

3 Multiple linear regression

So far we have dealt with a single explanatory variable. We now generalize the ideas to multiple explanatory variables. Thus, the data is now

$$(x_{i1}, \dots, x_{ip}, y_i), \quad i = 1, \dots, n.$$

The explanatory variable for the i th observation is a p -dimensional vector, $\mathbf{x}_i = (x_{i1}, \dots, x_{ip})$, and we assume throughout $p + 1 \leq n$ (we will see soon why the case $p + 1 > n$ is problematic).

In the multiple explanatory variables case it is generally challenging to visualize the data points, because we quickly run out of dimensions if we attempt to generate a scatterplot (the case $p = 2$ is still doable because we can draw a 3-dimensional plot, showing X_1, X_2 on the XY -plane and the response Y on the Z axis). However, the mathematical concepts from the simple regression case can still be extended to the general- p case. Thus, we may still try to predict y by a *linear* function of x_1, \dots, x_p , i.e., the analog of (1) will be

$$y = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2 + \dots + \hat{\beta}_p x_p = \sum_{j=0}^p \hat{\beta}_j x_j,$$

where for each observation we prepend $x_{i0} \equiv 1$ to the vector of explanatory variables, that is, we define $\mathbf{x}_i = (1, x_{i1}, \dots, x_{ip})$ instead of $\mathbf{x}_i = (x_{i1}, \dots, x_{ip})$.

For any method for fitting such a linear function (that is, for calculating $\hat{\beta}_1, \dots, \hat{\beta}_p$ as a function of the observed data), the *fitted values* are

$$\hat{y}_i = \sum_{j=0}^p \hat{\beta}_j x_{ij},$$

and the *residuals* are defined exactly as before,

$$e_i := y_i - \hat{y}_i$$

We can generalize also the least squares method, by seeking the set of coefficients $\hat{\beta}_1, \dots, \hat{\beta}_p$ which, as before, minimize the sum of squared errors. Formally, let $\mathbf{b} = (b_0, \dots, b_p)^\top \in \mathbb{R}^{p+1}$, and define

$$Q(\mathbf{b}) := \sum_{i=1}^n \left(y_i - \sum_{j=0}^p b_j x_{ij} \right)^2. \quad (8)$$

This is a differentiable function (now of $p + 1$ variables), and we can take partial derivatives and set them to zero. We have

$$\frac{\partial}{\partial b_r} \sum_{j=0}^p x_{ij} b_j = x_{ir},$$

and, therefore,

$$\frac{\partial Q}{\partial b_r} = -2 \sum_{i=1}^n \left(y_i - \sum_{j=0}^p x_{ij} b_j \right) \left(\frac{\partial}{\partial b_r} \sum_{j=0}^p x_{ij} b_j \right) = -2 \sum_{i=1}^n \left(y_i - \sum_{j=0}^p x_{ij} b_j \right) x_{ir}.$$

Setting this to zero for $r = 0, \dots, p$ yields the so-called *Normal equations*,

$$\sum_{j=0}^p \left(\sum_{i=1}^n x_{ir} x_{ij} \right) b_j = \sum_{i=1}^n x_{ir} y_i, \quad r = 0, \dots, p. \quad (9)$$

The solution to this set of $p + 1$ equations (in $p + 1$ variables), assuming it exists and is unique (we will see conditions for this to be the case), gives the least squares coefficients $\hat{\beta}_1, \dots, \hat{\beta}_p$.

The system of equations (9) can be written equivalently in *matrix* form. Thus, define

$$\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} \in \mathbb{R}^n, \quad \mathbf{X} := \begin{bmatrix} 1 & x_{11} & x_{12} & \cdots & x_{1p} \\ 1 & x_{21} & x_{22} & \cdots & x_{2p} \\ \vdots & \vdots & \vdots & \cdots & \vdots \\ 1 & x_{n1} & x_{n2} & \cdots & x_{np} \end{bmatrix} \in \mathbb{R}^{n \times (p+1)}.$$

The $n \times (p + 1)$ matrix \mathbf{X} is called the *design matrix*, or simply the X -matrix. For $i = 1, \dots, n$, the i th row of \mathbf{X} is the $(p + 1)$ -dimensional row vector \mathbf{x}_i^\top , the feature vector for the i th sample point (transposed). For $j = 0, \dots, p$, the j th column of \mathbf{X} is the n -dimensional column vector $\mathbf{X}_j := (x_{1j}, \dots, x_{nj})^\top \in \mathbb{R}^n$, the vector of values of the j th feature for each of the n observations. Note that $\mathbf{X}_0 = (1, \dots, 1)^\top =: \mathbf{1}_n$.

With this notation,

$$\sum_{i=1}^n x_{ir} x_{ij} = (\mathbf{X}^\top \mathbf{X})_{rj}$$

and

$$\sum_{i=1}^n x_{ir} y_i = (\mathbf{X}^\top \mathbf{y})_r,$$

so we can write (2.2) as

$$\sum_{j=0}^p (\mathbf{X}^\top \mathbf{X})_{rj} b_j = (\mathbf{X}^\top \mathbf{y})_r, \quad r = 0, \dots, p$$

or, equivalently,

$$\mathbf{X}^\top \mathbf{X} \mathbf{b} = \mathbf{X}^\top \mathbf{y}$$

Now, $\mathbf{X}^\top \mathbf{X}$ is a $(p + 1) \times (p + 1)$ matrix, and it is invertible if and only if the columns of \mathbf{X} are linearly independent (prove this!). If $p + 1 > n$ the columns of \mathbf{X} are necessarily linearly dependent (because column rank $\leq \min(n, p + 1) = n < p + 1$). If $p + 1 \leq n$, the columns of \mathbf{X} are linearly independent if no feature vector \mathbf{X}_j is a combination of the others, which is a reasonable assumption. If the columns of \mathbf{X} are linearly dependent, $\mathbf{X}^\top \mathbf{X}$ is singular and (2.3) has infinitely many solutions, i.e., the LS coefficients $\hat{\beta}$ are not unique (every solution for (2.3) is then said to be a LS solution). If the columns of \mathbf{X} are linearly independent, which we will generally assume from now on, then (2.3) has a unique solution given by

$$\hat{\beta} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}. \quad (10)$$

Geometric interpretation of the LS solution. There is in fact a more direct (and more intuitive) way to derive the LS solution using *geometric* interpretation. We first recall some basic definitions and facts from linear algebra.

1. For a $n \times m$ matrix \mathbf{A} , the *image* is the linear subspace of \mathbb{R}^n given by

$$\text{Im}(\mathbf{A}) := \{\mathbf{A}\mathbf{v} : \mathbf{v} \in \mathbb{R}^m\}$$

Recall that, if \mathbf{A}_j denotes the j -th column of \mathbf{A} , then $\mathbf{A}\mathbf{v} = \sum_j \mathbf{A}_j v_j$, so we have

$$\text{Im}(\mathbf{A}) = \text{sp}(\mathbf{A}_1, \dots, \mathbf{A}_m) \subseteq \mathbb{R}^n,$$

i.e., the image is the linear subspace (of \mathbb{R}^n) *spanned* by the columns $\mathbf{A}_1, \dots, \mathbf{A}_m$.

2. The (standard) *inner product* of vectors $\mathbf{u}, \mathbf{v} \in \mathbb{R}^n$ is $\mathbf{u}^\top \mathbf{v} = \mathbf{v}^\top \mathbf{u} = \sum_{i=1}^n u_i v_i$.
3. The *norm* of a vector $\mathbf{v} \in \mathbb{R}^n$ is $\|\mathbf{v}\| = (\mathbf{v}^\top \mathbf{v})^{1/2} = (\sum_{i=1}^n v_i^2)^{1/2}$. The *Euclidean distance* between two vectors $\mathbf{u}, \mathbf{v} \in \mathbb{R}^n$ is $\|\mathbf{u} - \mathbf{v}\|$. In $\mathbb{R}^2, \mathbb{R}^3$ this coincides with the usual geometric notion of a distance between two points.
4. For $\mathbf{u}, \mathbf{v} \in \mathbb{R}^n$, we say that \mathbf{v} is *orthogonal* to \mathbf{u} , and denote $\mathbf{u} \perp \mathbf{v}$, if $\mathbf{u}^\top \mathbf{v} = 0$. In $\mathbb{R}^2, \mathbb{R}^3$ this coincides with the usual geometric notion of perpendicularity.
5. The *orthogonal complement* of a subspace $M \subseteq \mathbb{R}^n$ is the subspace

$$M^\perp := \{\mathbf{v} : \mathbf{v}^\top \mathbf{u} = 0 \quad \forall \mathbf{u} \in M\}$$

(Exercise: verify that M^\perp is indeed a linear space.)

6. *Pythagorean theorem*: If $\mathbf{u} \perp \mathbf{v}$, then $\|\mathbf{v} + \mathbf{u}\|^2 = \|\mathbf{v}\|^2 + \|\mathbf{u}\|^2$.

We now derive the LS solution (10) from an alternative, geometric viewpoint. For any $\mathbf{b} = (b_0, \dots, b_p)^\top \in \mathbb{R}^{p+1}$ we can write the objective function (8) in vector form as

$$Q(\mathbf{b}) = \|\mathbf{y} - \mathbf{X}\mathbf{b}\|^2.$$

Thus, if $\hat{\beta}$ is the minimizer of $Q(\mathbf{b})$ over all $\mathbf{b} \in \mathbb{R}^{p+1}$, then $\hat{\mathbf{y}} = \mathbf{X}\hat{\beta}$ minimizes the squared norm (equivalently, it minimizes the Euclidean distance) between \mathbf{y} and \mathbf{z} among all vectors $\mathbf{z} \in \text{Im}(\mathbf{X})$. Then $\hat{\mathbf{y}} = \mathbf{X}\hat{\beta}$ must be the *projection* of \mathbf{y} onto $\text{colsp}(\mathbf{X})$ (this claim requires a proof, but it's intuitive geometrically). Now, a necessary and sufficient condition for $\hat{\mathbf{y}} = \mathbf{X}\hat{\beta}$ to be the projection of \mathbf{y} onto $\text{colsp}(\mathbf{X})$ is

$$\mathbf{X}^\top (\mathbf{y} - \mathbf{X}\hat{\beta}) = \mathbf{0}, \tag{11}$$

which simply requires that the dot product between the residual $\mathbf{e} = \mathbf{y} - \hat{\mathbf{y}}$ and each column of \mathbf{X} is zero. Rearranging (11), we get

$$\mathbf{X}^\top \mathbf{y} = \mathbf{X}^\top \mathbf{X} \hat{\beta}$$

which, assuming again that $\mathbf{X}^\top \mathbf{X}$ is invertible, yields (10). The vector

$$\hat{\mathbf{y}} = \mathbf{X}\hat{\beta} = \mathbf{X} \left(\mathbf{X}^\top \mathbf{X} \right)^{-1} \mathbf{X}^\top \mathbf{y}$$

is called the orthogonal projection of \mathbf{y} onto $\text{Im}(\mathbf{X})$, and the matrix

$$\mathbf{P}_\mathbf{X} := \mathbf{X} \left(\mathbf{X}^\top \mathbf{X} \right)^{-1} \mathbf{X}^\top$$

i.e., the matrix such that $\hat{\mathbf{y}} = \mathbf{P}_\mathbf{X} \mathbf{y}$, is called the projection matrix of \mathbf{X} (this is the matrix projecting any vector onto the linear space spanned by the columns of \mathbf{X}).

Remark. We are assuming, as before, that $\mathbf{X}^\top \mathbf{X}$ is invertible, although the projection matrix of \mathbf{X} is well-defined also when the columns of \mathbf{X} are linearly dependent; for example, this can be achieved by choosing a basis for $\text{Im}(\mathbf{X})$ and writing (2.5) using the "thinner" matrix instead of \mathbf{X}).

Linear algebra interlude: projection matrices and related results.

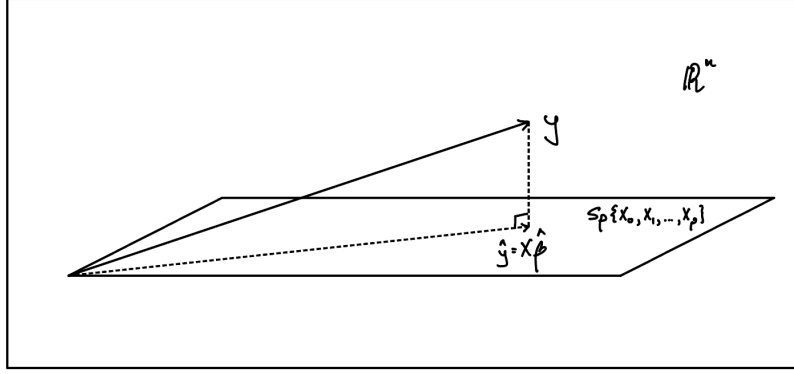


Figure 4: Geometric interpretation of the Least squares estimator

Proposition 4. Let X be an $n \times m$ matrix and assume that it has linearly independent columns (i.e., full column rank; remember that this implies $m \leq n$). Then the projection matrix P_X has the following properties.

1. P_X is symmetric
2. P_X is idempotent, $P_X^2 = P_X$
3. $P_X X = X$
4. $X^\top (I - P_X) = 0 \in \mathbb{R}^{m \times n}$
5. $P_X v \in \text{Im}(X)$ for all $v \in \mathbb{R}^n$
6. If $m = n$ and X is invertible, then $P_X = I$
7. $(I - P_X) v \in \text{Im}(X)^\perp$ for all $v \in \mathbb{R}^n$
8. If $w \in \text{Im}(X)$, then $P_X w = w$
9. If $w \in \text{Im}(X)^\perp$, then $P_X w = 0$
10. If Z is another $n \times m$ matrix s.t. $\text{Im}(Z) = \text{Im}(X)$, then $P_Z = P_X$. This means that P_X depends on X only through the span of its columns. Hence, for an arbitrary linear space M , we can define the projection matrix P_M onto M (an explicit form for P_M can be obtained by taking any basis of M and stacking its elements as columns in a matrix X , then forming $P_X := X (X^\top X)^{-1} X^\top$)
11. If L and M are two subspaces with $L \subseteq M$, then $P_M P_L = P_L P_M = P_L$.

Proof.

$$1. P_X^\top = \left[X (X^\top X)^{-1} X^\top \right]^\top = X (X^\top X)^{-1} X^\top \text{ where we used the fact } \left[(X^\top X)^{-1} \right]^\top = \left[(X^\top X)^\top \right]^{-1} = (X^\top X)^{-1}.$$

2. $P_X^2 = \left[X (X^\top X)^{-1} X^\top \right] \left[X (X^\top X)^{-1} X^\top \right] = X (X^\top X)^{-1} X^\top$
3. $P_X X = X (X^\top X)^{-1} X^\top X = X$
4. $X^\top (I - P_X) = [(I - P_X) X]^\top = [X - P_X X]^\top = [X - X]^\top = \mathbf{0}^\top \in \mathbb{R}^{n \times m} = \mathbf{0} \in \mathbb{R}^{m \times n}$,
where we used fact #3
5. $P_X v$ for all $P_X v = X (X^\top X)^{-1} X^\top v = X \left[(X^\top X)^{-1} X^\top v \right] \in \text{Im}(X)$
6. If P_X is (square and) invertible, $\left[X (X^\top X)^{-1} X^\top \right]^{-1} = \left[X^\top \right]^{-1} (X^\top X) X^{-1} = I_n$.
7. $(Xu)^\top (I - P_X) v = u^\top \underbrace{X^\top (I - P_X)}_{\mathbf{0} \in \mathbb{R}^{m \times n}} v = 0$
8. $P_X \underbrace{(Xu)}_w = (P_X X) u = \underbrace{Xu}_w$
9. $P_X w = X (X^\top X)^{-1} \underbrace{X^\top w}_{=0} = \mathbf{0}$
10. Denote X_j, Z_j for the j th columns of X, Z , respectively. Then $Z_j = X h_j$ for some $h_j \in \mathbb{R}^m$ because $Z_j \in \text{Im}(Z) = \text{Im}(X)$. Putting $H := [h_1 h_2 \cdots h_m] \in \mathbb{R}^{m \times m}$, we have $Z = XH$. Further, if $w \in \mathbb{R}^m$ is s.t. $Hw = \mathbf{0}$, then $Zw = XHw = \mathbf{0}$, which implies $w = \mathbf{0}$ because Z was assumed to have full column rank. Hence, H is invertible. Then

$$\begin{aligned} P_Z &= Z (Z^\top Z)^{-1} Z^\top = XH [(XH)^\top XH]^{-1} (XH)^\top = \\ &= XHH^{-1}X^{-1} [(XH)^\top]^{-1} (XH)^\top = X \left[(X^\top X)^{-1} \right] X^\top = P_X \end{aligned}$$

11. Let v be any vector. Then $P_M \underbrace{P_L v}_v \stackrel{s}{=} P_L v$, implying $P_M P_L = P_L$. Transposing both sides and $\in M$ using the fact that P_M, P_L are both symmetric, we obtain also $P_L P_M = P_L$.

□

Additionally, we have the following results related to projection matrices.

Proposition 5. Let M be a subspace of \mathbb{R}^n with $\dim(M) = m \leq n$. Then any vector $v \in \mathbb{R}^n$ can be uniquely represented as $v = w + z$ where $w \in M$ and $z \in M^\perp$, the orthogonal complement of M . Moreover, in this representation, $w = P_M v$, the projection of v onto M , and satisfies $w = \arg \min_{u \in M} \|v - u\|^2$.

Proof. Taking $w = P_M v, z = v - P_M v$, we have $w \in M, z \in M^\perp$ by the properties above, and $v = w + z$. We show that this representation is unique. Thus, suppose $v = w_1 + z_1, v = w_2 + z_2$ for $w_1, w_2 \in M, z_1, z_2 \in M^\perp$. Then $\mathbf{0} = (w_1 + z_1) - (w_2 + z_2) = (w_1 - w_2) + (z_1 - z_2)$, and, furthermore, $(w_1 - w_2) \perp (z_1 - z_2)$ because $(w_1 - w_2) \in M, (z_1 - z_2) \in M^\perp$ (remember that each of M, M^\perp is a subspace so closed under addition/subtraction). Therefore,

$$0 = \|\mathbf{0}\|^2 = \|(w_1 - w_2) + (z_1 - z_2)\|^2 = \|w_1 - w_2\|^2 + \|z_1 - z_2\|^2 \Rightarrow (w_1 - w_2) = \mathbf{0}, (z_1 - z_2) = \mathbf{0}.$$

It remains to show that \mathbf{w} minimizes the Euclidean distance from \mathbf{v} to M . Take any $\mathbf{u} \in M$. Then

$$\|\mathbf{v} - \mathbf{u}\|^2 = \|(\mathbf{w} + \mathbf{z}) - \mathbf{u}\|^2 = \|\mathbf{w} - \mathbf{u} + \mathbf{z}\|^2 = \|\mathbf{w} - \mathbf{u}\|^2 + \|\mathbf{z}\|^2 \geq \|\mathbf{z}\|^2,$$

and on the other hand, for $\mathbf{u} = \mathbf{w} = \mathbf{P}_M \mathbf{v}$ we have $\|\mathbf{v} - \mathbf{u}\|^2 = \|\mathbf{z}\|^2$. \square

Proposition 6. *We have*

1. $\mathbf{I} - \mathbf{P}_X = P_{\text{Im}(\mathbf{X})^\perp}$
2. if L and M are two subspaces of \mathbb{R}^n with $L \subseteq M$, then $\mathbf{P}_M - \mathbf{P}_L = \mathbf{P}_{M \cap L^\perp}$

Proof. (a) this part follows from the uniqueness of the representation in Proposition 1. (b) this also follows from the uniqueness of the representation in Proposition 1, taking the original space to be M (I.e., the space to which \mathbf{v} belongs) and the subspace (what's denoted M in Proposition 1) to be L . \square

Proposition 7. *Let \mathbf{Q} be an $n \times n$ matrix of rank $m \leq n$ which is symmetric and idempotent, $\mathbf{Q}^\top = \mathbf{Q}$, $\mathbf{Q}^2 = \mathbf{Q}$. Then $\mathbf{Q} = \mathbf{P}_M$ where $M := \text{Im}(\mathbf{Q})$.*

Proof. Exercise. \square

Diagonalizability and positive-semidefiniteness of a projection matrix. First, we recall some facts:

1. A matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ is diagonalizable (over \mathbb{R}) if there is an invertible matrix $\mathbf{P} \in \mathbb{R}^{n \times n}$ and a diagonal matrix \mathbf{D} such that $\mathbf{A} = \mathbf{P}\mathbf{D}\mathbf{P}^{-1}$.

Remark: \mathbf{P}^{-1} is the transition matrix from the standard basis of \mathbb{R}^n into the basis in the columns of \mathbf{P} .

2. If $\mathbf{A} \in \mathbb{R}^{n \times n}$ is symmetric, then there is an orthogonal matrix $\mathbf{U} \in \mathbb{R}^{n \times n}$ and a diagonal matrix \mathbf{D} such that $\mathbf{A} = \mathbf{U}\mathbf{D}\mathbf{U}^{-1} = \mathbf{U}\mathbf{D}\mathbf{U}^\top$.

Remark: (i) a square matrix \mathbf{U} is orthogonal if $\mathbf{U}^\top \mathbf{U} = \mathbf{I}$ (this implies $\mathbf{U}^\top = \mathbf{U}^{-1}$, because for any two square matrices \mathbf{A}, \mathbf{B} , we have $\mathbf{AB} = \mathbf{I} \Rightarrow \mathbf{BA} = \mathbf{I}$; prove this!). (ii) consistent with the above, you should have seen a proof in your linear algebra course that if \mathbf{A} is a (real) symmetric matrix, then it is diagonalizable, and all of its eigenvalues are real.

3. A symmetric matrix \mathbf{A} is *positive semidefinite* if all of its eigenvalues are nonnegative. It is called *positive definite* if all of its eigenvalues are positive.

Remark: this is equivalent to saying that the *quadratic form* given by $\mathbf{x}^\top \mathbf{A} \mathbf{x}$ is *nonnegative* for all vectors \mathbf{x} (positive semidefinite) or that $\mathbf{x}^\top \mathbf{A} \mathbf{x}$ is *positive* (positive definite).

4. Let \mathbf{A} be a positive semidefinite matrix. Then:

- (i) there is \mathbf{B} square such that $\mathbf{B}^2 = \mathbf{A}$, and in fact this representation is unique: take $\mathbf{B} = \mathbf{U}\mathbf{D}^{1/2}\mathbf{U}^\top$, where $\mathbf{D}^{1/2} := \text{diag}(d_1, \dots, d_n)$, so $\mathbf{B}^2 = \mathbf{U}\mathbf{D}^{1/2}\mathbf{U}^\top \mathbf{U}\mathbf{D}^{1/2}\mathbf{U}^\top = \mathbf{A}$ (this \mathbf{B} is often called the square root of \mathbf{A}).
- (ii) there is \mathbf{B} square such that $\mathbf{B}\mathbf{B}^\top = \mathbf{A}$: take $\mathbf{B} = \mathbf{U}\mathbf{D}^{1/2}$, then $\mathbf{B}^2 = \mathbf{U}\mathbf{D}^{1/2}(\mathbf{U}\mathbf{D}^{1/2})^\top = \mathbf{A}$

Now suppose that $\mathbf{Q} \in \mathbb{R}^{n \times n}$ is a projection matrix onto a subspace $M \subseteq \mathbb{R}^n$ of dimension $\dim(M) = m$, then \mathbf{Q} is positive semidefinite, and there is a representation

$$\mathbf{Q} = \mathbf{U}\mathbf{D}\mathbf{U}^\top, \quad \mathbf{D} = \text{diag}(\underbrace{1, \dots, 1}_m, \underbrace{0, \dots, 0}_{n-m}),$$

where \mathbf{U} is orthogonal, such that the first m columns of \mathbf{U} are an orthonormal basis of M , and the last $n - m$ columns of \mathbf{U} are an orthonormal basis of M^\perp .

Remark: consistent with the previous item, note that $\mathbf{Q} = \mathbf{B}\mathbf{B}^\top$ for $\mathbf{B} = \mathbf{U}\mathbf{D}^{1/2} = \tilde{\mathbf{U}} := [\mathbf{U}_1, \dots, \mathbf{U}_m]$.