# Regression Graphics: Added-Variable and Component+Residual Plots As Implemented in the car and effects packages for R

last modified: 2022-11-03

John Fox

McMaster University

CANSSI Statistical Software Conference 2022

# Outline

1. Introduction

2. Added-Variable Plots

3. Component+Residual Plots

4. References

# Outline

# Introduction

- I'll describe and illustrate two kinds of regression graphs and some of their extensions: added-variable (AV) plots and component+residual (C+R) plots.
- These graphs are implemented in the `avPlots()`, `crPlots()`, `avPlot3d()`, and `crPlot3d()` functions and their relatives in the **car** package for R, and in the `predictorEffects()` and `Effect()` functions and their relatives in the **effects** package for R.
- I'll focus on linear least-squares estimation of the regression model $y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \cdots + \beta_k x_{ik}$, $i = 1, \ldots, n$, where $\varepsilon_i \sim \mathrm{NID}(0, \sigma^2)$, but AV and C+R plots extend to other classes of regression models, such as generalized linear models estimated by maximum-likelihood; see, e.g., Wang (1987) for AV plots, and Cook and Croos-Dabrera (1998) for C+R plots.
- Most of the work that I'll discuss was undertaken jointly with Sandy Weisberg.
- General references for the methods in this presentation are Fox (2020), Fox and Weisberg (2019), Fox and Weisberg (2018), and Cook (1998).

# Outline

---

# Added-Variable Plots
Definition

- A very useful influence graph is the *added-variable plot*, also called a *partial-regression plot*.
  - Let $y_i^{(1)}$ represent the residuals from the least-squares regression of $y$ on all of the $x$s with the exception $x_1$, that is, the residuals from the fitted model

$$y_i = b_0^{(1)} + b_2^{(1)} x_{i2} + \cdots + b_k^{(1)} x_{ik} + y_i^{(1)}$$

  - Likewise, $x_i^{(1)}$ are residuals from the least-squares regression of $x_1$ on the other $x$s:

$$x_{1i} = c_0^{(1)} + c_2^{(1)} x_{i2} + \cdots + c_k^{(1)} x_{ik} + x_i^{(1)}$$

  - The notation emphasizes the interpretation of the residuals $y^{(1)}$ and $x^{(1)}$ as the parts of $y$ and $x_1$ that remain when the linear dependence of these variables on $x_2, \ldots, x_k$ is removed.
- The AV plot for $x_1$ is the scatterplot of $y^{(1)}$ versus $x^{(1)}$.
- We repeat the procedure for each $x_j$, $j = 0, 1, \ldots, k$ (where $x_0 = 1$).
- In effect the $(k + 1)$-dimensional scatterplot for $y$ and $x_1, \ldots, x_k$ is reduced to a sequence of $k + 1$ 2D AV plots.

# Added-Variable Plots
Properties

- The AV plots therefore visualize leverage and influence on each of the regression coefficients.
- The added-variable plot for $x_1$ has the following very interesting properties:
  - The slope of the least-squares *simple*-regression line of $y^{(1)}$ on $x^{(1)}$ is the same as the least-squares slope $b_1$ for $x_1$ in the full *multiple* regression.
  - The residuals from this simple regression are the same as the residuals $e_i$ from the full regression.
  - Consequently, the standard deviation of the residuals in the added-variable plot is $s$ from the multiple regression (if we use residual degrees of freedom $= n - k - 1$ to compute $s$).
  - The standard error of $b_1$ in the *multiple* regression is then $\text{SE}(b_1) = s/\sqrt{\sum x_i^{(1)^2}}$.
  - Because the $x_i^{(1)}$ are residuals, they are less variable than $x_1$ if $x_1$ is correlated with the other $x$s. The added-variable plot therefore shows how collinearity can degrade the precision of estimation by decreasing the conditional variation of an $x$.

# Added-Variable Plots
Extension to 3D

- AV plots extend straightforwardly to two regressors, say $x_1$ and $x_2$:
  1. Regress each of $x_1$, $x_2$, and $y$ on the other regressors, $x_3, \ldots, x_k$, obtaining residuals $x^{(1)}, x^{(2)}$, and $y^{(1,2)}$.
  2. Draw the 3D scatterplot of $y^{(1,2)}$ versus $x^{(1)}$ and $x^{(2)}$.
- The properties of the resulting 3D AV plot are entirely analogous to those of traditional 2D AV plots.

# Outline

# Component+Residual Plots
Lack-of-Fit: "Nonlinearity"

- The assumption that the average regression error is 0 everywhere implies that the regression surface captures the dependency of the conditional mean of $y$ on the $x$s.
- Violating the assumption of linearity implies that the model fails to represent the relationship between the average response and the explanatory variables.
  - For example, a partial relationship specified to be linear may be nonlinear, or two explanatory variables specified to have additive partial effects may interact in determining $y$.
  - Sometimes the fitted model may still be a useful approximation to the true regression surface $E(y)$.
  - But in other instance, the model may be extremely misleading.
- I think of nonlinearity (fitting the wrong equation to the data) as potentially the most serious problem with a regression model.
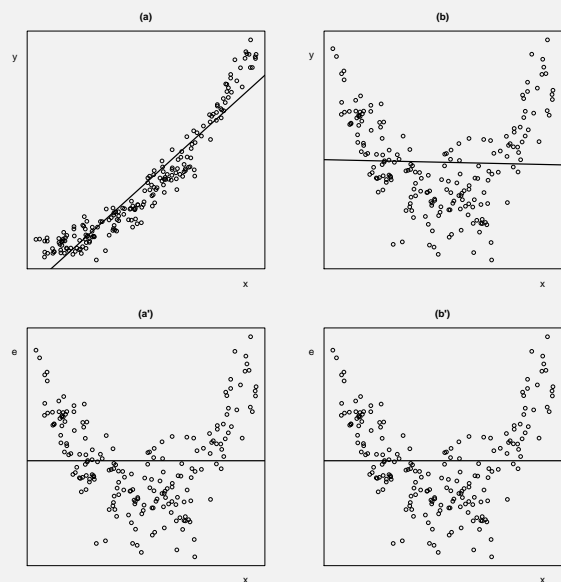
# Component+Residual Plots

- *Component+residual (C+R) plots* are the primary graphical device for diagnosing nonlinearity.
- The regression surface is generally high dimensional, even after accounting for regressors (such as polynomial terms, regression splines, dummy variables, and interactions) that are functions of a smaller number of explanatory variables.
- C+R plots and their relatives represent 2D views of the $(k+1)$D point-cloud of cases $\{y_i, x_{i1}, \ldots, x_{ik}\}$—similar to, but distinct from, AV plots.
- With modern computer graphics, these ideas here can be extended to 3D graphs (e.g., Cook, 1998).
- Even so, 2D and 3D projections of the data can fail to capture their systematic structure.

# Component+Residual Plots
## Limitations of Marginal Plots and Residual Plots

- It is useful to plot $y$ against each $x$ but these plots do not tell the whole story—and can be misleading.
  - Our interest centers on the *partial* relationship between $y$ and each $x$, controlling for the other $x$s, not on the *marginal* relationship between $y$ and a single $x$.
- Plotting residuals against each $x$ is helpful for detecting departures from linearity, but residual plots cannot distinguish between monotone and nonmonotone nonlinearity.
  - Case (a) might be modeled by $y = \beta_0 + \beta_1 x^2 + \varepsilon$ (a transformation of $x$), but (b) needs $y = \beta_0 + \beta_1 x + \beta_2 x^2 + \varepsilon$ (a quadratic).

# Component+Residual Plots
## Definition

- Added-variable plots, for detecting influential data, are partial plots, but they don't work well for detecting nonlinearity because they are biased towards linearity (Cook, 1998: Sec. 14.5).
- *Component+residual plots*, also called *partial-residual plots*, are often an effective alternative:
    1. Define the *partial residuals* for the $j$th regressor as
    $$e_i^{(j)} = b_j x_{ij} + e_i$$
    $e_i$ may include an unmodeled nonlinear component.
    2. Plot $e^{(j)}$ versus $x_j$.
    3. Repeat for $j = 1, \ldots, k$.
- By construction, the *multiple*-regression coefficient $b_j$ is the slope of the *simple* linear regression of $e^{(j)}$ on $x_j$, but nonlinearity may be apparent in the plot as well.
- This essentially simple idea was suggested independently by Larsen and McClearly (1972) and Wood (1973), and can be traced to work by Ezekial (1930).

# Component+Residual Plots
## Addressing Nonlinearity

- In multiple regression, we generally prefer to transform an $x$ rather than $y$, unless we see a common pattern of nonlinearity in the partial relationships of $y$ to several $x$s.
    - Transforming $y$ changes the shape of its relationship to *all* of the $x$s, and also changes the shape of the residual distribution.
- After an $x_j$ is transformed, we can plot partial residuals $e^{(j)}$ against the transformed $x_j$, say $t(x_j)$, in which case we want the plot to be linear, or against the original $x$, in which case we want the plot to follow the *partial fit* $\widehat{f}_j(x_j) = b_j t(x_j)$.
- This idea extends to partial fits based on more than one regressor, such as a quadratic, $b_{j1} x_j + b_{j2} x_j^2$: We can plot partial residuals $e_i^{(j)} = b_{j1} x_{ij} + b_{j2} x_{ij}^2 + e_i$ against the partial fit $b_{j1} x_{ij} + b_{j2} x_{ij}^2$, in which case the plot should be linear, or against $x_{ij}$, in which case the plot should follow the partial fit.

# Component+Residual Plots
## When are Component+Residual Plots Accurate?

- Cook (1993) explored the circumstances under which component+residual plots accurately visualize the unknown partial-regression function $f_1(x_1)$ in the model $y_i = \beta_0 + f_1(x_{i1}) + \beta_2 x_{i2} + \cdots + \beta_k x_{ik} + \varepsilon_i$ where $E(\varepsilon_i) = 0$.
  - The partial-regression function $f_1(x_1)$ may be nonlinear.
- Instead of fitting this "true" model, we fit the *working model* $y_i = \beta_0' + \beta_1' x_{i1} + \beta_2' x_{i2} + \cdots + \beta_k' x_{ik} + \varepsilon_i'$.
  - The partial residuals for the working model estimate $\varepsilon_i^{(1)} = \beta_1' x_{i1} + \varepsilon_i'$ rather than $f_1(x_{i1}) + \varepsilon_i$.
  - We hope that any nonlinear part of $f_1(x_1)$ is captured in the $\varepsilon_i'$.
  - Cook showed that $\varepsilon_i^{(1)} = f_1(x_{i1}) + \varepsilon_i$ either if the partial-regression function $f_1(x_1)$ is linear after all or if the other $x$s are each linearly related to $x_1$.
  - We can then legitimately smooth the scatterplot of the partial residuals versus $x_1$ to estimate $f_1(x_1)$.

# Component+Residual Plots
## When are Component+Residual Plots Accurate?

- There's therefore an advantage in having linearly related $x$s, a goal that's promoted, for example, by transforming the $x$s towards multivariate normality.
- In practice, it's only *strongly* nonlinearly related $x$s that seriously threaten the validity of C+R plots.
- A problem can also arise if $y$ is nonlinearly related to a *different $x$* (say, $x_2$) rather than to $x_1$:
  - Correlation between $x_1$ and $x_2$ can induce spurious nonlinearity in the C+R plot for $x_1$.
  - This suggests trying to correct nonlinearity for one $x$ at a time, but in my experience, it's rarely necessary to proceed sequentially.
- There are more robust versions of C+R plots that allow more complex relationships among the $x$s (see Mallows, 1986, and Cook, 1993) but these usually produce results very similar to simple C+R plots.
  - These methods are implemented in the `crPlots()` and `ceresPlots()` functions in the **car** package.

# Component+Residual Plots
Component+Residual Plots for Interactions

- Fox and Weisberg (2018) describe a framework for C+R plots that accommodates not only nonlinear terms, such as polynomials, regression splines, and transformations of $x$s, but also interactions of arbitrary complexity.
- This framework applies to *predictor effect plots*, which focus serially on each explanatory variable ("predictor") in a regression model, partitioning the other explanatory variables into two subsets:
  1. *conditioning predictors*, which interact with the focal predictor, either individually or in combination;
  2. *fixed predictors*, which simply are to be controlled statistically.

# Component+Residual Plots
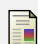Component+Residual Plots for Interactions

- The focal predictor ranges over its values in the data on the horizontal axis of a multi-panel array of 2D scatterplots of partial residuals versus the focal predictor for a combination of values of the conditioning predictors, while the fixed predictors are set to typical values.
- Conditioning is straightforward for factors, which simply take on in turn each of their various levels; numeric conditioning predictors are set successively to each of several representative values over their ranges.
- Conversely, fixing predictors is straightforward for numerical explanatory variables, which are typically set to their means; categorical fixed predictors are typically set to their distribution in the data.
- The regression surface is graphed by computing the fitted values under the model for the combinations of predictors in each panel.
- Each case in the data is allocated to one panel, and the residual for the case is added to its fitted value, forming a partial residual.

# Outline

# Regression Graphics: Added-Variable and Component+Residual Plots
## References

📄 R. D. Cook and R. Croos-Dabrera, *Partial residual plots in generalized linear models*, Journal of the American Statistical Association **93** (1998), 730–739.

📄 R. D. Cook, *Exploring partial residual plots*, Technometrics **35** (1993), 351–362.

📄 _____, *Regression graphics: Ideas for studying regressions through graphics*, Wiley, New York, 1998.

📄 M. Ezekial, *Methods of correlation analysis*, Wiley, New York, 1930.

📄 J. Fox, *Regression diagnostics: An introduction*, second ed., Sage, Thousand Oaks CA, 2020.

📄 J. Fox and S. Weisberg, *Visualizing fit and lack of fit in complex regression models with predictor effect plots and partial residuals*, Journal of Statistical Software **87** (2018), no. 9, 1–27.

# Regression Graphics: Added-Variable and Component+Residual Plots
## References

📄 _____, *An R companion to applied regression*, third ed., Sage, Thousand Oaks CA, 2019.

📄 W. A. Larsen and S. J. McClearly, *The use of partial residual plots in regression analysis*, Technometrics **14** (1972), 781–790.

📄 C. L. Mallows, *Augmented partial residuals*, Technometrics **28** (1986), 313–319.

📄 P. C. Wang, *Residual plots for detecting nonlinearity in generalized linear models*, Technometrics **29** (1987), 435–438.

📄 F. S. Wood, *The use of individual effects and residuals in fitting equations to data*, Technometrics **15** (1973), 677–695.