35, 73–101.
- Huber, P. J. (1967). The behavior of maximum likelihood estimates under non-standard conditions. In *Proceedings of the Fifth Berkeley Symposium on Mathematical Statistics and Probability*, Volume 1, pp. 221–233.
- Friedman, J. H., T. Hastie, and R. Tibshirani (2010). Regularization paths for generalized linear models via coordinate descent. *Journal of Statistical Software* 33, 1–22.
- Khan, J. A., S. Van Aelst, and R. H. Zamar (2007). Robust linear model selection based on least angle regression. *Journal of the American Statistical Association* 102, 1289–1299.
- Lin, D., D. P. Foster, and L. H. Ungar (2011). VIF regression: A fast regression algorithm for large data. *Journal of the American Statistical Association* 106, 232–247.
- Machado, J. A. F. (1993). Robust model selection and m -estimation. *Econometric Theory* 9, 478–493.
- Mallows, C. L. (1973). Some comments on c_p . *Technometrics* 15, 661–675.
- Marquardt, D. W. (1970). Generalized inverses, ridge regression, biased linear estimation, and nonlinear estimation. *Technometrics* 12, 591–612.
- Renaud, O. and M.-P. Victoria-Feser (2010). A robust coefficient of determination for regression. *Journal of Statistical Planning and Inference* 140, 1852–1862.
- Ronchetti, E. (1982). *Robust Testing in Linear Models: The Infinitesimal Approach*. Ph. D. thesis, ETH, Zurich, Switzerland.
- Ronchetti, E., C. Field, and W. Blanchard (1997). Robust linear model selection by cross-validation. *Journal of the American Statistical Association* 92, 1017–1023.
- Ronchetti, E. and R. G. Staudte (1994). A robust version of Mallows’s C_p . *Journal of the American Statistical Association* 89, 550–559.
- Rouse, C.E. (1995). Democratization or Diversion? The Effect of Community Colleges on Educational Attainment. *Journal of Business and Economic Statistics* 12, 217–224.
- Schwarz, G. (1978). Estimating the dimension of a model. *The Annals of Statistics* 6, 461–464.
- Tibshirani, R. (1996). Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society, Series B* 58, 267–288.
- Zhang, T. (2009). Adaptive forward-backward greedy algorithm for sparse learning with linear models. *Advances in Neural Information Processing Systems* 21, 1921–1928.
- Zhou, J., D. P. Foster, R. A. Stine, and L. H. Ungar (2006). Streamwise feature selection. *Journal of Machine Learning Research* 7, 1861–1885.