

On the relationship between the Method of Least Squares and Gram-Schmidt orthogonalization

Hilmar Drygas, University of Kassel, Germany

Summary The method of Least Squares is due to Carl Friedrich Gauss. The Gram-Schmidt orthogonalization method is of much younger date. A method for solving Least Squares Problems is developed which automatically results in the appearance of the Gram-Schmidt orthogonalizers. Given these orthogonalizers an induction-proof is available for solving Least Squares Problems.

Keywords: Linear models, Method of Least Squares, simple regression, Steiner-theorem, orthogonal projection, Gram-Schmidt orthogonalization.

1 Introduction

The method of Least consist in the following problem. Given vectors $y, x_1, \dots, x_k \in \mathbb{R}^n$ find numbers β_1, \dots, β_k such that

$$\left\| y - \sum_{i=1}^k \beta_i x_i \right\| \quad (1)$$

is minimized. The underlying linear model is $y = x_1\beta_1 + \dots + x_k\beta_k + \epsilon$, where ϵ is a disturbance term. Mostly, it is assumed that ϵ is a random vector with expectation 0 and covariance-matrix $\sigma^2 I_n$, where $\sigma > 0$ is unknown parameter. The method of Least Squares is therefore also described by

$$\|\epsilon\| = \text{Min}. \quad (2)$$

The simplest linear model is $y = a \mathbf{1}_n + \epsilon$, where $\mathbf{1}_n$ is the all one-vector. The Least Squares Problem for estimating a can be solved by Steiners theorem.

1.1 Steiners Theorem:

$$\sum_{i=1}^n w_i (y_i - a)^2 = \sum_{i=1}^n w_i (y_i - \bar{y}_{wgh})^2 + \left(\sum_{i=1}^n w_i \right) (a - \bar{y}_{wgh})^2, \text{ where } w_i \geq 0, \sum_{i=1}^n w_i > 0.$$

$$y_{wgh} = \left(\sum_{i=1}^n w_i \right)^{-1} \left(\sum_{i=1}^n w_i y_i \right).$$

The proof follows from Pythagoras theorem since

$$\sum_{i=1}^n w_i (y_i - \bar{y}_{wgh})(a - \bar{y}_{wgh}) = 0. \quad (3)$$

Thus $a = \bar{y}_{wgh}$ solves the Least Squares Problem $\sum_{i=1}^n w_i (y_i - a)^2 = \text{Min}$.

This theorem can also be used to solve the regression model $y_i = \alpha + \beta x_i + \epsilon_i$, $i = 1, \dots, n$, $y = \alpha 1_n + \beta x + \epsilon$. The task consists in minimizing

$$Q = \sum_{i=1}^n (y_i - \alpha - \beta x_i)^2. \quad (4)$$

By Steiners Theorem we get the solution

$$\hat{\alpha} = \bar{y} - \beta \bar{x}, \quad \bar{x} = \frac{1}{n} \sum_{i=1}^n x_i, \quad \bar{y} = \frac{1}{n} \sum_{i=1}^n y_i. \quad (5)$$

By plugging in we get

$$\begin{aligned} Q &= \sum_{i=1}^n \left(y_i - \bar{y} - \beta(x_i - \bar{x}) \right)^2 \\ &= \sum_{i=x_i \neq \bar{x}} (x_i - \bar{x})^2 \left(\frac{y_i - \bar{y}}{x_i - \bar{x}} - \beta \right)^2 + \sum_{i:x_i = \bar{x}} (y_i - \bar{y})^2. \end{aligned} \quad (6)$$

According to Steiners theorem the minimizing β is a weighted mean of the slopes $\frac{y_i - \bar{y}}{x_i - \bar{x}}$, namely

$$\hat{\beta} = \frac{\sum_{i=x_i \neq \bar{x}} (x_i - \bar{x})^2 \frac{(y_i - \bar{y})}{(x_i - \bar{x})}}{\sum_{i=x_i \neq \bar{x}} (x_i - \bar{x})^2} = \frac{\sum_{i=1}^n (x_i - \bar{x}) y_i}{\sum_{i=1}^n (x_i - \bar{x})^2}. \quad (7)$$

If all x_i are equal to \bar{x} , then β is arbitrary since Q does not depend on β .

This method of successive solution and plugging in can be extended to the general case as will be shown in the next section.

2 Generalization of successive estimation

2.1 Generalized Steiner Theorem:

$$\|y - ax\|^2 = \left\| y - \frac{(y,x)}{(x,x)}x \right\|^2 + \|x\|^2 \left(a - \frac{(x,y)}{(x,x)} \right)^2 \text{ if } x \neq 0.$$

Proof: $\left(y - \frac{(y,x)}{(x,x)}x \right)$ and x are orthogonal. Pythagoras Theorem therefore yields the result. \square

2.2 Corollary:

$a = \frac{(y,x)}{(x,x)}$ is the Least Squares-solution of $\|y - ax\| = \text{Min.}$ and $a = 0$ yields a proof of the Cauchy-Schwarz inequality.

Now we want to minimize

$$\left\| y - \sum_{i=1}^k \beta_i x_i \right\|^2. \quad (1)$$

If $x_1 = 0$, then β_1 does not appear in (2.1) and it is therefore arbitrary. If $x_1 \neq 0$, then according to theorem 2.1

$$\hat{\beta}_1 = \frac{(y - \sum_{i=2}^k \beta_i x_i, x_1)}{(x_1, x_1)}. \quad (2)$$

By plugging in we get the new minimization problem

$$\left\| y^{(2)} - \sum_{i=2}^k \beta_i x_i^{(2)} \right\| = \text{Min}, \quad (3)$$

where

$$y^{(2)} = y - \frac{(y, x_1)}{(x_1, x_1)}x_1 = P_{\{x_1\}^\perp}y, \quad x_i^{(2)} = x_i - \frac{(x_i, x_1)}{(x_1, x_1)}x_1 = P_{\{x_1\}^\perp}x_i. \quad (4)$$

If $x_2^{(2)} \neq 0$ - otherwise β_2 is arbitrary - we obtain

$$\hat{\beta}_2 = \frac{(y^{(2)} - \sum_{i=3}^k \beta_i x_i^{(2)}, x_2^{(2)})}{(x_2^{(2)}, x_2^{(2)})} = \frac{(y - \sum_{i=3}^k \beta_i x_i, x_2^{(2)})}{(x_2^{(2)}, x_2^{(2)})}. \quad (5)$$

and again by plugging in we get a new problem with $y^{(3)}, x_i^{(3)}, i = 3, \dots, k$. Continuing we get successively the solutions ($j = 3, \dots, k$)

$$\hat{\beta}_j = \frac{(y - \sum_{i=j+1}^k \beta_i x_i^{(j)}, x_j^{(j)})}{(x_j^{(j)}, x_j^{(j)})}, \text{ if } x_j^{(j)} \neq 0 \quad (6)$$

and finally

$$\hat{\beta}_k = \frac{(y, x_k^{(k)})}{(x_k^{(k)}, x_k^{(k)})}, \text{ if } x_k^{(k)} \neq 0. \quad (7)$$

In order to simplify the notation we define

$$q_1 = x_1, q_j = x_j^{(j)}, j = 2, \dots, k. \quad (8)$$

Then

$$\begin{aligned} x_i^{(l)} &= x_i^{(l-1)} - \frac{(x_i^{(l-1)}, q_{l-1})}{(q_{l-1}, q_{l-1})} q_{l-1} \\ &= P_{\{q_{l-1}\}^\perp} x_i^{(l-1)}, i = l, \dots, k, l = 1, \dots, k_i \end{aligned} \quad (9)$$

where, of course, $x_i^{(1)} = x_i, i = 1, \dots, k$. Therefore

$$q_l = P_{\{q_{l-1}\}^\perp} x_i^{(l-1)} \quad (10)$$

and

$$x_i^{(l)} = P_{\{q_{l-1}\}^\perp} P_{\{q_{l-2}\}^\perp} \dots P_{\{q_1\}^\perp} x_i \quad (11)$$

$$q_l = P_{\{q_{l-1}\}^\perp} \dots P_{\{q_1\}^\perp} x_l, l = 2, \dots, k. \quad (12)$$

The next step consists in proving that

$$\prod_{j=1}^{i-1} P_{\{q_{i-j}\}^\perp} = P_{\{q_1, \dots, q_{i-1}\}^\perp}. \quad (13)$$

By Achieser/Glasman, 1981, p. 97 pp. the product of projections is a projector iff the projectors commute. By page 189 in Rao/Mitra, 1971, the projection onto the intersection of the subspaces M and N is given by

$$2P(P+Q)^-Q \quad (14)$$

where P is the projection onto M and Q the projection onto N . There must be a simple formula for the generalized inverse of $(P+Q)$, namely $(P+Q)^+$. This formula will be given by the following theorem:

2.3 Theorem:

If $PQ = QP$, then $(P+Q)^+ = P+Q - \frac{3}{2}PQ$.

Proof: The proof follows from verification. An alternative is that P and Q are jointly diagonalizable if $PQ = QP$. $P = C \text{diag}(\lambda_1, \dots, \lambda_n)C'$, $Q = C \text{diag}(\mu_1, \dots, \mu_n)C'$ and the λ_i and μ_i are either 0 or 1. Then $P+Q = C(\text{diag}(\lambda_1 + \mu_1), \dots, (\lambda_n + \mu_n))C'$, $(P+Q)^+ = C \text{diag}((\lambda_1 + \mu_1)^+, \dots, (\lambda_n + \mu_n)^+)C'$. But

$$(\lambda_i + \mu_i)^+ = \lambda_i + \mu_i - \frac{3}{2} \lambda_i \mu_i \quad (15)$$

in all possible cases. \square

2.4 Theorem:

PQ is the projection onto $\text{im}(P) \cap \text{im}(Q)$ iff $QM^\perp \subseteq M^\perp$. Sufficient for this is $M^\perp \subseteq N$.

Proof: PQ is the projection onto $M \cap N$ iff it is the identity on $N \cap M$ and vanishes on $(M \cap N)^\perp = M^\perp + N^\perp$. Since the other properties are obvious only $PQM^\perp = 0$ must be considered. This is equivalent to $QM^\perp \subseteq M^\perp$. This condition met if $M^\perp \subseteq N$. \square

2.5 Theorem:

$$\prod_{j=1}^{i-1} P_{\{q_{i-j}\}^\perp} = P_{\{q_1, \dots, q_{i-1}\}^\perp} \quad \text{and} \quad q_i \in \{q_1, \dots, q_{i-1}\}^\perp. \quad (16)$$

Proof: Mathematical induction. The first assertion of the theorem is correct for $i = 2$ and $q_2 = P_{\{q_1\}^\perp} x_2 \in \{q_1\}^\perp$. Let by induction assumption

$$\prod_{j=1}^{i-1} P_{\{q_j\}^\perp} = P_{\{q_1, \dots, q_{i-1}\}^\perp} \quad \text{and} \quad q_i \in \{q_1, \dots, q_{i-1}\}^\perp. \quad (17)$$

Then

$$\prod_{j=1}^i P_{\{q_{i-j}\}^\perp} = P_{\{q_i\}^\perp} P_{\{q_1, \dots, q_{i-1}\}^\perp}. \quad (18)$$

Since $q_i \in \{q_1, \dots, q_{i-1}\}^\perp$ it follows from theorem ($M = \{q_i\}^\perp$, $M^\perp = \{\lambda q_i; \lambda \in \mathbb{R}\}$) that

$$\prod_{j=1}^i P_{\{q_j\}^\perp} = P_{\{q_i\}^\perp} P_{\{q_1, \dots, q_{i-1}\}^\perp} = P_{\{q_i\}^\perp \cap \{q_1, \dots, q_{i-1}\}^\perp} = P_{\{q_1, \dots, q_i\}^\perp}. \quad (19)$$

Since $q_{i+1} = P_{\{q_1, \dots, q_i\}^\perp} x_{i+1} \in \{q_1, \dots, q_i\}^\perp$ also the second assertion is proved. \square

2.6 Corollary:

If $q_0 = 0$, then $q_i = P_{\{q_0, \dots, q_{i-1}\}^\perp} x_i, i=1, \dots, k$ and $x_i^{(l)} = P_{\{q_0, q_1, \dots, q_{l-1}\}^\perp} x_i$.

Since

$$q_i = P_{\{q_0, \dots, q_{i-1}\}^\perp} x_i = x_i - P_{\text{span}\{q_1, \dots, q_{i-1}\}} x_i = x_i - \sum_{j: q_j \neq 0}^{i-1} \frac{(q_j, x_i)}{(q_j, q_j)} q_j \quad (20)$$

the q_i describe the Gram-Schmidt orthogonalization procedure. It follows that from the principle of Least Squares the Gram-Schmidt orthogonalization procedure could be invented.

3 An induction proof

Since the Gram-Schmidt orthogonalization procedure is well-known now the Least Squares Solutions can also be proved by mathematical induction. The induction is on the number m of linear independent vectors among x_1, \dots, x_k . We assume that x_1, \dots, x_m are linearly independent and $\text{Rank} \{x_1, \dots, x_k\} = m$. Therefore x_{m+1}, \dots, x_k are linear combinations of x_1, \dots, x_m . As we have seen in the last section

$$\hat{\beta}_1 = \frac{(y - \sum_{i=2}^k \beta_i x_i, x_1)}{(x_1, x_1)} \quad (1)$$

and by plugging in we get the new minimization problem

$$\text{Minimize } \| y^{(2)} - \sum_{i=2}^k \beta_i x_i^{(2)} \|^2 \quad (2)$$

where

$$y^{(2)} = P_{\{x_1\}^\perp} y, \quad x_i^{(2)} = P_{\{x_1\}^\perp} x_i, \quad i = 2, \dots, k. \quad (3)$$

3.1 Lemma:

Let $\text{Rank} (x_1, \dots, x_k)$ be equal m and let x_1, \dots, x_m be linearly independent. Then $x_i^{(2)}$, $i = 2, \dots, m$ are linearly independent and the $x_i^{(2)}$, $i > m$ are linear combinations of the $x_i^{(2)}$, $i = 2, \dots, m$.

Proof:

a) From $\sum_{i=2}^m \lambda_i x_i^{(2)} = P(\sum_{i=2}^m \lambda_i x_i) = 0$ – we write P for short instead of

$P_{\{x_1\}^\perp}$ – it follows that $\sum_{i=2}^m \lambda_i x_i \in \text{span}\{x_1\}$ and hence $\lambda_2 = \dots = \lambda_m = 0$ from the linear independence of x_1, \dots, x_m .

b) For $i > m$ we get $x_i^{(2)} = \sum_{j=2}^m \lambda_{ij} x_j^{(2)}$ if $x_i = \sum_{j=1}^m \lambda_{ij} x_j$. □

3.2 Theorem:

Let $\text{Rank} (x_1, \dots, x_k) = m$ and let, moreover, x_1, \dots, x_m be linearly independent. Furthermore, let q_1, \dots, q_m be the pairwise orthogonal vectors obtained from (x_1, \dots, x_m) by applying the Gram-Schmidt orthogonalization procedure. Then the Least Squares solutions $\hat{\beta}_1, \dots, \hat{\beta}_m$ are recursively given by

$$\hat{\beta}_m = \frac{(q_m, y - \sum_{i=m+1}^k \beta_i x_i)}{(q_m, q_m)} \quad (4)$$

$$\hat{\beta}_i = \frac{(q_i, y - \sum_{j=i+1}^m \hat{\beta}_j x_j - \sum_{j=m+1}^k \beta_j x_j)}{(q_i, q_i)} \quad (5)$$

$i = m - 1, m - 2, \dots, 1$. Here $\beta_{m+1}, \dots, \beta_k$ are completely arbitrary.

Moreover

$$y - \sum_{i=1}^m \hat{\beta}_i x_i - \sum_{j=m+1}^k \beta_j x_j$$

does not depend on $\beta_{m+1}, \dots, \beta_k$.

Proof: Mathematical induction on m . If $m = 1$, then

$$\hat{\beta}_1 = \frac{(x_1, y - \sum_{i=2}^k \beta_i x_i)}{(x_1, x_1)} \quad (6)$$

and

$$y - \hat{\beta}_1 x_1 - \sum_{i=2}^k \beta_i x_i = y^{(2)} - \sum_{i=2}^k \beta_i x_i^{(2)}. \quad (7)$$

But since $x_i \in \text{span}\{x_1\}$ it follows that $x_i^{(2)} = 0$, $i = 2, \dots, k$ and therefore

$$y - \hat{\beta}_1 x_1 - \sum_{i=2}^k \beta_i x_i = y^{(2)} \quad (8)$$

which does not depend on β_2, \dots, β_k .

We now arrive at the problem of minimizing

$$\| y^{(2)} - \sum_{i=2}^k \beta_i x_i^{(2)} \| . \quad (9)$$

By the induction assumption using that $x_2^{(2)}, \dots, x_k^{(2)}$ are linearly independent and $\text{Rank}(x_2^{(2)}, \dots, x_k^{(2)}) = m - 1$ we get that the solutions are as follows:

$$(q_m^{(2)}, q_m^{(2)}) \hat{\beta}_m = (q_m^{(2)}, y^{(2)} - \sum_{i=m+1}^k \beta_i x_i^{(2)}) \quad (10)$$

and

$$(q_i^{(2)}, q_i^{(2)}) \hat{\beta}_i = (q_i^{(2)}, y^{(2)} - \sum_{j=i+1}^m \hat{\beta}_j x_j - \sum_{j=m+1}^k \beta_j x_j) \quad (11)$$

$i = m - 1, \dots, 2$ Here $\beta_{m+1}, \dots, \beta_k$ are arbitrary numbers and $q_2^{(2)}, \dots, q_m^{(2)}$ are obtained by applying the Gram-Schmidt orthogonalization procedure to $x_2^{(2)}, \dots, x_m^{(2)}$. Moreover,

$y^{(2)} - \sum_{i=2}^m \hat{\beta}_i x_i - \sum_{i=m+1}^k \beta_i x_i$ does not depend on $\beta_{m+1}, \dots, \beta_k$. From this it follows that

$$y - \sum_{i=1}^m \hat{\beta}_i x_i - \sum_{j=m+1}^k \beta_j x_j = y^{(2)} - \sum_{i=2}^m \hat{\beta}_i x_i^{(2)} - \sum_{i=m+1}^k \beta_i x_i^{(2)} \quad (12)$$

as well does not depend on $\beta_{m+1}, \dots, \beta_k$.

We now prove by mathematical induction that $q_i^{(2)} = q_i$, $i = 2, \dots, m$. This is indeed correct for $i = 2$ since $x_2^{(2)} = q_2^{(2)} = x_2 - \frac{(x_1, x_2)}{(x_1, x_1)} x_1 = q_2$ and by using the induction assumption we get

$$q_i^{(2)} = x_2^{(2)} - \sum_{j=2}^{i-1} \frac{(x_i^{(2)}, q_j^{(2)})}{(q_j^{(2)}, q_j^{(2)})} q_j^{(2)} = x_i - \frac{(x_i, x_1)}{(x_1, x_1)} x_1 - \sum_{j=2}^{i-1} \frac{(x_i^{(2)}, q_j)}{(q_j, q_j)} q_j. \quad (13)$$

Since $(x_i^{(2)}, q_j) = (x_i, q_j)$ for $j \geq 2$, it follows indeed that $q_i^{(2)} = q_i$, $i = 2, \dots, m$.

From $(q_m, q_1) = 0$ for $i \geq 2$ we finally get

$$\hat{\beta}_m = \frac{(q_m, y^{(2)} - \sum_{i=m+1}^k \beta_i x_i^{(2)})}{(q_m, q_m)} = \frac{(q_m, y - \sum_{i=m+1}^k \beta_i x_i)}{(q_m, q_m)} \quad (14)$$

and for $i = m-1, \dots, 2$

$$\begin{aligned} \hat{\beta}_i &= \frac{(q_i, y^{(2)} - \sum_{j=i+1}^m \hat{\beta}_j x_j^{(2)} - \sum_{j=m+1}^k \beta_j x_j^{(2)})}{(q_i, q_i)} \\ &= \frac{(q_i, y - \sum_{j=i+1}^m \hat{\beta}_j x_j - \sum_{j=m+1}^k \beta_j x_j)}{(q_i, q_i)}. \end{aligned} \quad (15)$$

This is completed by

$$\hat{\beta}_1 = \frac{(x_1, y - \sum_{j=2}^m \hat{\beta}_j x_j - \sum_{j=m+1}^k \beta_j x_j)}{(x_1, x_1)} = \frac{(q_1, y - \sum_{j=2}^m \hat{\beta}_j x_j - \sum_{j=m+1}^k \beta_j x_j)}{(q_1, q_1)}. \quad (16)$$

□

3.3 References

- Achieser-Glasmann (1981), Theorie der linearen Operatoren im Hilbert-Raum, Akademie-Verlag Berlin.
- Drygas, H. (2008), QR-decomposition from the statistical point of view, Recent Advances in Linear Models and Related Areas, Essays in Honour of Helge Toutenburg, Shalab and Heumann (Eds) p. 293-311, Physica-Verlag, Springer, Heidelberg.
- Drygas, H. (2009), Statistical Analysis of Diabetes Mellitus. *Discussiones Mathematicae, Probability and Statistics*, 29 (2009), p. 69 - 90, University of Zielona Góra, Poland.
- Gauss, C. F. "Theoria combinationis Observationum Erroribus Minimis Obnoxae" *Comm. recent. Soc. Gött.* 5 (1819 - 1822) oder *Werke*, Bd IV, Leipzig, 1879, p. 3 - 53. A German translation under the title "Abhandlungen zur Methode der kleinsten Quadrate" was edited by Börsch and Simon in 1987, Druck and Verlag von P. Stankiewicz Buchdruckerei, Berlin.
- Linnik, J. W. (1961), *Methode der kleinsten Quadrate in Moderner Darstellung*, VEB Deutscher Verlag der Wissenschaften, Berlin.
- Rao, C. R. and Mitra, S. K. (1971), *Generalized Inverse of Matrices and its Applications*, Wiley, New York - London - Sydney - Toronto.
- Richter, H., Mammitzsch, V. (1973), *Methode der kleinsten Quadrate mit Übungen und Aufgaben*.