Correlation of Residuals in Successive Fittings with Least-Squares

 Robert Jacobsen

## Abstract

This paper derives the covariance relations of the residuals in successive least-squares fits, with application to tests of heteroscedasticity.

------------------------------

Biometrics Unit, Department of Plant Breeding and Biometry, Cornell University, Ithaca, New York.

Correlation of Residuals in Successive Fittings with Least-Squares

BU-242-M                      Robert Jacobsen                      June, 1967

We give some simplified proofs and extensions of results in A. Hedayat's paper No. BU-135.

Let V be the observation space of dim. N, $\mathcal{Y}$ the observed point.

Let $E_\theta \mathcal{Y} = X\,\theta$, $\theta \in \ominus = R^p$, X: $\ominus \to V$ linear.

Let $\Omega$ denote the mean space, Im X, and Cov $\mathcal{Y} = D$, where D is diagonal with respect to the orthonormal standard basis $e_1, \cdots, e_N$.

Denote $V_i$ = the span of $\{e_1, \cdots, e_i\}$, and $\Omega_i = P_{V_i}\Omega$, where $P_W$ denotes orthogonal projection onto $W \subset V$.

We are concerned with computing the covariance relations among the least-squares estimates of $E\mathcal{Y}$ and the residuals based on different numbers of observations.

(1)   Now cov $\left[ \ (e_k,\ P_{V_i - \Omega_i}\, P_{V_i}\, \mathcal{Y}\ ),\ (e_\ell,\ P_{V_j - \Omega_j}\, P_{V_j}\, \mathcal{Y}\ )\ \right]$

$$k = 1, \cdots, i;\quad \ell = 1, \cdots, j;\quad 1 \le i \le j \le N,$$

is the covariance between the $k^{th}$ coordinate of the residual vector, based on a fit to the $1^{st}$ i observations, and the $\ell^{th}$ coordinate of the residual based on the $1^{st}$ j observations.

---

Biometrics Unit, Department of Plant Breeding and Biometry, Cornell University, Ithaca, New York.

$(e_k, P_{V_i - \Omega_i} P_{V_i} \mathcal{Y}) = (P^1_{V_i - \Omega_i} P^1_{V_i} e_k, \mathcal{Y}) = (P_{V_i - \Omega_i} e_k, \mathcal{Y})$, as a projection

is self-adjoint.

$P_{V_i - \Omega_i} e_k = (I - P_{\Omega_i}) e_k$, as $e_k \in V_i$. So (1) becomes $(P_{V_i - \Omega_i} e_k, D P_{V_j - \Omega_j} e_\ell)$
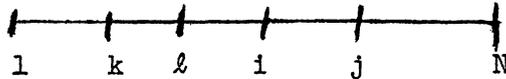
$= (e_k, D e_\ell) - (e_k, D P_{\Omega_j} e_\ell) - (D e_\ell, P_{\Omega_i} e_k) + (P_{\Omega_i} e_k, D P_{\Omega_j} e_\ell)$, \hfill (2)

by the definition of cov $\mathcal{Y}$.

### Evaluation of (2).

Assume $D = \sigma^2 I$.

Case 1.



$1 \leq k \leq i < \ell \leq j \leq N$

Write $e_k = P_{\Omega_i} e_k + P_{V_i - \Omega_i} e_k$

Now $V_i - \Omega_i \perp \Omega_i$ and $V_j - V_i$.

But $\Omega_j \subset \Omega_i \oplus V_j - V_i$.

So $V_i - \Omega_i \perp \Omega_j$.

(4) Hence, $(e_k, P_{\Omega_j} e_\ell) = (P_{\Omega_i} e_k, P_{\Omega_j} e_\ell) + (P_{V_i - \Omega_i} e_k, P_{\Omega_j} e_\ell)$

$\qquad\qquad = (P_{\Omega_i} e_k, P_{\Omega_j} e_\ell)$.

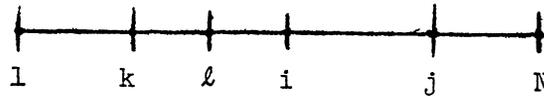(5) Further, as $i < \ell$, $(e_\ell, P_{\Omega_i} e_k) = 0$.

(6) And, as $k < \ell$, $(e_k, e_\ell) = 0$.

So (2) becomes 0.

Therefore, any component of the residual based on the first i observations is uncorrelated with any component > i of the residual based on the first j observations, $j > i$, in the homoscedastic case.
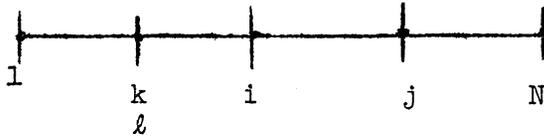
Case 2.

$1 \leq k < \ell \leq i \leq j \leq N$



(4) and (6) still hold.

So (2) $= -(e_\ell, P_{\Omega_i} e_k) \sigma^2$, which doesn't depend on j.

Case 3.

$1 \leq k = \ell \leq i \leq j \leq N$



(4) still holds.

So (2) $= \sigma^2 \left\{ (e_\ell, e_\ell) - (e_\ell, P_{\Omega_i} e_\ell) \right\}$, which doesn't depend on j.

(7) Thus, $\text{var} (e_\ell, P_{V_i - \Omega_i} P_{V_i} \mathcal{Y}) = (1 - \|P_{\Omega_i} e_\ell\|^2) \sigma^2$

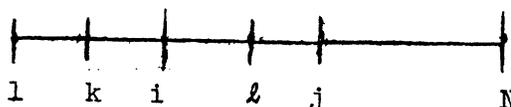A formula for the correlation between two residuals can be given.

(8) $\rho_{k,i,\ell,j} = \dfrac{-(e_\ell, P_{\Omega_i} e_k)}{\left[ 1 - (e_\ell, P_{\Omega_i} e_\ell) \right]^{\frac{1}{2}} \left[ 1 - (e_k, P_{\Omega_j} e_k) \right]^{\frac{1}{2}}}$

for $1 \leq k < \ell \leq i \leq j \leq N$.

## Evaluation of (2).

Assume $D = \begin{pmatrix} \sigma_1^2 & & 0 \\ & \ddots & \\ 0 & & \sigma_N^2 \end{pmatrix}$ .
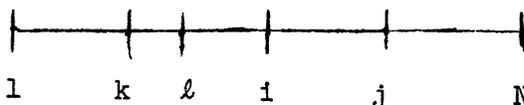
Case 1.

$$\begin{array}{cccccc} \vdash & \mid & \mid & \mid & \mid & \mid \\ 1 & k & i & \ell & j & N \end{array}$$

$1 \leq k \leq i < \ell \leq j \leq N.$

The $1^{st}$ and $3^{rd}$ terms of (2) vanish. (2) becomes

$$(P_{\Omega_i} e_k, D P_{\Omega_j} e_\ell) - (e_k, D P_{\Omega_j} e_\ell) = (D P_{\Omega_i} e_k, P_{\Omega_j} e_\ell) - \sigma_k^2(e_k, P_{\Omega_j} e_\ell) \qquad (9)$$
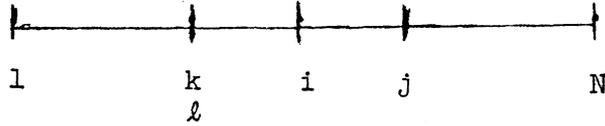
which is not, in general, zero.

Case 2.

$$\begin{array}{cccccc} \vdash & \mid & \mid & \mid & \mid & \mid \\ 1 & k & \ell & i & j & N \end{array}$$

$1 \leq k < \ell \leq i \leq j \leq N.$

The $1^{st}$ term of (2) vanishes. (2) becomes

$$- (e_k, D P_{\Omega_j} e_\ell) - (D e_\ell, P_{\Omega_i} e_k) + (P_{\Omega_i} e_k, D P_{\Omega_j} e_\ell), \text{ which depends on } j.$$
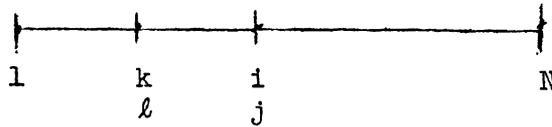
Case 3.

$$1 \leq k = \ell \leq i < j \leq N.$$

(2) remains unchanged.

Case 4.

$$1 \leq k = \ell \leq i = j \leq N.$$

The $2^{nd}$ and $3^{rd}$ terms of (2) become identical. (2) becomes

$$(10) \quad \sigma_\ell^2 - 2(e_\ell, D P_{\Omega_i} e_\ell) + (P_{\Omega_i} e_\ell, D P_{\Omega_i} e_\ell) = \text{var}(e_\ell, P_{V_i - \Omega_i} P_{V_i} y)$$

Now to investigate

$$(11) \quad \text{cov}\left[ (e_k, P_{\Omega_i} P_{V_i} y), (e_\ell, P_{V_j - \Omega_j} P_{V_j} y) \right]$$

the covariance between the $k^{th}$ coordinate of the estimated mean vector based on the first $i$ observations and the $\ell^{th}$ coordinate of the residual based on the first $j$ observations.

$$k = 1, \cdots, i; \quad \ell = 1, \cdots, j; \quad 1 \leq i \leq N; \quad 1 \leq j \leq N.$$
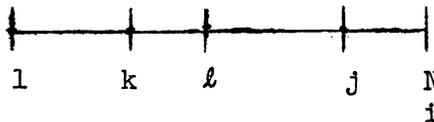
(11) becomes $(P_{\Omega_i} e_k, D P_{V_j - \Omega_j} e_\ell) = (P_{\Omega_i} e_k, D(I - P_{\Omega_j}) e_\ell)$

$$= (P_{\Omega_i} e_k, D e_\ell) - (P_{\Omega_i} e_k, D P_{\Omega_j} e_\ell) \tag{12}$$

## Evaluation of (12).

Assume $D = \sigma^2 I$.

Case 1.



$1 \le k,\ \ell,\ j \le N$

$\quad i = N, \ell \le j$

(12) becomes

$$\sigma^2 \left\{ (P_{\Omega_N} e_k,\ e_\ell) - (P_{\Omega_N} e_k,\ P_{\Omega_j} e_\ell) \right\} = \sigma^2 (P_{\Omega_N} e_k,\ P_{V_j - \Omega_j} e_\ell)$$

But $V_j - \Omega_j \perp \Omega_j$ and $\perp V_N - V_j$. And $\Omega_N \subset \Omega_j \oplus V_N - V_j$. So $V_j - \Omega_j \perp \Omega_N$. Hence above equals O.

## Evaluation of (12).

Assume $D = \begin{pmatrix} \sigma_1^2 & & 0 \\ & \ddots & \\ 0 & & \sigma_N^2 \end{pmatrix}$.

Case 1.

$$1 \le k, \ell, j \le N$$
$$i = N, \ell \le j.$$

(12) becomes

$$\sigma_\ell^2 (P_{\Omega_N} e_k, e_\ell) - (P_{\Omega_N} e_k, D P_{\Omega_j} e_\ell),$$

which is not zero, in general.

Let $f_n = (e_n, P_{V_n - \Omega_n} P_{V_n} y)$

Var $f_n = \left\{ 1 - \| P_{\Omega_n} e_n \|^2 \right\} \sigma^2 = C_n \sigma^2$, under homoscedasticity assumption.

Var $f_n = \sigma_n^2 - 2\sigma_n^2 (e_n, P_{\Omega_n} e_n) + (P_{\Omega_n} e_n, D P_{\Omega_n} e_n) = C_n' \sigma_n^2$, under

heteroscedasticity assumption.

Let $d_n = \dfrac{f_n}{\sqrt{c_n}}$ .

Var $d_n = \sigma^2$, under homoscedasticity assumption.

$$= \sigma_n^2 \frac{c_n'}{c_n} , \text{ under heteroscedasticity assumption.}$$

Under homoscedasticity assumption, the $d_n$ are uncorrelated, with constant var. $\sigma^2$, $n = r + 1, \cdots, N$, where $r$ = rank of X. Under heteroscedasticity assumption, the $d_n$ are correlated, with $\text{cov}(d_n, d_{n+1}) = \dfrac{1}{\sqrt{c_n c_{n+1}}} \text{cov}(f_n, f_{n+1})$,

and var $d_n = \sigma_n^2 \dfrac{C_n'}{C_n}$.

The d's have expectation 0, under both hypotheses. If heteroscedasticity holds,

$$P\left\{\, |d_{n+1}| > |d_n|\, \right\} > \frac{1}{2} \quad \text{if}$$

$$\left| \; \text{cor}\,(d_{n+1},\, d_n) \; \sqrt{\dfrac{\text{var } d_{n+1}}{\text{var } d_n}} \; \right| \; > 1.$$

Above equals $\left| \dfrac{\text{cov}\,(d_{n+1},\, d_n)}{\text{var } d_n} \right| = \left| \dfrac{\text{cov}\,(f_{n+1},\, f_n)}{\sigma_n^2 \dfrac{C_n'}{C_n} \sqrt{C_{n+1}\, C_n}} \right|$

$$= \left| \dfrac{\text{cov}\,(f_{n+1},\, f_n)}{\text{var } f_n \; \dfrac{\sqrt{C_{n+1}}}{C_n}} \right| \; .$$

Thus, a sufficient condition that $P\left\{|d_{n+1}| > |d_n|\right\} > \dfrac{1}{2}$ $\quad n = r + 1, \cdots, N$

is that the absolute value of

$$\dfrac{(D\, P_{\Omega_n} e_n,\; P_{\Omega_{n+1}} e_{n+1}) - \sigma_n^2\, (e_n,\; P_{\Omega_{n+1}} e_{n+1})}{\left[\, \sigma_n^2 - 2\sigma_n^2\, (e_n,\; P_{\Omega_n} e_n) + (P_{\Omega_n} e_n,\; D\, P_{\Omega_n} e_n)\,\right] \left[\, \dfrac{1 - \|P_{\Omega_{n+1}} e_{n+1}\|^2}{1 - \|P_{\Omega_n} e_n\|^2}\,\right]^{\frac{1}{2}}}$$

be $\geq 1$.

This condition could then be used to insure power against alternatives in the Goldfeld, Quandt peak-test.

## References

[1] Hedayat, Abdossamad (1966). Homoscedasticity in Linear Regression Analysis with Equally Spaced x's. M.S. Thesis, Cornell University, Ithaca, New York.