



Graphical Interpretation of Variance Inflation Factors

Author(s): Robert A. Stine

Source: *The American Statistician*, Vol. 49, No. 1 (Feb., 1995), pp. 53-56

Published by: [American Statistical Association](#)

Stable URL: <http://www.jstor.org/stable/2684812>

Accessed: 26/10/2011 13:37

Your use of the JSTOR archive indicates your acceptance of the Terms & Conditions of Use, available at <http://www.jstor.org/page/info/about/policies/terms.jsp>

JSTOR is a not-for-profit service that helps scholars, researchers, and students discover, use, and build upon a wide range of content in a trusted digital archive. We use information technology and tools to increase productivity and facilitate new forms of scholarship. For more information about JSTOR, please contact support@jstor.org.



American Statistical Association is collaborating with JSTOR to digitize, preserve and extend access to *The American Statistician*.

<http://www.jstor.org>

- Feder, P. I. (1968), "On the Distribution of the Log Likelihood Ratio Statistic When the True Parameter is Close to the Boundary of the Hypothesis Regions," *Annals of Mathematical Statistics*, 39, 2044–2055.
- Fuller, W. A. (1987), *Measurement Error Models*, New York: John Wiley.
- Hahn, G. J., and Meeker, W. Q. (1991), *Statistical Intervals: A Guide for Practitioners*, New York: John Wiley.
- Kalbfleisch, J. D., and Prentice, R. L. (1980), *The Statistical Analysis of Failure Time Data*, New York: John Wiley.
- Kalbfleisch, J. D., and Sprott, D. A. (1970), "Applications of Likelihood Methods to Models Involving Large Numbers of Parameters" (with discussion), *Journal of the Royal Statistical Society, Ser. B*, 32, 175–208.
- Lawless, J. F. (1982), *Statistical Models and Methods for Life Time Data*, New York: John Wiley.
- Lehmann, E. L. (1983), *Theory of Point Estimation*, New York: John Wiley.
- (1986), *Testing Statistical Hypotheses* (2nd ed.), New York: John Wiley.
- Meeker, W. Q. (1987), "Limited Failure Population Life Tests: Application to Integrated Circuit Reliability," *Technometrics*, 29, 51–65.
- Nelson, W. (1990), *Accelerated Testing: Statistical Models, Test Plans, and Data Analyses*, New York: John Wiley.
- Ostrouchov, G., and Meeker, W. Q. (1988), "Accuracy of Approximate Confidence Bounds Computed from Interval Censored Weibull and Lognormal Data," *Journal of Statistical Computation and Simulation*, 29, 43–76.
- Rao, C. R. (1973), *Linear Statistical Inference and Its Applications*, New York: John Wiley.
- Searle, S. R., Casella, G., and McCulloch, C. E. (1992), *Variance Components*, New York: John Wiley.
- Seber, G. A. F., and Wild, C. J. (1989), *Nonlinear Regression*, New York: John Wiley.
- Severini, T. A. (1991), "On the Relationship Between Bayesian and Non-Bayesian Interval Estimates," *Journal of the Royal Statistical Society, Ser. B*, 53, 611–618.
- Sprott, D. A. (1973), "Normal Likelihoods and Their Relation to Large Sample Theory of Estimation," *Biometrika*, 60, 457–465.
- (1975), "Application of Maximum Likelihood Methods to Finite Samples," *Sankhyā*, 37, 259–270.
- (1980), "Maximum Likelihood in Small Samples: Estimation in the Presence of Nuisance Parameters," *Biometrika*, 67, 515–523.
- Stuart, A., and Ord, K. (1991), *Kendall's Advanced Theory of Statistics, Volume 2: Classical Inference and Relationship* (5th ed.), Oxford, U.K.: Oxford University Press.
- Vander Wiel, S. A., and Meeker, W. Q. (1990), "Accuracy of Approximate Confidence Bounds Using Censored Weibull Regression Data from Accelerated Life Tests," *IEEE Transactions on Reliability*, 39, 346–351.
- Venzon, D. J., and Moolgavkar, S. H. (1988), "A Method of Computing Profile-Likelihood-Based Confidence Intervals," *Applied Statistics*, 37, 87–94.

Graphical Interpretation of Variance Inflation Factors

Robert A. STINE

A dynamic graphical display is proposed for uniting partial regression and partial residual plots. This animated display helps students understand multicollinearity and interpret the variance inflation factor. The variance inflation factor is presented as the square of the ratio of t -statistics associated with the partial regression and partial residual plots. Examples using two small data sets illustrate this approach.

KEY WORDS: Collinearity; Interactive plots; Regression diagnostics

1. INTRODUCTION

This article focuses on the connection between the variance inflation factor (VIF) and two diagnostic plots for least squares regression, partial regression plots, and partial residual plots (added-variable plots and component-plus-residual plots). To help students master regression diagnostics, I have found it useful to point out explicitly the connections among them. Introductions to regression diagnostics at the level of Chatterjee and Price (1991) or Fox (1991) offer the student a variety of numerical and graphical diagnostics for judging the adequacy of a regression model. There are diagnostics for specification error, outliers, multicollinearity, nonlinearity, heteroscedasticity, and other faults. Rather than present each diagnostic individually, I find it useful to describe the connections

among them, much as one needs to do in presenting the various types of random variables in an introductory course.

The presentation offered here is relatively elementary. The level is appropriate for students who do not know linear algebra, and I have found it useful in more advanced courses as well. The presentation relies upon imbedding the three diagnostics in a single dynamic plot. At one extreme of a slider control, this plot is the partial residual plot, which shows none of the effects of collinearity. As the control moves to the other extreme, it becomes the partial regression plot, which conveys the effects of multicollinearity. The plot dynamically updates its coordinates to suggest the effects of intermediate levels of multicollinearity.

2. THE DIAGNOSTICS

The VIF measures how much multicollinearity has increased the variance of a slope estimate. Suppose that we write the full-rank regression model for n independent observations as

$$y_i = \beta_0 + \beta_1 x_{i1} + \cdots + \beta_k x_{ik} + \epsilon_i, \quad i = 1, \dots, n,$$

where $\text{var}(\epsilon_i) = \sigma^2$. In vector form, the model is $Y = X\beta + \epsilon$ where X is the $n \times (k + 1)$ matrix with columns X_0, X_1, \dots, X_k and X_0 is a column vector of 1s. The name of this diagnostic arises from writing the variance of the least squares estimator $\hat{\beta}_j$ ($j = 1, \dots, k$) as (e.g., Belsley 1991, sec. 2.3)

$$\begin{aligned} \text{var}(\hat{\beta}_j) &= \sigma^2 (X'X)_{jj}^{-1} \\ &= \frac{\sigma^2}{SS_j} \text{VIF}_j, \end{aligned}$$

Robert A. Stine is Associate Professor, Department of Statistics, The Wharton School of the University of Pennsylvania, Philadelphia, PA 19104-6302.

where $SS_j = \sum_i (x_{ij} - \bar{x}_j)^2$ and

$$VIF_j = \frac{1}{1 - R_j^2} \quad (1)$$

R_j^2 is the R^2 statistic from the regression of X_j on the other covariates. Unfortunately, there is no well-defined critical value for what is needed to have a “large” VIF. Some authors, such as Chatterjee and Price (1991), suggest 10 as being large enough to indicate a problem.

The variance inflation factor is closely tied to the difference between two added variable plots for a regression. The partial regression plot for the j th variable shows two sets of residuals, those from regressing Y and X_j on the other covariates. The associated simple regression has slope $\hat{\beta}_j$ and the same residuals $\hat{\epsilon} = Y - X\hat{\beta}$ as the multiple regression. Indeed, with an adjustment for degrees of freedom, the variance of the slope estimate based on the partial regression plot is the same as that for $\hat{\beta}_j$ in the multiple regression,

$$\begin{aligned} \text{var}_j^{\text{regr}} &= \frac{n - k - 1}{n - 2} \hat{\sigma}^2 (X'X)_{jj}^{-1} \\ &= \frac{n - k - 1}{n - 2} \frac{\hat{\sigma}^2}{SS_j(1 - R_j^2)}, \end{aligned} \quad (2)$$

where $\hat{\sigma}^2 = \sum_i \hat{\epsilon}_i^2 / (n - k - 1)$. While seldom useful for detecting nonlinearity (neither axis shows an observed variable), these plots identify influential observations, reveal multiple outliers (masking), and show the effects of multicollinearity.

In contrast, partial residual plots offer a means for identifying nonlinearity. The partial residual plot corresponding to X_j shows $\hat{\epsilon} + \hat{\beta}_j X_j$ versus X_j . These plots ignore the effects of multicollinearity and convey a misleading impression of the significance of the fit, as noted by various authors including Atkinson (1985), Chatterjee and Hadi (1988), and Cook and Weisberg (1982). Although the associated simple regression again has slope $\hat{\beta}_j$ and residuals $\hat{\epsilon}$, the estimated variance of the fitted slope is

$$\text{var}_j^{\text{res}} = \frac{n - k - 1}{n - 2} \frac{\hat{\sigma}^2}{SS_j}. \quad (3)$$

The variance equations (2) and (3) are well known (e.g., Cook and Weisberg 1982, eq. 2.3.12 and 2.3.13). Noting the form of the VIF in (1), it is immediate (although not explicitly in this reference or elsewhere in the literature) that

$$VIF_j = \frac{\text{var}_j^{\text{regr}}}{\text{var}_j^{\text{res}}}. \quad (4)$$

In other words, VIF_j is the square of the ratio of the t -statistics from fits in the partial residual plot and partial regression plot.

3. THE DYNAMIC PLOT

A single dynamic plot ties these diagnostics together. Let $P_{(-j)}$ denote the projection matrix associated with all of the covariates but X_j . Following Cook and Weisberg (1982) or Chatterjee and Hadi (1988), define

$$\hat{\epsilon}_j(\lambda) = (I - \lambda P_{(-j)})(X_j - \bar{X}_j)$$

and

$$\begin{aligned} \hat{\epsilon}_Y(\lambda) &= (I - \lambda P_{(-j)})(Y - \bar{Y}) \\ &= \hat{\epsilon} + (I - \lambda P_{(-j)})\hat{\beta}_j(X_j - \bar{X}_j) \\ &= \hat{\epsilon} + \hat{\beta}_j \hat{\epsilon}_j(\lambda). \end{aligned}$$

The dynamic plot of $\hat{\epsilon}_Y(\lambda)$ on $\hat{\epsilon}_j(\lambda)$ allows the viewer to manipulate $0 \leq \lambda \leq 1$ using slider tools like those in Lisp-Stat (Tierney 1990).

The animation opens with $\lambda = 0$, which offers the greatest variation in the x axis and is the (centered) partial residual plot. Intuitively, this is the relationship between Y and X_j were X_j uncorrelated with the other covariates. As λ varies from zero to one, the animation shows how the points move in response to the changing amounts of collinearity. As λ approaches one, the plot approaches the partial regression plot and shows the full impact of the multicollinearity present in the data. The simplicity of the calculations makes real-time animation possible on personal computers. The display also gives the effective variance inflation factor associated with the plotted data,

$$VIF_j(\lambda) = \frac{1}{1 - \lambda R_j^2}.$$

The following examples using two small data sets illustrate the use of this plot. Cook and Weisberg (1989) present other dynamic regression diagnostics. As in their examples, I give a sequence of several frames which attempt to represent the animated display.

An implementation of this dynamic graphic is available from the author via e-mail. It requires that the user have Lisp-Stat. The code consists of several methods that enhance the standard regression model object in this package.

4. TWO EXAMPLES

The first example considers a time-series regression that has substantial collinearity. The regression considers the dependence of domestic U.S. crude oil production (OUTPUT) upon gross national product (GNP), price, a time trend (YEAR), and level of wildcat drilling activity during the 31 years 1948–1978. The data appear in problem 7.17 of Gujarati (1988). The OLS fit including the VIF's for this model appear in Table 1. The VIF for GNP is 62.1—clearly a “large” VIF.

The sequence of six frames shown in Figure 1 conveys a sense of how the plot changes as the value of λ ranges over the interval $[0, 1]$. The year 1973, the year of the first oil embargo, is an influential outlier and is highlighted throughout. Initially, with $\lambda = 0$, the fit looks quite good

Table 1. Summary of the Least Squares Regression Model Fitted to the Oil Production Data for 31 Observations 1948–1978. The Square of the Multiple Correlation is $R^2 = .94$ and $\hat{\sigma} = .29$.

Variable	Estimate	Standard Error	t-Statistic	VIF
Constant	2.62	1.52	1.7	n.a.
GNP	.0011	.0015	.7	62.1
YEAR	.0960	.044	2.2	58.1
PRICE	-.699	.070	-9.9	1.2
WILDCATS	.094	.029	3.2	1.7

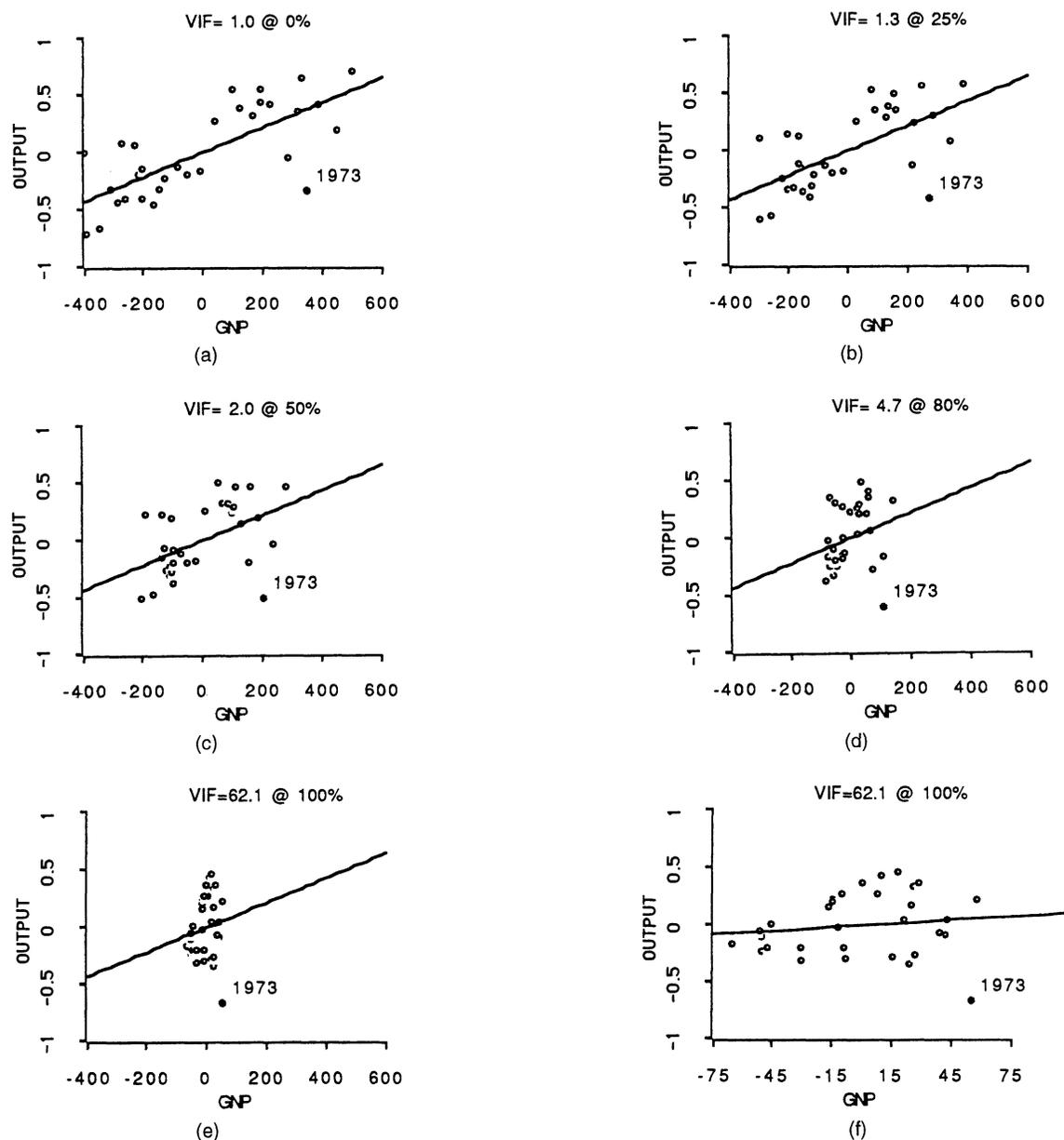


Figure 1. Frames From the Dynamic Plot for GNP in the Model for Oil Production.

in Figure 1a, the partial residual plot. As λ increases, collinearity compacts the variation on the x axis and the fit grows weaker. With $\lambda = .5$ in Figure 1c, the points are halfway to the partial regression plot, but $VIF(.5) = 2$ remains small. As λ nears one, the apparent VIF rapidly grows, reflecting the nonlinear definition of $VIF(\lambda)$. The sequence of plots shows that most of the “damage” is done by the time $VIF(\lambda)$ reaches the range 5–10, supporting the intuitive cutoff but allowing students to form their own opinions. Because the points move parallel to the fitted line, the plot also reinforces the notion that multicollinearity does not directly affect the residuals. The outlier year 1973 is just as far from the fitted line in Figure 1a as in Figure 1e. Figure 1f repeats Figure 1e, but with the x axis expanded to reveal the structure of the partial regression plot.

The second example demonstrates how collinearity affects a model with much less correlation among covariates. This example uses the data for 22 jet fighters reported in

Cook and Weisberg (1982, p. 47). The dependent variable is the log of the number of months after January 1940 of the first flight of the particular model of aircraft. The covariates are SPR (power per unit weight), RGF (range), PLF (payload as fraction of total weight), SLF (sustained load factor), and CAR (a dummy that is 1 if the plane can

Table 2. Summary of the Least Squares Regression for the Jet Fighter Data. The Response is $\text{LOG}(\text{FFD})$, the Log of the First Flight Date in Months after January 1940. $R^2 = .83$ and $\hat{\sigma} = .16$.

Variable	Estimate	Standard Error	t-Statistic	VIF
Constant	3.72	.27	14	n.a.
SPR	.085	0.022	3.9	1.45
RGF	.22	0.062	3.6	1.32
PLF	-.48	0.47	-1.0	1.15
SLF	.084	0.046	1.8	1.31
CAR	-.23	0.088	-2.7	1.27

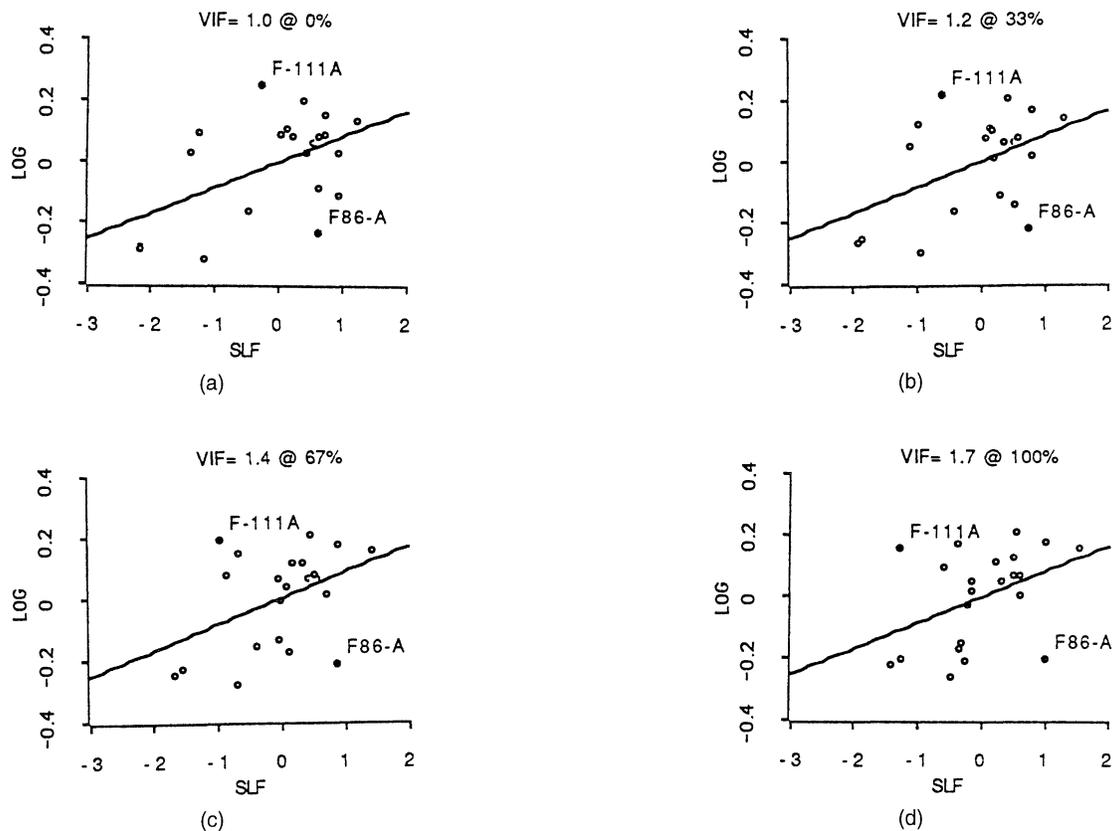


Figure 2. Frames From the Dynamic Plot for SLF in the Model for the log of the First Flight Date of Jet Fighters.

land on an aircraft carrier). Table 2 summarizes the fitted model. In contrast to the first example, all of the VIF's are less than 1.5.

While less dramatic than the first example, the dynamic plot for SLF is still interesting. In their analysis of these data, Cook and Weisberg note that two aircraft are outliers whose effects are disguised in the partial residual plot, but evident in the partial regression plot. The sequence of four frames of the dynamic VIF plot in Figure 2 shows how collinearity moves these two from being relatively innocent in Figure 2a to being quite influential (attenuating the slope) in Figure 2d.

[February 1993. January 1994.]

REFERENCES

- Atkinson, A. C. (1985), *Plots, Transformations, and Regression*, Oxford, U.K.: Oxford Publications.
- Belsley, D. A. (1991), *Conditioning Diagnostics*, New York: John Wiley.
- Chatterjee, S., and Price, B. (1991), *Regression Diagnostics*, New York: John Wiley.
- Chatterjee, S., and Hadi, A. S. (1988), *Sensitivity Analysis in Linear Regression*, New York: John Wiley.
- Cook, R. D., and Weisberg, S. (1982), *Residuals and Influence in Regression*, New York: Chapman and Hall.
- (1989), "Regression Diagnostics with Dynamic Graphics" (with discussion), *Technometrics*, 31, 277–311.
- Fox, J. (1991), *Regression Diagnostics*, Newbury Park, CA: Sage.
- Gujarati, D. N. (1988), *Basic Econometrics*, New York: McGraw-Hill.
- Tierney, L. (1990), *Lisp-Stat*, New York: John Wiley.

The Role of Geometry in Pairwise and Mutual Independence

Neil C. SCHWERTMAN and Terry L. KISER

The concepts of pairwise and mutual independence, while fundamental to probability theory, are sometimes difficult for students to differentiate. For pairwise independent random variables uniformly distributed over their region of support, a "rectangular support"—a Cartesian product of intervals—for the joint density is a necessary and sufficient condition for mutual independence. The use of geometric

illustrations will help students visualize this result and provide new insight into the difference between these two concepts.

KEY WORDS: Cartesian product space; Joint density; Marginal density; Support.

1. INTRODUCTION

The concepts of independence and marginal probability are fundamental to understanding probability theory

Neil C. Schwertman is Professor of Statistics and Terry L. Kiser is Associate Professor of Mathematics, Department of Mathematics and Statistics, California State University, Chico, CA 95929-0525.