# CLONING DATA: GENERATING DATASETS WITH EXACTLY THE SAME MULTIPLE LINEAR REGRESSION FIT

## S. J. Haslett and K. Govindaraju*

*Massey University*

## Summary

This paper presents a simple computational procedure for generating 'matching' or 'cloning' datasets so that they have exactly the same fitted multiple linear regression equation. The method is simple to implement and provides an alternative to generating datasets under an assumed model. The advantage is that, unlike the case for the straight model-based alternative, parameter estimates from the original data and the generated data do not include any model error. This distinction suggests that 'same fit' procedures may provide a general and useful alternative to model-based procedures, and have a wide range of applications. For example, as well as being useful for teaching, cloned datasets can provide a model-free way of confidentializing data.

*Key words*: data cloning; multiple linear regression; orthogonality.

## 1. Introduction

Anscombe (1973) presented four diverse datasets to illustrate the importance of the graphical exploration of data. These four very well-known datasets have the same simple regression fit, $R^2$, and estimated residual standard deviation, $s$, etc. Chatterjee & Firat (2007) provided a method of generating diverse (bivariate) datasets with identical summary statistics but dissimilar graphs by using a genetic-algorithm-based approach. Govindaraju & Haslett (2007) presented a rather simpler procedure based on the orthogonality principle to obtain several naturally matching datasets giving the same simple regression equation. Regressing towards the mean is a natural concept that forms part of the material in Govindaraju & Haslett (2007), and hence the matching datasets there are not as arbitrary or as diverse as those of Anscombe (1973) and Chatterjee & Firat (2007). One possible application of such naturally matching datasets is for confidentializing sensitive real data for publication purposes, where having datasets with exactly the same fit as the original data is a major advantage, and major departures in data structure, such as are apparent in Anscombe's dataset and inherent in Chaterjee & Firat's, are disadvantageous.

## 2. Methodology

We first consider the case of independent and identically distributed data, but extend the methods to an arbitrary error covariance structure. Let the multiple regression model be

$$Y = X\beta + \varepsilon, \tag{1}$$

where $Y$ is the $n \times 1$ vector of responses, $X = (X_1 : X_2 : \ldots : X_p)$ is the $n \times p$ covariate data matrix, $\beta$ is the unknown $p \times 1$ vector of parameters, and $\varepsilon$ is the $n \times 1$ vector of errors. When $X$ has full (column) rank, the ordinary least squares (OLS) estimate of $\beta$ is

$$b = (X^\top X)^{-1} X^\top Y,$$

and the fitted multiple regression equation is

$$\hat{Y} = Xb. \tag{2}$$

Let $\mathbf{1}$ denote a vector of ones and $\bar{Y}$ and $\bar{X}_i$ be the sample means of the elements in $Y$ and $X_i$ for $i = 1, 2, \ldots, p$, respectively. Then, (2) can be written in a mean-corrected form as

$$\hat{y} = b_1 x_1 + b_2 x_2 + \cdots + b_p x_p, \tag{3}$$

where $\hat{y} = \hat{Y} - \bar{Y}$, with $\bar{Y} = \bar{y}\mathbf{1}$, and $x_i = X_i - \bar{X}_i$, with $\bar{X}_i = \bar{x}_i \mathbf{1}$ for $i = 1, 2, \ldots, p$ (and where, if $X$ contains $\mathbf{1}$ as one of its columns, then because that column becomes a column of zeros, it is dropped and $p$ in (3) is reset to one less than the number of columns in $X$, without any loss of generality).

Our problem is to obtain a new vector of the response variable, $Y_{\text{new}}$, and a new data matrix, $X_{\text{new}}$, such that

$$b = \left(X_{\text{new}}^\top X_{\text{new}}\right)^{-1} X_{\text{new}}^\top Y_{\text{new}}.$$

In other words, what we require is a matching multivariate dataset $(Y_{\text{new}}, X_{\text{new},1}, X_{\text{new},2}, \ldots, X_{\text{new},j}, \ldots, X_{\text{new},p})$ that will lead to exactly the same multiple linear regression equation as that for the original dataset $(Y, X_1, X_2, \ldots, X_j, \ldots, X_p)$. Later we consider the extension to the case where the $\varepsilon$ are correlated, with a specified (i.e. known or estimated) covariance matrix. This extension will include generating cloned data in linear mixed models.

Returning to the independent identically distributed case for now, we show below how the generation of matching or cloned data can be achieved by manipulating any one of the covariates, say $x_j$, using the following steps.

(1) First fit the multiple regression model as in (3), using mean corrected data.
(2) Choose any component $x_j$, where $j = 1, 2, \ldots, p$.
(3) Let $\quad \hat{y} = k + b_j x_j$, where $\quad k = b_1 x_1 + b_2 x_2 + \cdots + b_{j-1} x_{j-1} + b_{j+1} x_{j+1} + \cdots + b_p x_p = \hat{y} - b_j x_j$. Perform the simple regression of $y_k = y - k$ on $x_j$ and obtain the fitted values $\hat{y}_k$. In addition, perform the inverse simple regression of $x_j$ on $y_k$ and obtain the fitted values $\hat{x}_j$.
(4) Regress each $x_i (i \neq j)$ on $x_j$ and obtain $\hat{x}_i = x_j (x_j^\top x_j)^{-1} x_j^\top x_i$. Also obtain $x_{i,j} = x_i - x_j (x_j^\top x_j)^{-1} x_j^\top x_i = (I - x_j (x_j^\top x_j)^{-1} x_j^\top) x_i$, where $I$ is the identity matrix.
(5) Form $y_{k,\text{new}} = \hat{y}_k + \sum_{i \neq j} b_i x_{i,j}$.
(6) Carry out the multiple regression of $y_{k,\text{new}}$ simultaneously on all the newly obtained $\{x_{i,j} (i \neq j)\}$ and $\hat{x}_j$, where $\hat{x}_j = y(y^\top y)^{-1} y^\top x_j$ in which $y = (\mathbf{1} : \hat{y})$ is $n \times 2$.
(7) If preferred, add back $\bar{y}$ and $\bar{x}_i$ for $i = 1, 2, \ldots, p$ to the cloned data, and/or multiply all the cloned data by the same scale factor.

(8) Repeat steps 1 to 7 above to prepare a sequence of datasets, all with the exactly the same regression coefficients. The choice of a possibly different value of $j$ can be made at each iteration.

Step 2 ensures that $b_j$ is unchanged (see Govindaraju & Haslett 2007, for a proof). For the multivariate case, the regression coefficient for $\hat{x}_j$ must be $b_j$ because $\{x_{i,j}(i \neq j)\}$ are orthogonal to $\hat{x}_j$. Similarly, for each of the $x_{i,j}$, the corresponding regression coefficient must be $b_i$, again because of the orthogonal construction of $\{x_{i,j}(i \neq j)\}$.

When the errors $\boldsymbol{\varepsilon}$ in the multiple regression equation (1) are correlated, with the known or estimated full-rank covariance matrix $V$, say, an eigen-analysis or Choleski decomposition will provide a full-rank matrix $V^{1/2}$ with an inverse $V^{-1/2}$ that has the property $V^{-1/2}(V^{-1/2})^\top = V^{-1}$, the inverse of $V$. Pre-multiplication of $Y$ and $X$ in (1) by $V^{-1/2}$ will then yield an equation for which the ordinary least squares estimator (OLSE) used in the algorithm above is the best linear unbiased estimator (BLUE) of the model with error covariance $V$. When $V$ is specified in a mixed model, and the random parameters are incorporated into $V$, the estimates of the fixed parameters in the mixed model will remain unchanged, and, in addition, the random parameter estimates will remain unchanged by construction. The algorithm above consequently provides a method of generating a sequence of multivariate datasets all with the same fixed parameter estimates (i.e. BLUEs) even when the original data are correlated.

When the design matrix $X$ in (1) is not of full rank, as is often the case in experimental design, the following preliminary and consequent steps (which reduce $X$ to canonical form and later back-transform it) allow the same algorithm to be used. Suppose $X$ (which is $n \times p$) has rank $r$, where $r < p$. Find a basis for (the columns of) $X$, or more simply choose $r$ linearly independent columns of $X$. Arrange these as an $n \times r$ matrix, denoted as $X_r$, say. Define $\boldsymbol{\beta}_r$ via $X\boldsymbol{\beta} = X_r\boldsymbol{\beta}_r$. Then by construction, as $X$ and $X_r$ have the same column space, we have $X = X_r L_p$ where $L_p = (X_r^\top X_r)^{-1}X_r^\top X$, and $X_r = XL_r$ with $L_r = (X^\top X)^- X^\top X_r$. Here $L_r$ and $(X^\top X)^-$ are generalized inverses of $L_p$ and $(X^\top X)$, respectively, the choice of $(X^\top X)^-$ simply reflecting how $X_r$ has been specified. This specification allows the original transformation from $X$ to $X_r$ to be reversed. Hence transforming the problem from $y = X\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ to $y = X_r\boldsymbol{\beta}_r + \boldsymbol{\varepsilon}$ before applying the algorithm above and reversing the process afterwards gives the adjustment required if $X$ is not full rank. Explicitly $b = L_r b_r$, where $b_r$ is the solution to the transformed problem.

## 3. An example

Using uncorrelated data and a full-rank design matrix, consider the fictitious variables $X_1$, $X_2$ and $Y$ shown in Table 1 andresulting in the multiple regression fit

$$\hat{Y} = 1.48 + 0.967X_1 + 0.179X_2. \tag{4}$$

Steps 1 to 5 given above (manipulating $X_2$) will yield the matching or cloned data shown in Table 2, for which the fitted multiple regression equation is exactly the same as in (4).

The matching data in $X_{2,\text{new}}$ include negative values and average to zero, which may not be realistic in some contexts. This can be easily corrected by adding constant values to $X_{2,\text{new}}$ at the same time as adjusting for the cloned data average values of $Y$ and $Y_{\text{new}}$. Table 3 shows this adjustment. The data shown in Table 3 again yield the same regression fit as in (4).

TABLE 1

*A fictitious three-variable dataset*

| $X_1$ | $X_2$ | $Y$ |
|---|---|---|
| 3 | 4.1 | 5 |
| 4 | 6.2 | 7 |
| 4 | 4.5 | 6 |
| 3 | 3.1 | 5 |
| 2 | 3.9 | 4 |
| 5 | 6.1 | 8 |
| 5 | 7.2 | 7 |
| 7 | 8.6 | 10 |
| 6 | 7.5 | 8 |
| 6 | 8.7 | 9 |

TABLE 2

*Cloned data for the data in Table 1*

| $X_{1,\text{new}}$ | $X_{2,\text{new}}$ | $Y_{\text{new}}$ |
|---|---|---|
| 2.99351 | −0.12 | 4.35652 |
| 4.56021 | 0.80 | 5.48777 |
| 3.88909 | −0.90 | 5.18349 |
| 3.16617 | −1.12 | 4.17754 |
| 2.06341 | 0.86 | 3.56536 |
| 5.54211 | −0.48 | 6.22523 |
| 4.38755 | 0.62 | 6.42212 |
| 7.03974 | −0.34 | 8.18343 |
| 5.30039 | −0.26 | 7.23118 |
| 6.05784 | 0.94 | 7.44597 |

TABLE 3

*Matching data for Table 1 data are obtained from Table 2 cloned data by adding 1.07214 to each $Y_{\text{new}}$ element and 5.99 to the $X_{2,\text{new}}$ elements in order to preserve the same mean response for the raw and matching data*

| $X_{1,\text{new}}$ | $X_{2,\text{new}}$ | $Y_{\text{new}}$ |
|---|---|---|
| 2.99350 | 5.87 | 5.42866 |
| 4.56020 | 6.79 | 6.55991 |
| 3.88908 | 5.09 | 6.25563 |
| 3.16616 | 4.87 | 5.24968 |
| 2.06340 | 6.85 | 4.63750 |
| 5.54210 | 5.51 | 7.29737 |
| 4.38754 | 6.61 | 7.49426 |
| 7.03973 | 5.65 | 9.25557 |
| 5.30038 | 5.73 | 8.30332 |
| 6.05783 | 6.93 | 8.51811 |

Figure 1 shows the matrix plot of the raw and matching data given in Tables 1 and 3, respectively. The orthogonal manipulation done on $X_2$ in the algorithm above, applied once, gives $X_{2,\text{new}}$. Figure 1 shows the effect of this manipulation. The strength of the bivariate relationship between $X_{2,\text{new}}$ and $Y_{\text{new}}$ is much weaker than that of the relationship between $X_2$
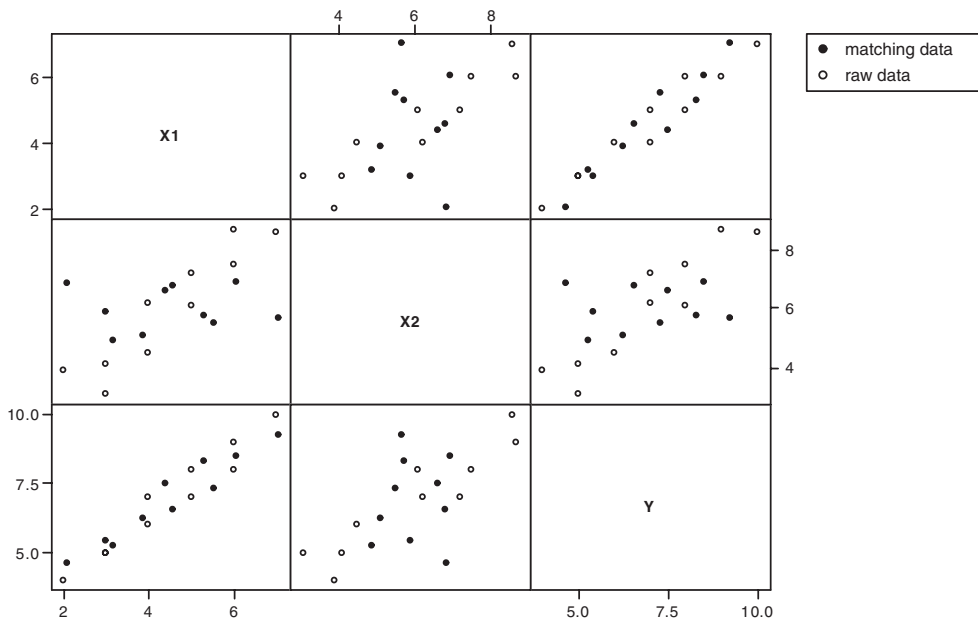
Figure 1. Matrix plot of raw and matching datasets, based on Table 1 raw data and Table 3 matching data. Data cloning was performed by manipulating $X_2$, and hence the matching $X_{2,\text{new}}$ data show a weaker relationship with the other variables.

and $Y$. However, this is not the case with $X_{1,\text{new}}$ and $Y_{\text{new}}$, because the orthogonal manipulation was not done with $X_1$.

## 4. Discussion

The usefulness of the algorithm in teaching to produce any number of cloned datasets with exactly the same multiple regression is clear. Perhaps less obviously, the algorithm provides an alternative way of confidentializing data so that a multiple regression analysis (which may include fitting mixed models, such as occur in small-area estimation, for example) has exactly the same fit in the cloned as in the original data, even though the data have been changed to be no longer confidential. This would seem to have important applications where model fitting to confidential unit records is required but the original data cannot be released. In this context, further research on whether methods can be developed so that cloning is possible for a collection of models would be a useful future research topic. More generally, given their parallels with multiple regression, data cloning for multivariate statistics may also warrant further research.

## References

ANSCOMBE, F.J. (1973). Graphs in statistical analysis. *Amer. Statist.* **27**, 17–21.

CHATTERJEE, S. & FIRAT, A. (2007). Generating data with identical statistics but dissimilar graphics: A follow up to the Anscombe dataset. *Amer. Statist.* **61**, 248–254.

GOVINDARAJU, K. & HASLETT, S.J. (2007). Illustration of regression towards the means. *Int. J. Math. Edu. Sci. Tech.* **39**, 544–550.