

# Least squares

From Wikipedia, the free encyclopedia

The method of **least squares** is a standard approach in regression analysis to the approximate solution of overdetermined systems, i.e., sets of equations in which there are more equations than unknowns. "Least squares" means that the overall solution minimizes the sum of the squares of the residuals made in the results of every single equation.

The most important application is in data fitting. The best fit in the least-squares sense minimizes *the sum of squared residuals* (a residual being: the difference between an observed value, and the fitted value provided by a model). When the problem has substantial uncertainties in the independent variable (the *x* variable), then simple regression and least-squares methods have problems; in such cases, the methodology required for fitting errors-in-variables models may be considered instead of that for least squares.

Least-squares problems fall into two categories: linear or ordinary least squares and nonlinear least squares, depending on whether or not the residuals are linear in all unknowns. The linear least-squares problem occurs in statistical regression analysis; it has a closed-form solution. The nonlinear problem is usually solved by iterative refinement; at each iteration the system is approximated by a linear one, and thus the core calculation is similar in both cases.

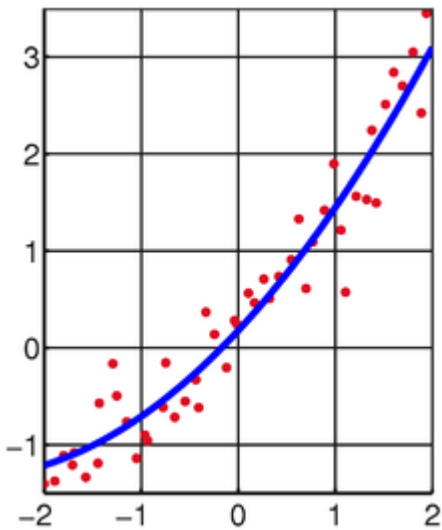
Polynomial least squares describes the variance in a prediction of the dependent variable as a function of the independent variable and the deviations from the fitted curve.

When the observations come from an exponential family and mild conditions are satisfied, least-squares estimates and maximum-likelihood estimates are identical.<sup>[1]</sup> The method of least squares can also be derived as a method of moments estimator.

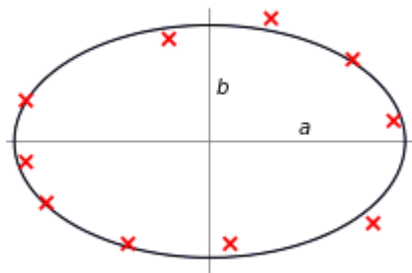
The following discussion is mostly presented in terms of linear functions but the use of least squares is valid and practical for more general families of functions. Also, by iteratively applying local quadratic approximation to the likelihood (through the Fisher information), the least-squares method may be used to fit a generalized linear model.

For the topic of approximating a function by a sum of others using an objective function based on squared distances, see least squares (function approximation).

The least-squares method is usually credited to Carl Friedrich Gauss (1795),<sup>[2]</sup> but it was first published by Adrien-Marie Legendre.<sup>[3]</sup>



The result of fitting a set of data points with a quadratic function



Conic fitting a set of points using least-squares approximation

## Contents

- 1 History
  - 1.1 Context
  - 1.2 The method
- 2 Problem statement
- 3 Limitations

- 4 Solving the least squares problem
  - 4.1 Linear least squares
  - 4.2 Non-linear least squares
  - 4.3 Differences between linear and nonlinear least squares
- 5 Least squares, regression analysis and statistics
- 6 Weighted least squares
- 7 Relationship to principal components
- 8 Regularized versions
  - 8.1 Tikhonov regularization
  - 8.2 Lasso method
- 9 See also
- 10 References
- 11 Further reading

## History

### Context

The method of least squares grew out of the fields of astronomy and geodesy, as scientists and mathematicians sought to provide solutions to the challenges of navigating the Earth's oceans during the Age of Exploration. The accurate description of the behavior of celestial bodies was the key to enabling ships to sail in open seas, where sailors could no longer rely on land sightings for navigation.

The method was the culmination of several advances that took place during the course of the eighteenth century.<sup>[4]</sup>

- The combination of different observations as being the best estimate of the true value; errors decrease with aggregation rather than increase, perhaps first expressed by Roger Cotes in 1722.
- The combination of different observations taken under the *same* conditions contrary to simply trying one's best to observe and record a single observation accurately. The approach was known as the method of averages. This approach was notably used by Tobias Mayer while studying the librations of the moon in 1750, and by Pierre-Simon Laplace in his work in explaining the differences in motion of Jupiter and Saturn in 1788.
- The combination of different observations taken under *different* conditions. The method came to be known as the method of least absolute deviation. It was notably performed by Roger Joseph Boscovich in his work on the shape of the earth in 1757 and by Pierre-Simon Laplace for the same problem in 1799.
- The development of a criterion that can be evaluated to determine when the solution with the minimum error has been achieved. Laplace tried to specify a mathematical form of the probability density for the errors and define a method of estimation that minimizes the error of estimation. For this purpose, Laplace used a symmetric two-sided exponential distribution we now call Laplace distribution to model the error distribution, and used the sum of absolute deviation as error of estimation. He felt these to be the simplest assumptions he could make, and he had hoped to obtain the arithmetic mean as the best estimate. Instead, his estimator was the posterior median.

### The method

The first clear and concise exposition of the method of least squares was published by Legendre in 1805.<sup>[5]</sup> The technique is described as an algebraic procedure for fitting linear equations to data and Legendre demonstrates the new method by analyzing the same data as Laplace for the shape of the earth. The value of Legendre's method of least squares was immediately recognized by leading astronomers and geodesists of the time.

In 1809 Carl Friedrich Gauss published his method of calculating the orbits of celestial bodies. In that work he claimed to have been in possession of the method of least squares since 1795. This naturally led to a priority dispute with Legendre. However, to Gauss's credit, he went beyond Legendre and succeeded in connecting the method of least squares with the principles of probability and to the normal distribution. He had managed to

complete Laplace's program of specifying a mathematical form of the probability density for the observations, depending on a finite number of unknown parameters, and define a method of estimation that minimizes the error of estimation. Gauss showed that arithmetic mean is indeed the best estimate of the location parameter by changing both the probability density and the method of estimation. He then turned the problem around by asking what form the density should have and what method of estimation should be used to get the arithmetic mean as estimate of the location parameter. In this attempt, he invented the normal distribution.



Carl Friedrich Gauss

An early demonstration of the strength of Gauss' method came when it was used to predict the future location of the newly discovered asteroid Ceres. On 1 January 1801, the Italian astronomer Giuseppe Piazzi discovered Ceres and was able to track its path for 40 days before it was lost in the glare of the sun. Based on these data, astronomers desired to determine the location of Ceres after it emerged from behind the sun without solving Kepler's complicated nonlinear equations of planetary motion. The only predictions that successfully allowed Hungarian astronomer Franz Xaver von Zach to relocate Ceres were those performed by the 24-year-old Gauss using least-squares analysis.

In 1810, after reading Gauss's work, Laplace, after proving the central limit theorem, used it to give a large sample justification for the method of least square and the normal distribution. In 1822, Gauss was able to state that the least-squares approach to regression analysis is optimal in the sense that in a linear model where the errors have a mean of zero, are uncorrelated, and have equal variances, the best linear unbiased estimator of the coefficients is the least-squares estimator. This result is known as the Gauss–Markov theorem.

The idea of least-squares analysis was also independently formulated by the American Robert Adrain in 1808. In the next two centuries workers in the theory of errors and in statistics found many different ways of implementing least squares.<sup>[6]</sup>

## Problem statement

The objective consists of adjusting the parameters of a model function to best fit a data set. A simple data set consists of  $n$  points (data pairs)  $(\mathbf{x}_i, \mathbf{y}_i)$ ,  $i = 1, \dots, n$ , where  $\mathbf{x}_i$  is an independent variable and  $\mathbf{y}_i$  is a dependent variable whose value is found by observation. The model function has the form  $\mathbf{f}(\mathbf{x}, \boldsymbol{\beta})$ , where  $m$  adjustable parameters are held in the vector  $\boldsymbol{\beta}$ . The goal is to find the parameter values for the model that "best" fits the data. The least-squares method finds its optimum when the sum,  $S$ , of squared residuals

$$S = \sum_{i=1}^n r_i^2$$

is a minimum. A residual is defined as the difference between the actual value of the dependent variable and the value predicted by the model. Each data point has one residual. Both the sum and the mean of the residuals are equal to zero.

$$r_i = y_i - f(x_i, \boldsymbol{\beta}).$$

An example of a model is that of the straight line in two dimensions. Denoting the y-intercept as  $\beta_0$  and the slope as  $\beta_1$ , the model function is given by  $\mathbf{f}(\mathbf{x}, \boldsymbol{\beta}) = \beta_0 + \beta_1 \mathbf{x}$ . See linear least squares for a fully worked out example of this model.

A data point may consist of more than one independent variable. For example, when fitting a plane to a set of height measurements, the plane is a function of two independent variables,  $x$  and  $z$ , say. In the most general case there may be one or more independent variables and one or more dependent variables at each data point.

## Limitations

This regression formulation considers only residuals in the dependent variable. There are two rather different contexts in which different implications apply:

- Regression for prediction. Here a model is fitted to provide a prediction rule for application in a similar situation to which the data used for fitting apply. Here the dependent variables corresponding to such future application would be subject to the same types of observation error as those in the data used for fitting. It is therefore logically consistent to use the least-squares prediction rule for such data.
- Regression for fitting a "true relationship". In standard regression analysis, that leads to fitting by least squares, there is an implicit assumption that errors in the independent variable are zero or strictly controlled so as to be negligible. When errors in the independent variable are non-negligible, models of measurement error can be used; such methods can lead to parameter estimates, hypothesis testing and confidence intervals that take into account the presence of observation errors in the independent variables.<sup>[7]</sup> An alternative approach is to fit a model by total least squares; this can be viewed as taking a pragmatic approach to balancing the effects of the different sources of error in formulating an objective function for use in model-fitting.

## Solving the least squares problem

The minimum of the sum of squares is found by setting the gradient to zero. Since the model contains  $m$  parameters, there are  $m$  gradient equations:

$$\frac{\partial S}{\partial \beta_j} = 2 \sum_i r_i \frac{\partial r_i}{\partial \beta_j} = 0, \quad j = 1, \dots, m,$$

and since  $r_i = y_i - f(x_i, \beta)$ , the gradient equations become

$$-2 \sum_i r_i \frac{\partial f(x_i, \beta)}{\partial \beta_j} = 0, \quad j = 1, \dots, m.$$

The gradient equations apply to all least squares problems. Each particular problem requires particular expressions for the model and its partial derivatives.

## Linear least squares

A regression model is a linear one when the model comprises a linear combination of the parameters, i.e.,

$$f(x, \beta) = \sum_{j=1}^m \beta_j \phi_j(x),$$

where the function  $\phi_j$  is a function of  $x$ .

Letting

$$X_{ij} = \frac{\partial f(x_i, \beta)}{\partial \beta_j} = \phi_j(x_i),$$

we can then see that in that case the least square estimate (or estimator, in the context of a random sample),  $\beta$  is given by

$$\hat{\beta} = (X^T X)^{-1} X^T y.$$

For a derivation of this estimate see Linear least squares (mathematics).

## Non-linear least squares

There is, in some cases, a closed-form solution to a non-linear least squares problem – but in general there is not. In the case of no closed-form solution, numerical algorithms are used to find the value of the parameters  $\beta$  that minimizes the objective. Most algorithms involve choosing initial values for the parameters. Then, the parameters are refined iteratively, that is, the values are obtained by successive approximation:

$$\beta_j^{k+1} = \beta_j^k + \Delta\beta_j,$$

where a superscript  $k$  is an iteration number, and the vector of increments  $\Delta\beta_j$  is called the shift vector. In some commonly used algorithms, at each iteration the model may be linearized by approximation to a first-order Taylor series expansion about  $\beta^k$ :

$$\begin{aligned} f(x_i, \beta) &= f^k(x_i, \beta) + \sum_j \frac{\partial f(x_i, \beta)}{\partial \beta_j} (\beta_j - \beta_j^k) \\ &= f^k(x_i, \beta) + \sum_j J_{ij} \Delta\beta_j. \end{aligned}$$

The Jacobian  $\mathbf{J}$  is a function of constants, the independent variable *and* the parameters, so it changes from one iteration to the next. The residuals are given by

$$r_i = y_i - f^k(x_i, \beta) - \sum_{k=1}^m J_{ik} \Delta\beta_k = \Delta y_i - \sum_{j=1}^m J_{ij} \Delta\beta_j.$$

To minimize the sum of squares of  $r_i$ , the gradient equation is set to zero and solved for  $\Delta\beta_j$ :

$$-2 \sum_{i=1}^n J_{ij} \left( \Delta y_i - \sum_{k=1}^m J_{ik} \Delta\beta_k \right) = 0,$$

which, on rearrangement, become  $m$  simultaneous linear equations, the **normal equations**:

$$\sum_{i=1}^n \sum_{k=1}^m J_{ij} J_{ik} \Delta\beta_k = \sum_{i=1}^n J_{ij} \Delta y_i \quad (j = 1, \dots, m).$$

The normal equations are written in matrix notation as

$$(\mathbf{J}^T \mathbf{J}) \Delta\beta = \mathbf{J}^T \Delta\mathbf{y}.$$

These are the defining equations of the Gauss–Newton algorithm.

## Differences between linear and nonlinear least squares

- The model function,  $f$ , in LLSQ (linear least squares) is a linear combination of parameters of the form  $f = X_{i1}\beta_1 + X_{i2}\beta_2 + \dots$ . The model may represent a straight line, a parabola or any other linear

combination of functions. In NLLSQ (nonlinear least squares) the parameters appear as functions, such as  $\beta^2$ ,  $e^{\beta x}$  and so forth. If the derivatives  $\partial f / \partial \beta_j$  are either constant or depend only on the values of the independent variable, the model is linear in the parameters. Otherwise the model is nonlinear.

- Algorithms for finding the solution to a NLLSQ problem require initial values for the parameters, LLSQ does not.
- Like LLSQ, solution algorithms for NLLSQ often require that the Jacobian can be calculated. Analytical expressions for the partial derivatives can be complicated. If analytical expressions are impossible to obtain either the partial derivatives must be calculated by numerical approximation or an estimate must be made of the Jacobian.
- In NLLSQ non-convergence (failure of the algorithm to find a minimum) is a common phenomenon whereas the LLSQ is globally concave so non-convergence is not an issue.
- NLLSQ is usually an iterative process. The iterative process has to be terminated when a convergence criterion is satisfied. LLSQ solutions can be computed using direct methods, although problems with large numbers of parameters are typically solved with iterative methods, such as the Gauss–Seidel method.
- In LLSQ the solution is unique, but in NLLSQ there may be multiple minima in the sum of squares.
- Under the condition that the errors are uncorrelated with the predictor variables, LLSQ yields unbiased estimates, but even under that condition NLLSQ estimates are generally biased.

These differences must be considered whenever the solution to a nonlinear least squares problem is being sought.

## Least squares, regression analysis and statistics

The method of least squares is often used to generate estimators and other statistics in regression analysis.

Consider a simple example drawn from physics. A spring should obey Hooke's law which states that the extension of a spring  $y$  is proportional to the force,  $F$ , applied to it.

$$y = f(F, k) = kF$$

constitutes the model, where  $F$  is the independent variable. To estimate the force constant,  $k$ , a series of  $n$  measurements with different forces will produce a set of data,  $(F_i, y_i)$ ,  $i = 1, \dots, n$ , where  $y_i$  is a measured spring extension. Each experimental observation will contain some error. If we denote this error  $\varepsilon$ , we may specify an empirical model for our observations,

$$y_i = kF_i + \varepsilon_i.$$

There are many methods we might use to estimate the unknown parameter  $k$ . Noting that the  $n$  equations in the  $m$  variables in our data comprise an overdetermined system with one unknown and  $n$  equations, we may choose to estimate  $k$  using least squares. The sum of squares to be minimized is

$$S = \sum_{i=1}^n (y_i - kF_i)^2.$$

The least squares estimate of the force constant,  $k$ , is given by

$$\hat{k} = \frac{\sum_i F_i y_i}{\sum_i F_i^2}.$$

Here it is assumed that application of the force **causes** the spring to expand and, having derived the force constant by least squares fitting, the extension can be predicted from Hooke's law.

In regression analysis the researcher specifies an empirical model. For example, a very common model is the straight line model which is used to test if there is a linear relationship between dependent and independent variable. If a linear relationship is found to exist, the variables are said to be correlated. However, correlation does not prove causation, as both variables may be correlated with other, hidden, variables, or the dependent variable may "reverse" cause the independent variables, or the variables may be otherwise spuriously correlated. For example, suppose there is a correlation between deaths by drowning and the volume of ice cream sales at a particular beach. Yet, both the number of people going swimming and the volume of ice cream sales increase as the weather gets hotter, and presumably the number of deaths by drowning is correlated with the number of people going swimming. Perhaps an increase in swimmers causes both the other variables to increase.

In order to make statistical tests on the results it is necessary to make assumptions about the nature of the experimental errors. A common (but not necessary) assumption is that the errors belong to a normal distribution. The central limit theorem supports the idea that this is a good approximation in many cases.

- The Gauss–Markov theorem. In a linear model in which the errors have expectation zero conditional on the independent variables, are uncorrelated and have equal variances, the best linear unbiased estimator of any linear combination of the observations, is its least-squares estimator. "Best" means that the least squares estimators of the parameters have minimum variance. The assumption of equal variance is valid when the errors all belong to the same distribution.
- In a linear model, if the errors belong to a normal distribution the least squares estimators are also the maximum likelihood estimators.

However, if the errors are not normally distributed, a central limit theorem often nonetheless implies that the parameter estimates will be approximately normally distributed so long as the sample is reasonably large. For this reason, given the important property that the error mean is independent of the independent variables, the distribution of the error term is not an important issue in regression analysis. Specifically, it is not typically important whether the error term follows a normal distribution.

In a least squares calculation with unit weights, or in linear regression, the variance on the  $j$ th parameter, denoted  $\text{var}(\hat{\beta}_j)$ , is usually estimated with

$$\text{var}(\hat{\beta}_j) = \sigma^2 ([X^T X]^{-1})_{jj} \approx \frac{S}{n - m} ([X^T X]^{-1})_{jj},$$

where the true error variance  $\sigma^2$  is replaced by an estimate based on the minimised value of the sum of squares objective function  $S$ . The denominator,  $n - m$ , is the statistical degrees of freedom; see effective degrees of freedom for generalizations.

Confidence limits can be found if the probability distribution of the parameters is known, or an asymptotic approximation is made, or assumed. Likewise statistical tests on the residuals can be made if the probability distribution of the residuals is known or assumed. The probability distribution of any linear combination of the dependent variables can be derived if the probability distribution of experimental errors is known or assumed. Inference is particularly straightforward if the errors are assumed to follow a normal distribution, which implies that the parameter estimates and residuals will also be normally distributed conditional on the values of the independent variables.

## Weighted least squares

A special case of generalized least squares called **weighted least squares** occurs when all the off-diagonal entries of  $\Omega$  (the correlation matrix of the residuals) are null; the variances of the observations (along the covariance matrix diagonal) may still be unequal (heteroscedasticity).

The expressions given above are based on the implicit assumption that the errors are uncorrelated with each other and with the independent variables and have equal variance. The Gauss–Markov theorem shows that, when this is so,  $\hat{\boldsymbol{\beta}}$  is a best linear unbiased estimator (BLUE). If, however, the measurements are uncorrelated but have different uncertainties, a modified approach might be adopted. Aitken showed that when a weighted sum of squared residuals is minimized,  $\hat{\boldsymbol{\beta}}$  is the BLUE if each weight is equal to the reciprocal of the variance of the measurement

$$S = \sum_{i=1}^n W_{ii} r_i^2, \quad W_{ii} = \frac{1}{\sigma_i^2}$$

The gradient equations for this sum of squares are

$$-2 \sum_i W_{ii} \frac{\partial f(x_i, \boldsymbol{\beta})}{\partial \beta_j} r_i = 0, \quad j = 1, \dots, n$$

which, in a linear least squares system give the modified normal equations,

$$\sum_{i=1}^n \sum_{k=1}^m X_{ij} W_{ii} X_{ik} \hat{\beta}_k = \sum_{i=1}^n X_{ij} W_{ii} y_i, \quad j = 1, \dots, m.$$

When the observational errors are uncorrelated and the weight matrix,  $\mathbf{W}$ , is diagonal, these may be written as

$$(\mathbf{X}^T \mathbf{W} \mathbf{X}) \hat{\boldsymbol{\beta}} = \mathbf{X}^T \mathbf{W} \mathbf{y}.$$

If the errors are correlated, the resulting estimator is the BLUE if the weight matrix is equal to the inverse of the variance-covariance matrix of the observations.

When the errors are uncorrelated, it is convenient to simplify the calculations to factor the weight matrix as  $w_{ii} = \sqrt{W_{ii}}$ . The normal equations can then be written in the same form as ordinary least squares:

$$(\mathbf{X}'^T \mathbf{X}') \hat{\boldsymbol{\beta}} = \mathbf{X}'^T \mathbf{y}'$$

where we define the following scaled matrix and vector:

$$\begin{aligned} \mathbf{X}' &= \text{diag}(\mathbf{w}) \mathbf{X}, \\ \mathbf{y}' &= \text{diag}(\mathbf{w}) \mathbf{y} = \mathbf{y} \oslash \boldsymbol{\sigma}. \end{aligned}$$

This is a type of whitening transformation; the last expression involves an entrywise division.

For non-linear least squares systems a similar argument shows that the normal equations should be modified as follows.

$$(\mathbf{J}^T \mathbf{W} \mathbf{J}) \Delta \boldsymbol{\beta} = \mathbf{J}^T \mathbf{W} \Delta \mathbf{y}.$$

Note that for empirical tests, the appropriate  $\mathbf{W}$  is not known for sure and must be estimated. For this feasible generalized least squares (FGLS) techniques may be used.

## Relationship to principal components

The first principal component about the mean of a set of points can be represented by that line which most closely approaches the data points (as measured by squared distance of closest approach, i.e. perpendicular to the line). In contrast, linear least squares tries to minimize the distance in the  $\mathbf{y}$  direction only. Thus, although



the two use a similar error metric, linear least squares is a method that treats one dimension of the data preferentially, while PCA treats all dimensions equally.

## Regularized versions

### Tikhonov regularization

In some contexts a regularized version of the least squares solution may be preferable. Tikhonov regularization (or ridge regression) adds a constraint that  $\|\beta\|^2$ , the  $L_2$ -norm of the parameter vector, is not greater than a given value. Equivalently, it may solve an unconstrained minimization of the least-squares penalty with  $\alpha\|\beta\|^2$  added, where  $\alpha$  is a constant (this is the Lagrangian form of the constrained problem). In a Bayesian context, this is equivalent to placing a zero-mean normally distributed prior on the parameter vector.

### Lasso method

An alternative regularized version of least squares is *Lasso* (least absolute shrinkage and selection operator), which uses the constraint that  $\|\beta\|$ , the  $L_1$ -norm of the parameter vector, is no greater than a given value.<sup>[8][9][10]</sup> (As above, this is equivalent to an unconstrained minimization of the least-squares penalty with  $\alpha\|\beta\|$  added.) In a Bayesian context, this is equivalent to placing a zero-mean Laplace prior distribution on the parameter vector.<sup>[11]</sup> The optimization problem may be solved using quadratic programming or more general convex optimization methods, as well as by specific algorithms such as the least angle regression algorithm.


One of the prime differences between Lasso and ridge regression is that in ridge regression, as the penalty is increased, all parameters are reduced while still remaining non-zero, while in Lasso, increasing the penalty will cause more and more of the parameters to be driven to zero. This is an advantage of Lasso over ridge regression, as driving parameters to zero deselects the features from the regression. Thus, Lasso automatically selects more relevant features and discards the others, whereas Ridge regression never fully discards any features. Some feature selection techniques are developed based on the LASSO including Bolasso which bootstraps samples,<sup>[12]</sup> and FeaLect which analyzes the regression coefficients corresponding to different values of  $\alpha$  to score all the features.<sup>[13]</sup>

The  $L^1$ -regularized formulation is useful in some contexts due to its tendency to prefer solutions where more parameters are zero, which gives solutions that depend on fewer variables.<sup>[8]</sup> For this reason, the Lasso and its variants are fundamental to the field of compressed sensing. An extension of this approach is elastic net regularization.

## See also

- Adjustment of observations
- Bayesian MMSE estimator
- Best linear unbiased estimator (BLUE)
- Best linear unbiased prediction (BLUP)
- Gauss–Markov theorem
- $L_2$  norm
- Least absolute deviation
- Measurement uncertainty
- Orthogonal projection
- Proximal gradient methods for learning
- Quadratic loss function
- Root mean square
- Squared deviations

## References

- Charnes, A.; Frome, E. L.; Yu, P. L. (1976). "The Equivalence of Generalized Least Squares and Maximum Likelihood Estimates in the Exponential Family". *Journal of the American Statistical Association*. **71** (353): 169–171. doi:10.1080/01621459.1976.10481508 (https://doi.org/10.1080%2F01621459.1976.10481508).
- Bretscher, Otto (1995). *Linear Algebra With Applications* (3rd ed.). Upper Saddle River, NJ: Prentice Hall.
- Stigler, Stephen M. (1981). "Gauss and the Invention of Least Squares" (http://projecteuclid.org/euclid.aos/1176345451). *Ann. Stat.* **9** (3): 465–474. doi:10.1214/aos/1176345451 (https://doi.org/10.1214%2Faoas/1176345451).
- Stigler, Stephen M. (1986). *The History of Statistics: The Measurement of Uncertainty Before 1900*. Cambridge, MA: Belknap Press of Harvard University Press. ISBN 0-674-40340-1.
- Legendre, Adrien-Marie (1805), *Nouvelles méthodes pour la détermination des orbites des comètes* (http://books.google.com/books/about/Nouvelles\_m%C3%A9thodes\_pour\_la\_d%C3%A9terminati.html?id=FRcOAAAAQAAJ) [*New Methods for the Determination of the Orbits of Comets*] (in French), Paris: F. Didot
- Aldrich, J. (1998). "Doing Least Squares: Perspectives from Gauss and Yule". *International Statistical Review*. **66** (1): 61–81. doi:10.1111/j.1751-5823.1998.tb00406.x (https://doi.org/10.1111%2Fj.1751-5823.1998.tb00406.x).
- For a good introduction to error-in-variables, please see Fuller, W. A. (1987). *Measurement Error Models*. John Wiley & Sons. ISBN 0-471-86187-1.
- Tibshirani, R. (1996). "Regression shrinkage and selection via the lasso". *Journal of the Royal Statistical Society, Series B*. **58** (1): 267–288. JSTOR 2346178 (https://www.jstor.org/stable/2346178).
- Hastie, Trevor; Tibshirani, Robert; Friedman, Jerome H. (2009). "The Elements of Statistical Learning" (http://www-stat.stanford.edu/~tibs/ElemStatLearn/) (second ed.). Springer-Verlag. ISBN 978-0-387-84858-7.
- Bühlmann, Peter; van de Geer, Sara (2011). *Statistics for High-Dimensional Data: Methods, Theory and Applications*. Springer. ISBN 9783642201929.
- Park, Trevor; Casella, George (2008). "The Bayesian Lasso". *Journal of the American Statistical Association*. **103** (482): 681–686. doi:10.1198/016214508000000337 (https://doi.org/10.1198%2F016214508000000337).
- Bach, Francis R (2008). "Bolasso: model consistent lasso estimation through the bootstrap" (http://dl.acm.org/citation.cfm?id=1390161). *Proceedings of the 25th international conference on Machine learning*: 33–40. doi:10.1145/1390156.1390161 (https://doi.org/10.1145%2F1390156.1390161).
- Zare, Habil (2013). "Scoring relevancy of features based on combinatorial analysis of Lasso with application to lymphoma diagnosis" (http://www.biomedcentral.com/1471-2164/14/S1/S14). *BMC Genomics*. **14**: S14. PMC 3549810 (https://www.ncbi.nlm.nih.gov/pmc/articles/PMC3549810)  PMID 23369194 (https://www.ncbi.nlm.nih.gov/pubmed/23369194). doi:10.1186/1471-2164-14-S1-S14 (https://doi.org/10.1186%2F1471-2164-14-S1-S14).

## Further reading

- Björck, Å. (1996). *Numerical Methods for Least Squares Problems*. SIAM. ISBN 0-89871-360-9.
- Kariya, T.; Kurata, H. (2004). *Generalized Least Squares*. Hoboken: Wiley. ISBN 0-470-86697-7.
- Luenberger, D. G. (1997) [1969]. "Least-Squares Estimation" (https://books.google.com/books?id=IZU0CAH4RccC&pg=PA78). *Optimization by Vector Space Methods*. New York: John Wiley & Sons. pp. 78–102. ISBN 0-471-18117-X.
- Rao, C. R.; Toutenburg, H.; et al. (2008). *Linear Models: Least Squares and Alternatives* (https://books.google.com/books?id=3LK9JoGEyN4C). Springer Series in Statistics (3rd ed.). Berlin: Springer. ISBN 978-3-540-74226-5.
- Wolberg, J. (2005). *Data Analysis Using the Method of Least Squares: Extracting the Most Information from Experiments*. Berlin: Springer. ISBN 3-540-25674-1.

Retrieved from "https://en.wikipedia.org/w/index.php?title=Least\_squares&oldid=794892465"

- This page was last edited on 10 August 2017, at 17:37.
- Text is available under the Creative Commons Attribution-ShareAlike License; additional terms may apply. By using this site, you agree to the Terms of Use and Privacy Policy. Wikipedia® is a registered trademark of the Wikimedia Foundation, Inc., a non-profit organization.