



OXFORD JOURNALS
OXFORD UNIVERSITY PRESS

A Historical Note on the Method of Least Squares

Author(s): R. L. Plackett

Source: *Biometrika*, Dec., 1949, Vol. 36, No. 3/4 (Dec., 1949), pp. 458-460

Published by: Oxford University Press on behalf of Biometrika Trust

Stable URL: <https://www.jstor.org/stable/2332682>

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.

Your use of the JSTOR archive indicates your acceptance of the Terms & Conditions of Use, available at <https://about.jstor.org/terms>



and Oxford University Press are collaborating with JSTOR to digitize, preserve and extend access to *Biometrika*

JSTOR

If there is more than one unknown parameter, and the variance v of the observations z_i is finite but not a known function of the mean θ , v must be estimated as sampling proceeds, and a boundary cannot be specified in advance. Let us suppose that v can be estimated from a sample of m observations by an estimate \hat{v} that is asymptotically normal with variance $O(m^{-1})$ as $m \rightarrow \infty$. When, say, half the sampling has been done that will ultimately be seen to be needed, \hat{v} satisfies in probability the relation

$$\hat{v} = v(1 + O(n^{-1})), \quad (23)$$

and so the eventual sample size n and estimate $\hat{\theta}$ of θ satisfy

$$\text{var}(\hat{\theta}) \sim \hat{v}/n. \quad (24)$$

Thus the required value of n can be predicted at this half-way stage with asymptotic validity.

I am indebted to Mr D. V. Lindley for helpful criticism. The treatment given above is somewhat heuristic in places, notably in the argument used in arriving at equation (11). It seems to me, however, that the line of approach may prove to be worth pursuing.

REFERENCES

- BLACKWELL, D. & GIRSHICK, M. A. (1947). A lower bound for the variance of some unbiased sequential estimates. *Ann. Math. Statist.* **18**, 277.
 FINNEY, D. J. (1949). On a method of estimating frequencies. *Biometrika*, **36**, 233.
 HALDANE, J. B. S. (1945). On a method of estimating frequencies. *Biometrika*, **33**, 222.
 STEIN, C. & WALD, A. (1947). Sequential confidence intervals for the mean of a normal distribution with known variance. *Ann. Math. Statist.* **18**, 427.
 TWEEDIE, M. C. K. (1945). Letter in *Nature, Lond.*, **155**, 453.

A historical note on the method of least squares

By R. L. PLACKETT, *University of Liverpool*

1. The purposes of this note are:

- (i) to summarize the justifications by Laplace, Gauss and Markoff of the method of least squares;
- (ii) to suggest that Gauss was the first who justified least squares as giving those linear estimates which are unbiased of minimum variance;
- (iii) to modernize and extend his proof to cover a general theorem due to Aitken.

It is not my object to provoke controversy, and I have attempted to indicate where a personal opinion is intended.

2. The method of least squares has been in use now for over 150 years. During the nineteenth century the writings of Todhunter (1865), Merriman (1877) and others gave the impression that Laplace (1812–20, collected works 1886) was largely responsible for putting the method on a theoretical basis by means of the calculus of probability, whereas the contribution of Gauss (1821, collected works 1873) was minimized or ignored. Lately, the emphasis has changed, and in recent papers and text-books Markoff (1912) is credited with justifying the method without superfluous assumptions of normality. For these reasons, it seems desirable to disentangle the various justifications proposed and to allot credit in due proportion

3. In general let $\theta(s \times 1)$ be a vector of unknown parameters, $\mathbf{x}(n \times 1)$ a vector of observations, $\epsilon(n \times 1)$ a vector of errors and $A(n \times s)$ a matrix of known quantities; so that

$$A\theta - \mathbf{x} = \epsilon.$$

Further, suppose that $W(n \times n)$ is a diagonal matrix whose elements are the reciprocals of the error variances. It is required to form an estimate θ^* of θ . The method of least squares leads to estimates which satisfy

$$A'WA\theta^* = A'W\mathbf{x}.$$

Neither Laplace nor Gauss used matrix notation, but their results can immediately be written in that form.

4. Laplace (1812–20) discusses the method of least squares in Book 2, Chapter 4, and in the first three Supplements. He proves a series of results which are summarized—I hope fairly—in the following:

THEOREM. *Among all $s \times n$ matrices F leading to estimates of the form $FA\theta^* = Fx$, the expected values of the elements of $|\theta^* - \theta|$ are minimized as $n \rightarrow \infty$ when $F = \mu A'W$, μ being an arbitrary multiplier.*

The proof is long but runs on these lines: if u is the error of θ^* then $FAu = F\epsilon$; Laplace proceeds to determine the joint characteristic function of $F\epsilon$ and deduces that when all errors have the same distribution, symmetrical about zero, $F\epsilon$ has a multivariate normal distribution as $n \rightarrow \infty$; whence u also has a multivariate normal distribution, and the expected values of the elements of $|u|$ are determined; finally, he shows that $F = \mu A'W$ implies the vanishing of the differential coefficients of these expected values with respect to the elements of F .

In more detail, Laplace first takes $s = 1$ and maximizes the probability that his estimate lies between given limits; he then notes that this is the same as minimizing $\mathcal{E}|\theta^* - \theta|$, and continues to use this criterion when $s = 2$, stating that the result can be extended to greater values of s . In the first Supplement he considers the possibility of a bias in the observations and suggests its removal by introducing an additional parameter whose coefficient is unity in all equations.

5. Gauss presented his justification in 1821. The paper is written in Latin, but a French translation was published by Bertrand in 1855 and the fundamental theorem incorporated in Bertrand's own book of 1888. In the early sections of his paper, Gauss also considers the possibility of bias in his observations and makes it clear that the preferred estimates are those with minimum variance, although of course he does not use this terminology. He begins with errors of differing variance, and by choosing suitable multipliers presents the equations in a form where the errors have the same variance. The proof of the following theorem is in Art. 20; it is implicit that he is seeking unbiased estimates:

THEOREM. *Among all the systems of coefficients $B(s \times n)$ which give $B\epsilon = \theta - \theta^\dagger$, the estimate θ^\dagger being independent of θ , those for which the diagonal elements of BB' are minimized are provided by the method of least squares.*

Put $\xi = A'\epsilon$ so that $\xi = A'A\theta - A'x$.

The solution of these equations is $\theta = \theta^* + D\xi$, and with $DA' = E$ this becomes

$$E\epsilon = \theta - \theta^* \quad \text{so} \quad (\theta^* - \theta^\dagger) = (B - E)\epsilon.$$

If this is true for all θ then $(B - E)A = 0$, i.e. on post-multiplying by D' , $(B - E)E' = 0$, which implies $BB' = EE' + (B - E)(B - E)'$. It is now clear that the diagonal elements of BB' are minimized when $B = E$.

6. Matrix notation has been adopted for brevity in the preceding sections, but no matrix theorems have been assumed. Taking for granted the now familiar properties of matrices regarding associative products and inverses, the preceding demonstration can be modernized and shortened.

If $\theta^* = Bx$ is unbiased for all θ , then $BA = I$. With $C = A'A$ it follows that $C^{-1} = BAC^{-1}$, so

$$BB' = (C^{-1}A')(C^{-1}A')' + (B - C^{-1}A')(B - C^{-1}A')',$$

i.e. the diagonal elements of BB' are least when $B = C^{-1}A'$, which is the solution provided by least squares.

7. Markoff (1912) devotes Chapter 7 of his book to the method of least squares. He states that each observation is to be considered as a particular case of many, and as an unbiased estimate of some linear function of the unknown parameters. His determination of unbiased estimates of these parameters having minimum variance is closely followed in the paper by David & Neyman (1938).

8. Aitken (1934) has extended the theorem of Gauss by proving that with a known matrix V of variances and covariances of the observations, the minimum of $(A\theta - x)'V^{-1}(A\theta - x)$ provides estimates θ^* such that $\varphi^* = P\theta^*$ is an unbiased estimate of $\varphi = P\theta$ with minimum variance. Gauss's method can be used to prove this also.

If $\varphi^* = Bx$ then $BA = P$ and

$$[B]V[P(A'V^{-1}A)^{-1}A'V^{-1}]' = [P(A'V^{-1}A)^{-1}A'V^{-1}]V[P(A'V^{-1}A)^{-1}A'V^{-1}]',$$

consequently

$$\begin{aligned} BVB' &= [P(A'V^{-1}A)^{-1}A'V^{-1}]V[P(A'V^{-1}A)^{-1}A'V^{-1}]' \\ &\quad + [B - P(A'V^{-1}A)^{-1}A'V^{-1}]V[B - P(A'V^{-1}A)^{-1}A'V^{-1}]'. \end{aligned}$$

If we consider the diagonal elements here, the second term on the right gives a positive definite quadratic form, so minimum variance is attained when

$$B = P(A'V^{-1}A)^{-1}A'V^{-1},$$

the solution given by the method indicated above.

9. It is therefore my opinion that Laplace and Gauss proved theorems which are quite different; that the justification given by Gauss is preferable; and that Markoff, who refers to Gauss's work, may perhaps have clarified assumptions implicit there but proved nothing new. It is evident that Gauss's proof is valid for all values of n , entirely free from any assumption of normality, and capable of immediate development.

REFERENCES

- AITKEN, A. C. (1934). On least squares and linear combination of observations. *Proc. Roy. Soc. Edinb. A*, **55**, 42–7.
- BERTRAND, J. (1888). *Calcul des probabilités*. Paris.
- DAVID, F. N. & NEYMAN, J. (1938). Extension of the Markoff theorem on least squares. *Statist. Res. Mem.* **2**, 105–16.
- GAUSS, C. F. (1855). *Méthode des moindres carrés* (trans. J. Bertrand). Paris.
- GAUSS, C. F. (1873). *Theoria combinationis observationum erroribus minimis obnoxiae*. Pars prior. *Werke*, Band 4. Göttingen.
- LAPLACE, P. S., Marquis de (1886). *Théorie analytique des probabilités*, 3rd edition, Oeuvres, 7. Paris.
- MARKOFF, A. A. (1912). *Wahrscheinlichkeitsrechnung* (trans. H. Liebmann), 2nd edition. Leipzig and Berlin.
- MERRIMAN, M. (1877). A list of writings relating to the method of least squares, with historical and critical notes. *Trans. Conn. Acad. Arts Sci.* **4**, 151–232.
- TODHUNTER, I. (1865). *A History of the Mathematical Theory of Probability*. Macmillan.

The characteristic function of a weighted sum of non-central squares of normal variates subject to s linear restraints

By G. I. BATEMAN

1. Suppose x_1, x_2, \dots, x_n are independent normal variables with expectations a_1, a_2, \dots, a_n respectively and with unit variance, that is to say, we suppose that

$$p(x_j) = \frac{1}{\sqrt{(2\pi)}} \exp\left[-\frac{1}{2}(x_j - a_j)^2\right], \quad (j = 1, 2, \dots, n).$$

We consider a weighted non-central sum of squares of the type

$$\psi'^2 = \sum_{j=1}^n c_j x_j^2,$$

where the x 's are as defined and the c_j ($j = 1, 2, \dots, n$) are constants. It is assumed that the x_j are subject to s linear restraints

$$\sum_{j=1}^n b_{lj} x_j = \rho_l \quad (l = 1, 2, \dots, s),$$

b_{lj} and ρ_l ($l = 1, \dots, s$; $j = 1, 2, \dots, n$) being given constants. The characteristic function of the joint distribution of $\psi'^2, \rho_1, \rho_2, \dots, \rho_s$ may be written down immediately. We have

$$\phi(t, t_1, \dots, t_s) = \prod_{j=1}^n \left[\int_{-\infty}^{+\infty} \frac{1}{\sqrt{(2\pi)}} \exp\left\{-\frac{1}{2}(x_j - a_j)^2 + itc_j x_j^2 + \sum_{l=1}^s it_l b_{lj} x_j\right\} dx_j \right].$$

The evaluation of the integral is straightforward and we obtain

$$\phi(t, t_1, \dots, t_s) = \prod_{j=1}^n (1 - 2itc_j)^{-1} \exp\left(it \sum_{j=1}^n \frac{c_j a_j^2}{1 - 2itc_j} - \frac{1}{2} \sum_{l=1}^s \sum_{m=1}^s A_{lm} t_l t_m + i \sum_{l=1}^s B_l t_l\right), \quad (1)$$

where

$$A_{lm} = \sum_{j=1}^n \frac{b_{lj} b_{mj}}{1 - 2itc_j} \quad \text{and} \quad B_l = \sum_{j=1}^n \frac{a_j b_{lj}}{1 - 2itc_j}. \quad (2)$$