

INCREASED RESIDUAL VARIANCE AT DEVELOPMENTAL SWITCH POINTS: STATISTICAL ARTIFACT OR INDICATOR OF EXPOSED GENOTYPIC INFLUENCE?

MARCO DEL GIUDICE

Center for Cognitive Science, University and Polytechnic of Turin, Via Po 14, Turin, 10123, Italy
E-mail: delgiudice@psych.unito.it

Abstract.—Discussing the organization of developmental switches, West-Eberhard (2003) proposed the use of regression residual plots to locate the neutral point of the switch, which is characterized by maximum genotypic influence on the resulting phenotype. However, statistical artifacts due to measurement error in nonlinear models might account for a substantial proportion of the increase of residuals variance at the switch point. Simulations based on field data show that increases in residual variance occur as artifacts when normal amounts of measurement error are present, even in absence of any genotypic variance. The results suggest that interpretation of nonlinear variation in threshold traits is problematic and requires considering this statistical effect. A method to estimate the weight of genotypic contribution to residual variance is proposed, and its assumptions and limitations are discussed.

Key words.—Development, genotype, measurement error, phenotypic plasticity, residual plot.

Received July 16, 2005. Accepted October 29, 2005.

In her book on developmental plasticity, West-Eberhard (2003) elaborated the concept of a developmental switch point. As an example of hormonally mediated switch, she cited Emlen's work on horn size in dimorphic beetles (Emlen 1994, 1997). Some species of beetles (e.g., *Onthophagus acuminatus*, *O. taurus*) show a bimodal distribution of horn size, which is related to body size in a nonlinear fashion (Fig. 1a) and depends on the quantity of food ingested by the larva.

Fitting a sigmoidal regression model to the data allows for estimation of the neutral point (a body size of approximately 5.0 mm in Fig. 1), which is the value of the independent variable (in this case, body size) that acts on the dependent as a switch point between conditional development alternatives (short vs. long horn). By plotting the regression residuals (Fig. 1b), one can estimate the unexplained variance of the dependent at different locations on the switch. West-Eberhard (2003, p. 125) suggested that “one could use a plot of residuals to locate . . . the point of least environmental influence, and therefore the point where variation in genotypic influence on the threshold of a switch . . . is revealed.” She then proposed a seemingly useful way of interpreting the residuals plot: “Increased phenotypic variance near a switch point may be explained as follows: at extreme (high or low) values of the determining variable . . . there is a relatively clear signal to switch into one alternative pathway or the other At intermediate values, any genotypic differences in response threshold would be exposed.”

Although West-Eberhard's argument as a whole is quite compelling, there is a serious methodological shortcoming to this specific proposal. The assumption behind it is that one can expect residuals to distribute in an homogeneous fashion along the entire range of the distribution, and that departures from homoscedasticity near the flex point of the curve represent variance of possible genotypic origin. Residual plots could then be used to locate the point of minimum environmental influence (maximum genotypic influence). Unfortunately, this assumption holds only in the ideal case where the independent variable has no measurement error. When measurement error is present, no matter how small, it will

contribute to the residual variance; if the regression equation is nonlinear, measurement error will contribute in a nonlinear way.

Let ε_{mx} be the measurement error of the independent variable. In the case of a sigmoidal relationship, its contribution to the residuals on the dependent will be very small (approaching zero) for points of the curve far from the flex point and will increase with direct proportionality to the first derivative of the curve, reaching its maximum at the switch point (Fig. 2). The steeper the curve and the larger the error, the more ε_{mx} will lead to an increase of the residual variance near the switch point. This effect can be easily overlooked, since standard regression techniques assume that X is known with zero error. In standard biological data, measurement error is usually not trivial compared to the range of interest; this should often lead to dramatically increased variance around the switch point as a purely statistical artifact.

METHODS

To test the practical relevance of this concern, simulations were run based on actual biological data from a study of Moczek et al. (2002) on dung beetles (*O. taurus*). This dimorphic beetle shows a size distribution very similar to that reported by West-Eberhard (2003), and it shares the same size-dependent effect on horn development. The simulation was designed as follows.

Two hundred fifty sample points for the independent variable (body size in mm) were generated from a normally distributed population with mean = 5.03 mm, SD = 0.36 mm (values were taken from sample *NC1996*, Moczek et al. 2002; p. 593). Expected horn size was predicted using the same regression parameters obtained by the authors (p. 593). Then, a stochastic measurement error ε_m was added to both variables. ε_m was modeled as being normally distributed, with mean = 0 mm, SD = 0.05 mm. Although measurement error was not directly estimated in the cited study, this amount is a reasonably conservative estimate, since measures on the beetles were taken to the nearest 0.05 mm (Emlen 1994, in Moczek et al. 2002). Residuals for the dependent variable

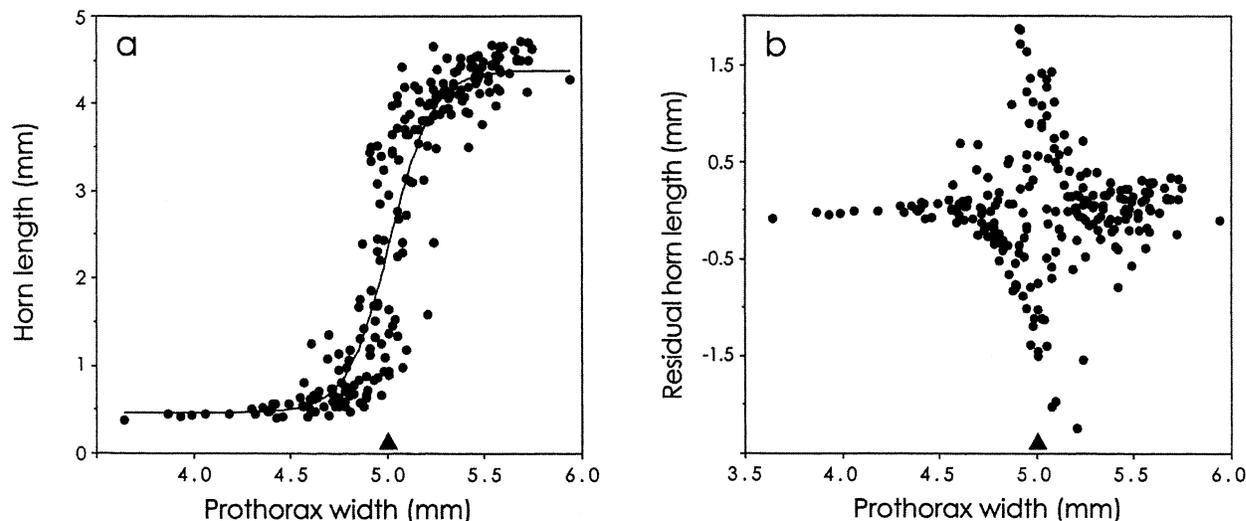


FIG. 1. (a) Horn length in relation to body size (indexed by prothorax width) in a sample of *Onthophagus acuminatus*. (b) Regression residuals plot. Residuals variance increases near the switch point (body size of about 5.0 mm). Reproduced with permission from Oxford University Press (West-Eberhard 2003).

were then calculated. The simulation was performed with Excel software (Microsoft Corp., Redmond, WA).

RESULTS

Figure 3 is the plot of a typical simulation result. As can be readily seen, it bears striking resemblance to the empirical one depicted in Figure 1. If ε_m is raised to a standard

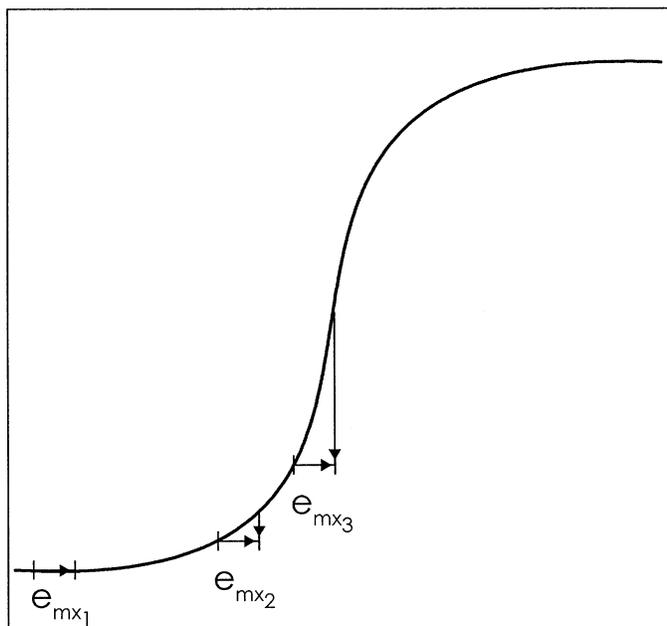


FIG. 2. Examples of the effect of independent variable measurement error on the regression residuals at different points of the curve. While ε_{mx} contribution to residuals is negligible in regions far from the flex (e_{mx1}), it increases steeply when approaching the switch point (e_{mx2} , e_{mx3}). In this region, even small deviations from the true value of the independent result in sizable departures from the expected value of the dependent.

deviation of about ± 0.08 mm, the resulting residual plots closely match the range and distribution of Figure 1b. It is important to note that the only stochastic element added to the regression model was a relatively small measurement error, with constant size along the whole measured range. Genotypic variability is not included in the model, yet the peaked pattern of residual variance near the switch point is reproduced as an artifact. Therefore, one cannot accept the idea that increased variance is entirely a sign of exposed genetic influence. When measurement error is not trivial, its effects on residuals could outweigh those due to genotypic variation.

DISCUSSION

West-Eberhard's (2003) interpretation of residual plots for the analysis of developmental switches is compromised by the fact that it does not take in account the artifacts generated by measurement error alone. However, the author discusses good reasons to expect a greater genotypic influence near the switch point, and Emlen (1996) effectively used regression residuals as a selection criterion in an artificial selection study. Useful as they might be, the above analysis argues strongly against an easy interpretation of residual plots. Even if they cannot be used to locate the neutral point, as West-Eberhard suggested, they could still prove useful to estimate the relative weight of genotypic contribution to residual variance at the switch point. To this aim, two problems need to be solved: (1) one must model the statistical artifact from measurement error; and (2) a way is needed to disentangle the unique contribution of genotypic variance from other unknown sources of residual error.

I will now outline a method to solve this problem in the case of a sigmoidal relationship between variables. The needed assumption is that additional, unknown sources of residual variance (e.g., environmental noise, stochastic processes) are uniformly distributed along the entire measured range. In

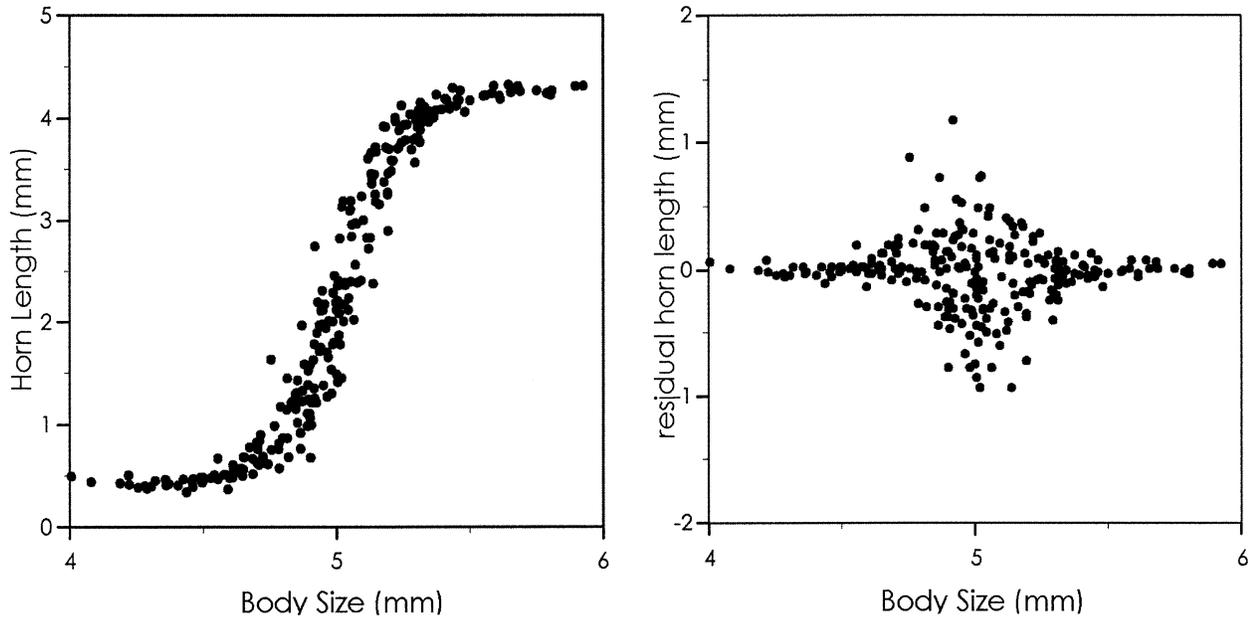


FIG. 3. Simulation of the body size–horn length relationship and residuals distribution in a sample of 250 beetles (*Onthophagus taurus*) in presence of measurement error, based on field data and equations from Moczek et al. (2002). Measurement error of both body size and horn length is set to ± 0.05 mm.

addition, both ε_{mx} and ε_{my} must be estimated with a high degree of precision; they need to be constant along the measured range, or at least their variation must be known (for ease of presentation, I will discuss the former case). If these assumptions are met, genotypic contribution to residuals can be estimated as follows.

First, we can decompose the total residual variance at a given point of the curve into separate components:

$$\sigma_{res}^2 = \sigma_{\varepsilon_{my}}^2 + f'(x)\sigma_{\varepsilon_{mx}}^2 + \sigma_{gen}^2 + \sigma_{\varepsilon}^2; \quad (1)$$

Total residual variance (σ_{res}^2) is the sum of measurement error variance of the dependent ($\sigma_{\varepsilon_{my}}^2$, constant along the range), measurement error variance of the independent ($\sigma_{\varepsilon_{mx}}^2$, constant along the range) multiplied by the function's derivative at a given point, $f'(x)$ (thus following a peaked distribution), variance from exposed genotypic variation (σ_{gen}^2 , expected to follow a similar peaked distribution for biological reasons), and additional variance of unknown source (σ_{ε}^2 , assumed to be constant).

At the range's extremes, the sigmoidal function can be approximated by a linear function with a slope of zero. This will result in the ε_{mx} contribution becoming zero. At the same time, the genotypic contribution is also expected to approach zero because of strong environmental canalization. In this section of the curve, the above expression will then simplify to

$$\sigma_{res_e}^2 = \sigma_{\varepsilon_{my}}^2 + \sigma_{\varepsilon}^2. \quad (2)$$

Because $S_{res_e}^2$ (sample residual variance at the extremes) and $S_{\varepsilon_{my}}^2$ (measurement error variance of the dependent) are known, we can estimate the additional component of variance σ_{ε}^2 by subtraction.

Now we turn to the central section of the curve. Here, the curve can be approximated by a linear function with slope

b_{sp} , equal to the derivative at the switch point. Total residual variance in this section ($\sigma_{res_{sp}}^2$) will then become

$$\sigma_{res_{sp}}^2 = \sigma_{\varepsilon_{my}}^2 + b_{sp}\sigma_{\varepsilon_{mx}}^2 + \sigma_{gen_{sp}}^2 + \sigma_{\varepsilon}^2. \quad (3)$$

Now, all components of variance are known with the exception of $\sigma_{gen_{sp}}^2$ (variance of genotypic origin at the switch point), which can then be estimated by subtraction. The statistical significance of the genotypic contribution can also be tested, comparing expected and empirical variance in the switch-point region under the null hypothesis that $\sigma_{gen_{sp}}^2 = 0$.

The main limitation of this method—and of any other method attempting to model the statistical artifact—is that it requires precise estimation of ε_{my} , ε_{mx} , and b_{sp} . A potential problem here is that the greater ε_{mx} , the more the residuals will be heteroskedastically distributed; this, in turn, will render precise determination of the equation parameters more difficult, affecting the estimate of b_{sp} . This method can be a viable approach if the study is carefully designed to minimize measurement error, and a great number of datapoints is collected (leading to high statistical power); residual plots, then, may offer a practical way to estimate the strength of exposed genotypic variance. When estimation of measurement error is difficult or too imprecise, however, more direct ways of assessing genotypic influence (such as breeding experiments) will probably be needed.

ACKNOWLEDGMENTS

I thank the Associate Editor, P. Phillips, and the anonymous reviewers for their valuable comments and suggestions.

LITERATURE CITED

Emlen, D. J. 1994. Environmental control of horn length dimorphism in the beetle *Onthophagus acuminatus* (Coleoptera: Scarabaeidae). *Proc. R. Soc. B* 256:131–136.

- . 1996. Artificial selection on horn length–body size allometry in the horned beetle *Onthophagus acuminatus* (Coleoptera: Scarabaeidae). *Evolution* 50:1219–1230.
- . 1997. Alternative reproductive tactics and male dimorphism in the horned beetle *Onthophagus acuminatus* (Coleoptera: Scarabaeidae). *Behav. Ecol. Soc.* 41:335–341.
- Moczek, A. P., J. Hunt, D. J. Emlen, and L. W. Simmons. 2002. Threshold evolution in exotic populations of a polyphenic beetle. *Evol. Ecol. Res.* 4:587–601.
- West-Eberhard, M. J. 2003. *Developmental plasticity and evolution*. Oxford Univ. Press, New York.

Corresponding Editor: P. Phillips