# Chapter 17: The Least Squares Procedure

## 17.1 Projection Approach

We observe the following data:

$$y = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix} \in \mathbb{R}^n; \quad X = \begin{pmatrix} x_{11} & \cdots & x_{1K} \\ \vdots & \ddots & \vdots \\ x_{n1} & \cdots & x_{nK} \end{pmatrix} = (x_1, \ldots, x_K).$$

We would like to find $\beta$ such that $y = X\beta$, but when $n > K$ this is not possible. So we do the next best thing, which is to find the best linear (in $X$) approximation to $y$.

### Definition 17.1.

The **(Ordinary) Least Squares (OLS)** procedure chooses $\hat{\beta}$ to minimize the quadratic form

$$S(\beta) = (y - X\beta)^T (y - X\beta) = \|y - X\beta\|^2$$

with respect to $\beta \in \mathbb{R}^K$.

The data, $y, x_1, \ldots, x_K$, can all be viewed as elements of the vector space $\mathbb{R}^n$. Define the column span of $X$,

$$C(X) = \{\alpha_1 x_1 + \cdots + \alpha_K x_K\} = \{X\alpha : \alpha \in \mathbb{R}^K\} \subset \mathbb{R}^n.$$

Then, $C(X)$ is a linear subspace of $\mathbb{R}^n$. The projection theorem says that there is a unique solution to the minimization problem, call it $\hat{y}$, which is characterized by the fact that $y - \hat{y} = \hat{\varepsilon}$ is orthogonal to the space $C(X)$. That is, we can write uniquely

$$y = \hat{y} + \hat{\varepsilon},$$

where (the "fitted value") $\hat{y} \in C(X)$ and (the "residual") $\hat{\varepsilon} \in C(X)^\perp$. This means that

$$\hat{y} = X\hat{\beta}, \text{ some } \hat{\beta} \in \mathbb{R}^K \text{ and } X^T \hat{\varepsilon} = 0.$$

**Intuition:** Imagine $y$ is a point in a high-dimensional space, and $C(X)$ is a subspace (like a plane or a line) within that space. The OLS method finds the point $\hat{y}$ in the subspace $C(X)$ that is closest to $y$. The vector connecting $\hat{y}$ and $y$ is $\hat{\varepsilon}$. The vector $\hat{\varepsilon}$ represents the part of $y$ that cannot be explained by the linear combination of the columns of $X$.

The orthogonality conditions can be written out explicitly:

$$\sum_{i=1}^{n} x_{1i}\hat{\varepsilon}_i = 0$$

$$\sum_{i=1}^{n} x_{2i}\hat{\varepsilon}_i = 0$$

$$\vdots$$

$$\sum_{i=1}^{n} x_{Ki}\hat{\varepsilon}_i = 0.$$

Note that if, as usual, $x_{1i} = 1$, then we have $\sum_{i=1}^{n}\hat{\varepsilon}_i = 0$.

The vector $\hat{y}$ is interpreted as the best linear fit to the vector $y$, which establishes a connection with the regression material of Chapter 7, although we have not yet specified any random structure generating the data; the results are data specific.

We may write $\hat{\varepsilon} = y - X\hat{\beta}$ so that the orthogonality conditions may be rewritten as

$$X^T(y - X\hat{\beta}) = 0. \tag{17.1}$$

We can rewrite this equation as the so-called **normal equations**

$$X^T X\hat{\beta} = X^T y \tag{17.2}$$

$$X^T X = \begin{pmatrix} \sum_{i=1}^{n} x_{1i}^2 & \cdots & \sum_{i=1}^{n} x_{1i}x_{Ki} \\ \vdots & \ddots & \vdots \\ \sum_{i=1}^{n} x_{1i}x_{Ki} & \cdots & \sum_{i=1}^{n} x_{Ki}^2 \end{pmatrix}; \quad X^T y = \begin{pmatrix} \sum_{i=1}^{n} x_{1i}y_i \\ \vdots \\ \sum_{i=1}^{n} x_{Ki}y_i \end{pmatrix}.$$

There always exists some $\hat{\beta}$ such that $S(\hat{\beta})$ is minimal, but it may not be unique. When $\text{rank}(X) = K$, then $X^T X$ is of full rank and hence invertible. Then $\hat{\beta}$ is uniquely defined for any $y$, i.e.,

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

How do we find $\hat{y}$ and $\hat{\varepsilon}$?

## Definition 17.2.

When $X$ is of full rank define the $n \times n$ **Projector matrices**

$$P_X = X(X^T X)^{-1} X^T; \quad M_X = I - X(X^T X)^{-1} X^T$$

which project onto $C(X)$ and onto the orthogonal complement of $C(X)$, denoted $C(X)^{\perp}$. For any $y$, we can uniquely write

$$y = \hat{y} + \hat{\varepsilon} = P_X y + M_X y.$$

We have $P_X X = X$ and $M_X X = 0$. The matrices $P_X$ and $M_X$ are symmetric and idempotent, i.e.,

$$P_X = P_X^T; \quad P_X^2 = P_X.$$

**Intuition:** $P_X$ takes any vector $y$ and projects it onto the subspace spanned by the columns of $X$. The resulting

vector is the fitted value $\hat{y}$. $M_X$ is the complement of $P_X$. It projects $y$ onto the orthogonal complement of $C(X)$. The resulting vector is the residual $\hat{\varepsilon}$.

After applying $P_X$ once you are ready in $C(X)$. This means that they have eigenvalues either 0 or 1 using the eigendecomposition. Let $P_X = U\Lambda U^T$ for orthonormal $U$ and diagonal $\Lambda$. Then

$$(P_X)^2 = (U\Lambda U^T)(U\Lambda U^T) = (U\Lambda^2 U^T)$$

so if $P_X^2 = P_X$, then $\Lambda^2 = \Lambda$, and eigenvalues have to be in $\{0, 1\}$. $P_X$ is of rank $K$ and has $K$ eigenvalues equal to 1 and $n - K$ eigenvalues equal to 0, with the reverse for $M_X$.

## Theorem 17.1.

The space $C(X)$ is invariant to nonsingular linear transforms. That is for $A_{K \times K}$ with $\det A \neq 0$, we have

$$C(XA) = C(X).$$

**Proof.** Let $v \in C(X)$. Then there exists an $\alpha \in \mathbb{R}^K$ such that $v = X\alpha$. Therefore,

$$v = XAA^{-1}\alpha = XA\gamma,$$

where $\gamma = A^{-1}\alpha \in \mathbb{R}^K$. Likewise for any $v = XA\gamma$, we can write $v = X\alpha$ for $\alpha = A^{-1}\gamma$. That is, $X \in C(XA)$ and so $C(X) = C(XA)$. $\square$

Since $C(X)$ is invariant to linear transformations, so are $\hat{y}$ and $\hat{\varepsilon}$ (but not $\hat{\beta}$). For example, rescaling of the components of $X$ does not affect the values of $\hat{y}$ and $\hat{\varepsilon}$.

$y$ on $(x_1, x_2, x_3)$ (17.3)

$y$ on $(x_1 + x_2, 2x_2 - x_3, 3x_1 - 2x_2 + 5x_3)$ (17.4)

in which case the transformation is

$$A = \begin{pmatrix} 1 & 0 & 3 \\ 1 & 2 & -2 \\ 0 & -1 & 5 \end{pmatrix}$$

which is of full rank. Therefore, (17.3) and (17.4) yield the same $\hat{y}$, $\hat{\varepsilon}$.

## Example 17.1.

**Dummy variables.** Suppose that $x_{ji} = 1$ if $j \in I_l$ and $x_{ji} = 0$ if $j \notin I_l$, where $\{I_l\}$ forms a partition of $\{1, \ldots, n\}$. Specifically,

$$I_l \cap I_{l'} = \emptyset, l \neq l'; \quad \bigcup_{l=1}^{L} I_l = \{1, \ldots, n\}.$$

For example, day of the week dummies. In this case $L = 5$

$$X = \begin{pmatrix} 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ \vdots & 0 & 1 & 0 & 0 \\ \vdots & 0 & 1 & 0 \\ \vdots & \vdots & 0 & 1 \\ \vdots & \vdots & \vdots & \vdots \end{pmatrix}.$$

We have a diagonal hat matrix

$$X^T X = \begin{pmatrix} n_1 & 0 & 0 & 0 & 0 \\ 0 & n_2 & 0 & 0 & 0 \\ 0 & 0 & n_3 & 0 & 0 \\ 0 & 0 & 0 & n_4 & 0 \\ 0 & 0 & 0 & 0 & n_5 \end{pmatrix}.$$

**Intuition:**

Dummy variables are used to represent categorical data. In the example, we are looking at days of the week, so we have $L = 5$ categories. Each column of $X$ represent one category. The $X^T X$ matrix is diagonal. The elements on the diagonal represent the number of observations in each category.

Emphasizing $C(X)$ rather than $X$ itself is called the **coordinate free approach**. Some aspects of model/estimate are properties of $C(X)$, the choice of coordinates is irrelevant.

When $X$ is not of full rank, the space $C(X)$ is still well defined, as is the projection from $y$ onto $C(X)$. The fitted value $\hat{y}$ and residual $\hat{\varepsilon}$ are uniquely defined in this case, but there is no unique coefficient vector $\hat{\beta}$. This case is often called **multicollinearity** in econometrics, although it also arises in big data where the number of covariates is large relative to the number of observations. Suppose that $X$ is of deficient rank $J < K$. The singular value decomposition of $X$ is

$$X = USV^T,$$

where $S = \mathrm{diag}\{s_1, s_2, \ldots, s_J, 0, \ldots, 0\}|_{0_{n-K \times K}}$. Then we may write for $\lambda_j > 0$

$$X^T X = U \Lambda U^T, \quad U^T X^T X U = \Lambda$$

**Intuition:** When $X$ is not full rank, it means that at least one of the columns of $X$ can be written as linear combination of the other columns. This implies that we have redundant information in our regressors.

## 17.2 Partitioned Regression

We next consider an important application of the projection idea. Partition

$$X = (X_{1n \times K_1}, X_{2n \times K_2}), \quad K_1 + K_2 = K,$$

and suppose we are interested only in $\hat{\beta}_1$ in the **long regression**.

We are going to show a formula for $\hat{\beta}_1$ that does not involve computing all of $\hat{\beta}$ and reading off the subvector $\hat{\beta}_1$. This will be called the **Frisch-Waugh-Lovell Theorem**. A key property of projection is given below.

## Theorem 17.2.

Suppose that $X_1$ and $X_2$ are orthogonal, i.e., $X_1^T X_2 = 0$. Then

$$P_X = P_{X_1} + P_{X_2}.$$

This can be verified algebraically, but also should be obvious geometrically.

We return to the general case. We have

$$\hat{y} = X\hat{\beta} = P_X y = P_{X_1} y + P_{X_2} y = X_1 \hat{\beta}_1 + X_2 \hat{\beta}_2.$$

This just says that if $X_1$ and $X_2$ were orthogonal, then we could get $\hat{\beta}_1$ by regressing $y$ on $X_1$ only, and $\hat{\beta}_2$ by regressing $y$ on $X_2$ only.

$$\begin{pmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \end{pmatrix} = \begin{pmatrix} X_1^T X_1 & 0 \\ 0 & X_2^T X_2 \end{pmatrix}^{-1} \begin{pmatrix} X_1^T y \\ X_2^T y \end{pmatrix} = \begin{pmatrix} (X_1^T X_1)^{-1} X_1^T y \\ (X_2^T X_2)^{-1} X_2^T y \end{pmatrix}.$$

It is rare that design matrices $X_1$ and $X_2$ are orthogonal, but we can construct equivalent regressors that *are* orthogonal. Suppose we have general $X_1$ and $X_2$, whose dimensions satisfy $K_1 + K_2 = K$. We make the following observations:

1. $(X_1, X_2)$ and $(M_2 X_1, X_2)$ span the same space. This follows because $X_1 = M_2 X_1 + P_2 X_1$, where $C(P_2 X_1) \subset C(X_2)$. Therefore, $C(M_2 X_1, X_2) = C(X_1, X_2)$.

2. $M_2 X_1$ and $X_2$ are orthogonal.

This says that if we regress $y$ on $(X_1, X_2)$ or $y$ on $(M_2 X_1, X_2)$ we get the same $\hat{y}$ and $\hat{\varepsilon}$, and that if we wanted the coefficients on $M_2 X_1$ from the second regression we could in fact just regress $y$ on $M_2 X_1$ only.

What are the coefficients on $M_2 X_1$? Recall that

$$\begin{aligned}
y &= X_1 \hat{\beta}_1 + X_2 \hat{\beta}_2 \\
&= (M_2 + P_2) X_1 \hat{\beta}_1 + X_2 \hat{\beta}_2 \\
&= M_2 X_1 \hat{\beta}_1 + X_2 [\hat{\beta}_2 + (X_2^T X_2)^{-1} X_2^T X_1 \hat{\beta}_1] \\
&= M_2 X_1 \hat{\beta}_1 + X_2 \hat{C},
\end{aligned}$$

where $\hat{C} = \hat{\beta}_2 + (X_2^T X_2)^{-1} X_2^T X_1 \hat{\beta}_1$. So the coefficient on $M_2 X_1$ is the original $\hat{\beta}_1$, while that on $X_2$ is some combination of $\hat{\beta}_1$ and $\hat{\beta}_2$. Note that $M_2 X_1$ are the residuals from a regression of $X_1$ on $X_2$.

## Theorem 17.3. Frisch-Waugh-Lovell.

The coefficient on $M_2 X_1$ is the original $\hat{\beta}_1$, i.e.,

$$\begin{aligned}
\hat{\beta}_1 &= (X_1^T M_2 X_1)^{-1} X_1^T M_2 y = ((X_1^T M_2)(M_2 X_1))^{-1} (X_1^T M_2) M_2 y \\
&= (\tilde{X}_1^T \tilde{X}_1)^{-1} \tilde{X}_1^T \tilde{y}
\end{aligned}$$

**PRACTICAL IMPLICATION.** If $K$ is large and we are primarily interested in first $K_1$ variables, then we can get $\hat{\beta}_1$ by regressing $y$ [or $M_2 y$ equivalently] on $M_2 X_1$ only, i.e.,

$$\hat{\beta}_1 = (X_1^T M_2 X_1)^{-1} X_1^T M_2 y = (X_1^T M_2 M_2 X_1)^{-1} X_1^T M_2 M_2 y.$$

This involves inversion of only $K_1 \times K_1$ and $K_2 \times K_2$ matrices, which involves less computing time than inverting $K \times K$ matrices, especially when $K$ is large [this computation can be as bad as $O(K^3)$].

**Example 17.2.**

Suppose that $X_2 = (1, 1, \ldots, 1)^T = i$ and $X_1 = x_1$, then

$$M_{X_2} = I_n - i(i^T i)^{-1} i^T = I_n - \frac{i i^T}{n}$$

$$M_{X_2} x_1 = x_1 - \frac{1}{n} \sum_{j=1}^{n} x_{ji} \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} = \begin{pmatrix} x_{11} - \bar{x}_1 \\ \vdots \\ x_{1n} - \bar{x}_1 \end{pmatrix}$$

$$\hat{\beta}_1 = \frac{\sum_{i=1}^{n} (x_{1i} - \bar{x}_1)(y_i - \bar{y})}{\sum_{i=1}^{n} (x_{1i} - \bar{x}_1)^2}$$

When regression includes an intercept, can first demean the $X$ variables (and the $y$'s) then do regression on the demeaned variables.

Other examples include seasonal components (dummy variables), trends, etc. Common practice is to deseasonalize data or detrend data before analyzing effect of interest. This argument shows it is justified.

## 17.3 Restricted Least Squares

Suppose we have $q$ linear restrictions on $\beta$, i.e.,

$$\begin{pmatrix} R_{11} & \cdots & R_{1K} \\ \vdots & & \vdots \\ R_{q1} & \cdots & R_{qK} \end{pmatrix} \begin{pmatrix} \beta_1 \\ \vdots \\ \beta_K \end{pmatrix} = R\beta = r = \begin{pmatrix} r_1 \\ \vdots \\ r_q \end{pmatrix} \tag{17.5}$$

where $R$ is $q \times K$ with $q < K$ and $R$ is of rank $q$. **Intuition** Sometimes economic theory dictates us to impose some restrictions on the coefficients.

**Example 17.3.**

Suppose that $\beta_1 + \cdots + \beta_K = 1$ (constant returns to scale). Then substituting in $\beta_K = 1 - \beta_1 - \cdots - \beta_{K-1}$, we obtain

$$X\beta = \sum_{j=1}^{K} x_j \beta_j = \sum_{j=1}^{K-1} x_j \beta_j + x_K (1 - \beta_1 - \cdots - \beta_{K-1})$$

$$= \sum_{j=1}^{K-1} (x_j - x_K) \beta_j + x_K.$$

This depends on only $\beta_1, \ldots, \beta_{K-1}$.

We want to minimize the quadratic form

$$S(\beta) = (y - X\beta)^T(y - X\beta)$$

subject to the restrictions. Let $\tilde{\beta}$ denote the solution.

Partition $X$, $\beta$, and $R$

$$X = (X_1 \quad X_2); \quad R = (R_1 \quad R_2); \quad \beta = \begin{pmatrix} \beta_1 \\ \beta_2 \end{pmatrix}$$

with dimensions $n \times (K - q)$, $n \times q$, $q \times (K - q)$, $q \times q$, $(K - q) \times 1$, $q \times 1$, respectively.

And

$$X_1\beta_1 + X_2\beta_2 = X\beta; \quad R_1\beta_1 + R_2\beta_2 = r,$$

where $R_2$ is of full rank $q$ and invertible. This may require some reordering. We can write

$$\beta_2 = R_2^{-1}(r - R_1\beta_1)$$

Therefore,

$$X\beta = X_1\beta_1 + X_2[R_2^{-1}(r - R_1\beta_1)] = (X_1 - X_2R_2^{-1}R_1)\beta_1 + X_2R_2^{-1}r.$$

In other words, we can find $\tilde{\beta}_1$ by minimizing

$$(y^* - X_1^*\tilde{\beta}_1)^T(y^* - X_1^*\tilde{\beta}_1)$$

with respect to $\beta_1$, where $y^* = y - X_2R_2^{-1}r$ and $X_1^* = X_1 - X_2R_2^{-1}R_1$ to get

$$\tilde{\beta}_1 = (X_1^{*T}X_1^*)^{-1}X_1^{*T}y^*$$
$$\tilde{\beta}_2 = R_2^{-1}(r - R_1\tilde{\beta}_1).$$

In fact there are other more convenient expressions for this.

## Theorem 17.4. Iterated Projection.

Suppose that $X = (X_1, X_2)$, and let $C(X_1)$ be the columns span of $X_1$ so that $C(X_1)$ is a subspace of $C(X)$. Then

$$P_{X_1} = P_{X_1}P_X.$$

This says that one can first project onto the bigger space and then onto the smaller space

**Proof.** Using the orthogonal representation of $C(X)$ we have

$$P_Xy = X_1\hat{\beta}_1 + M_{X_1}X_2\hat{\beta}_2$$

where $\hat{\beta}_1 = (X_1^TX_1)^{-1}X_1^Ty$, and

$$P_{X_1}P_Xy = P_{X_1}X_1\hat{\beta}_1 + P_{X_1}M_{X_1}X_2\hat{\beta}_2 = X_1\hat{\beta}_1 \quad \square$$

### Theorem 17.5. Consider the special case of (17.5) where $r = 0$. Then

$$\tilde{\beta} = \hat{\beta} - (X^TX)^{-1}R^T[R(X^TX)^{-1}R^T]^{-1}R\hat{\beta} = (I_K - (X^TX)^{-1}R^T[R(X^TX)^{-1}R^T]^{-1}R)\hat{\beta}$$

**Proof:** Define

$$S = \{z \in \mathbb{R}^n : z = X\beta, R\beta = 0\}$$

This is a subspace of $C(X)$, because if $z_1, z_2 \in S$, we have $\alpha_1 z_1 + \alpha_2 z_2 \in S$ for any scalars $\alpha_1, \alpha_2$. This is because there exist $\beta_1, \beta_2$ such that $z_1 = X\beta_1$, $z_2 = X\beta_2$ with $R\beta_1 = 0$, $R\beta_2 = 0$. Therefore,

$$\alpha_1 z_1 + \alpha_2 z_2 = \alpha_1 X\beta_1 + \alpha_2 X\beta_2 = X(\alpha_1\beta_1 + \alpha_2\beta_2)$$
$$R(\alpha_1\beta_1 + \alpha_2\beta_2) = \alpha_1 R\beta_1 + \alpha_2 R\beta_2 = 0.$$

Therefore, find the restricted least squares estimator using the principle of iterated projection - find $\tilde{y} \in S$ to minimize

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

The projection operator onto $S$ is

$$M_W = I_n - W(W^T W)^{-1}W^T; \quad W = X(X^T X)^{-1}R^T$$

Hence $M_W X\tilde{\beta} = X\tilde{\beta}$ and the result follows. $\square$

## 17.3.1 Backfitting in Linear Regression

We next consider an iterative approach to computing OLS estimators. Suppose that we have $y, x_1, \ldots, x_K$ vectors in $\mathbb{R}^n$. Define the projection operators $P_j = x_j(x_j^T x_j)^{-1}x_j^T$, $P_X = X(X^T X)^{-1}X^T$, $M_j = I_n - P_j$, and $M_X = I_n - P_X$, where $X = (x_1, \ldots, x_K)$. Thus

$$y = P_X y + M_X y = \hat{y} + \hat{\varepsilon}. \tag{17.6}$$

Suppose that one proceeds as follows

**Backfitting Algorithm**

1. First regresses $y$ on $x_1$ and get the residuals $M_1 y$.

2. Then regress $M_1 y$ on $x_2$ to get residuals $M_2 M_1 y$.

3. Continue doing one dimension regressions in the order $x_1$ then $x_2$ etc. until the process converges.

The residual after $K$ cycles of the backfitting algorithm is

$$\hat{\varepsilon}^K = T^K y; \quad T = M_K M_{K-1} \ldots M_1,$$

where $T : \mathbb{R}^n \to \mathbb{R}^n$.

## Theorem 17.6.

Suppose that $X$ is of full rank. Then, we have as $K \to \infty$

$$T^K y \to M_X y = \hat{\varepsilon}.$$

**Proof.** We prove that $T$ is a strict contraction mapping for $z \in C(X)$, that is $\|Tz\| < \|z\|$ for any $z \in C(X)$. This means that $T$ shrinks any such vectors. First, for any vector $z$ write $Tz = M_K M_{K-1} \ldots M_1 z = M_K v$ for $v = M_{K-1} \ldots M_1 z$. Then

$$\|Tz\|^2 = z^T T^2 z$$
$$= v^T M_K v$$
$$\leq \frac{v^T M_K v}{v^T v} v^T v$$
$$\leq \lambda_{max}(M_K)\|M_{K-1}\ldots M_1 z\|^2$$
$$\leq \|M_{K-1}\ldots M_1 z\|^2$$
$$\leq \|z\|^2$$

since $\lambda_{max}(M_j) = 1$ ($M_j$ are symmetric idempotent). Therefore, $T$ is a weak contraction. Furthermore, if $\|Tz\| = \|z\|$ then $\|M_1 z\| = \|z\|$ from the above argument. This implies that $z$ is orthogonal to the space spanned by $x_1$, i.e., $z \in C(x_1)^\perp$. Similarly one obtains that $z \in C(x_j)^\perp$, $j = 2, \ldots, K$. In other words, if $\|Tz\| = \|z\|$, then $z \in C(X)^\perp$. Therefore, if $z \in C(X)$, it must be that

$$\|Tz\| < \|z\| \leq (1 - \epsilon)\|z\|$$

for some $0 < \epsilon < 1$. Also, if $z \in C(X)$, then $Tz \in C(X)$. Hence

$$\|T^K z\| \leq (1 - \epsilon)\|T^{K-1} z\| \leq (1 - \epsilon)^K \|z\|.$$

Then combine with (17.6) we obtain the result. $\square$

This says that one can compute the regression of $y$ on $x_1, \ldots, x_K$ by computing a sequence of linear one dimensional regressions. This method is of interest when the number of covariates is large, since it avoids directly inverting the full design matrix.

## 17.3.2 Goodness of Fit

We may decompose the total sum of squares as follows using the orthogonality

$$y^T y = \hat{y}^T \hat{y} + \hat{\varepsilon}^T \hat{\varepsilon} \tag{17.7}$$

This is valid for any set of $X = (x_1, \ldots, x_K)$. It follows that

$$\hat{y}^T \hat{y} \leq y^T y, \tag{17.8}$$

i.e., the Euclidean norm of the projected vector is smaller than the norm of the original series. The notion of fit is widely used in practice. It captures the idea of how much of the variation in the data is explained by the covariates, i.e., by how much $\hat{y}^T \hat{y} \leq y^T y$. One possibility is to measure the fit by the residual sum of squares

$$S(\hat{\beta}) = RSS = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \|y - \hat{y}\|^2 = \hat{\varepsilon}^T \hat{\varepsilon} = y^T M_X y = y^T y - \hat{y}^T \hat{y},$$

which captures the Euclidean distance of the sample from the "fitted value". In general, the smaller the $RSS$ the better. However, the numerical value of $RSS$ depends on the units used to measure $y$ in so that one cannot compare across different $X$'s.

The idea is to compare the $RSS$ from a given $X$ with the $RSS$ of the $X$ that contains only ones. When $X = i$ ( $i = (1, \ldots, 1)^T$), that is, we just calculate the mean of $y$, we have

$$y^T y = \bar{y}^2 i^T i + \tilde{\varepsilon}^T \tilde{\varepsilon}$$

where $\tilde{\varepsilon} = y - \bar{y}i$

## Definition 17.3.

For a general $X = (x_1, \ldots, x_K)$, we define

$$R^2 = 1 - \frac{\hat{\varepsilon}^T \hat{\varepsilon}}{\tilde{\varepsilon}^T \tilde{\varepsilon}}.$$

**Remarks.**

1. When $X$ contains a column vector of ones, $0 \leq R^2 \leq 1$. If $X$ does not contain a column vector of ones, $R^2$ could be less than zero.

2. In the bivariate case, $R^2$ is the squared sample correlation between $y$ and $x$.

3. $R^2$ is invariant to some changes of units. If $y \to ay + b$ for any constants $a$, $b$, then $\hat{y}_i \to a\hat{y}_i + b$ and $\bar{y} \to a\bar{y} + b$, so $R^2$ is the same in this case. Clearly, if $X \to XA$ for a nonsingular matrix $A$, then $\hat{y}$ is unchanged, as is $\bar{y}$ and $\tilde{y}$.

4. $R^2$ always increases with the addition of variables. With $K = n$ we can make $R^2 = 1$.

5. $R^2$ can't be used to compare across different $y$.

# Exercises

## Exercise 1

Solution 1

Consider a dataset $(y, X)$ where $y = (2, 4, 6, 8)^T$ and $X = \begin{pmatrix} 1 & 1 \\ 1 & 2 \\ 1 & 3 \\ 1 & 4 \end{pmatrix}$. Calculate the OLS estimator $\hat{\beta}$ using the

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

## Exercise 2

Solution 2

Using the data from Exercise 1, calculate the fitted values $\hat{y}$ and the residuals $\hat{\varepsilon}$. Verify that $X^T \hat{\varepsilon} = 0$.

## Exercise 3

Solution 3

For the same dataset in Exercise 1, calculate the projection matrix $P_X$ and the annihilator matrix $M_X$. Verify that $P_X$ and $M_X$ are symmetric and idempotent.

## Exercise 4

Solution 4

Prove that for any matrix $X$, the matrices $P_X = X(X^TX)^{-1}X^T$ and $M_X = I - X(X^TX)^{-1}X^T$ are idempotent, that is, $P_X^2 = P_X$ and $M_X^2 = M_X$.

## Exercise 5

Explain in intuitive terms what the **column span** $C(X)$ of a matrix $X$ represents. Provide an example with a $3 \times 2$ matrix.

## Exercise 6

Explain the geometric interpretation of the **Ordinary Least Squares (OLS)** procedure. How does it relate to the projection of $y$ onto the column space of $X$?

## Exercise 7

Consider the matrices $X_1 = \begin{pmatrix} 1 \\ 2 \\ 3 \end{pmatrix}$ and $X_2 = \begin{pmatrix} 2 \\ 4 \\ 6 \end{pmatrix}$. Show that $C(X_1) = C(X_2)$. What does this imply about the relationship between $X_1$ and $X_2$?

## Exercise 8

State the **normal equations** in the context of OLS and explain how they are derived from the minimization of the sum of squared errors.

## Exercise 9

Explain the concept of **multicollinearity** and how it affects the uniqueness of the OLS estimator $\hat{\beta}$.

## Exercise 10

Given $X = \begin{pmatrix} 1 & 2 \\ 2 & 4 \end{pmatrix}$, show that $X^TX$ is not invertible. What does this indicate about the rank of $X$?

## Exercise 11

Consider a partitioned regression model where $X = (X_1, X_2)$. If $X_1^T X_2 = 0$, what can you say about relationship between $X_1$ and $X_2$ and what is the implication to $P_X$

## Exercise 12

State the **Frisch-Waugh-Lovell Theorem** and explain its practical implication for computing OLS estimates in a partitioned regression.

## Exercise 13

Suppose you have a regression model $y = X\beta + \varepsilon$, where $X = (X_1, X_2)$. You are only interested in estimating $\beta_1$. Using the **Frisch-Waugh-Lovell Theorem**, explain the steps to obtain $\hat{\beta}_1$.

## Exercise 14

Explain the **coordinate-free approach** in the context of OLS regression. What is the main advantage of this approach?

## Exercise 15

What does it mean for a matrix to be of **full rank**? How does the rank of the matrix $X$ affect the invertibility of $X^T X$?

## Exercise 16

Define the **residual sum of squares (RSS)** and explain how it is related to the goodness of fit of an OLS regression.

## Exercise 17

Define $R^2$ in the context of OLS regression. Explain its properties and potential limitations.

## Exercise 18

Explain the concept of **restricted least squares** and provide an example of a linear restriction on the coefficients.

## Exercise 19

State the concept of **iterated projection**.

## Exercise 20

Explain the process of **backfitting** in linear regression and its purpose.

# Solutions

## Solution 1

First, we calculate $X^T X$:

$$X^T X = \begin{pmatrix} 1 & 1 & 1 & 1 \\ 1 & 2 & 3 & 4 \end{pmatrix} \begin{pmatrix} 1 & 1 \\ 1 & 2 \\ 1 & 3 \\ 1 & 4 \end{pmatrix} = \begin{pmatrix} 4 & 10 \\ 10 & 30 \end{pmatrix}.$$

Next, we find the inverse of $X^T X$:

$$(X^T X)^{-1} = \frac{1}{(4)(30) - (10)(10)} \begin{pmatrix} 30 & -10 \\ -10 & 4 \end{pmatrix} = \frac{1}{20} \begin{pmatrix} 30 & -10 \\ -10 & 4 \end{pmatrix} = \begin{pmatrix} 1.5 & -0.5 \\ -0.5 & 0.2 \end{pmatrix}.$$

Now, we calculate $X^T y$:

$$X^T y = \begin{pmatrix} 1 & 1 & 1 & 1 \\ 1 & 2 & 3 & 4 \end{pmatrix} \begin{pmatrix} 2 \\ 4 \\ 6 \\ 8 \end{pmatrix} = \begin{pmatrix} 20 \\ 60 \end{pmatrix}.$$

Finally, we compute $\hat{\beta}$:

$$\hat{\beta} = (X^T X)^{-1} X^T y = \begin{pmatrix} 1.5 & -0.5 \\ -0.5 & 0.2 \end{pmatrix} \begin{pmatrix} 20 \\ 60 \end{pmatrix} = \begin{pmatrix} 0 \\ 2 \end{pmatrix}.$$

**Intuition:** The OLS estimator gives the coefficients of the linear combination of the columns of X that best approximates $y$.

## Solution 2

We have $\hat{\beta} = \begin{pmatrix} 0 \\ 2 \end{pmatrix}$ from Exercise 1. The fitted values are:

$$\hat{y} = X\hat{\beta} = \begin{pmatrix} 1 & 1 \\ 1 & 2 \\ 1 & 3 \\ 1 & 4 \end{pmatrix} \begin{pmatrix} 0 \\ 2 \end{pmatrix} = \begin{pmatrix} 2 \\ 4 \\ 6 \\ 8 \end{pmatrix}.$$

The residuals are:

$$\hat{\varepsilon} = y - \hat{y} = \begin{pmatrix} 2 \\ 4 \\ 6 \\ 8 \end{pmatrix} - \begin{pmatrix} 2 \\ 4 \\ 6 \\ 8 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \end{pmatrix}.$$

We verify that $X^T\hat{\varepsilon} = 0$:

$$X^T\hat{\varepsilon} = \begin{pmatrix} 1 & 1 & 1 & 1 \\ 1 & 2 & 3 & 4 \end{pmatrix} \begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}.$$

**Intuition:** The residuals represent the difference between observed and fitted values. The property $X^T\hat{\varepsilon} = 0$ reflects that residuals are orthogonal to the column space of X.

## Solution 3

[Exercise 3](#)

From Exercise 1, we have $(X^TX)^{-1} = \begin{pmatrix} 1.5 & -0.5 \\ -0.5 & 0.2 \end{pmatrix}$. The projection matrix $P_X$ is:

$$P_X = X(X^TX)^{-1}X^T = \begin{pmatrix} 1 & 1 \\ 1 & 2 \\ 1 & 3 \\ 1 & 4 \end{pmatrix} \begin{pmatrix} 1.5 & -0.5 \\ -0.5 & 0.2 \end{pmatrix} \begin{pmatrix} 1 & 1 & 1 & 1 \\ 1 & 2 & 3 & 4 \end{pmatrix} = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix}.$$

The annhilator matrix is given by:

$$M_X = I - P_X = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix} - \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix} = \begin{pmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix}.$$

$P_X$ and $M_X$ are symmetric by construction. It can be verified by computation that $P_X^2 = P_X$ and $M_X^2 = M_X$.
**Intuition:** $P_X$ projects any vector onto the column space of $X$, while $M_X$ projects any vector onto the orthogonal complement of the column space of $X$.

## Solution 4

[Exercise 4](#) To show $P_X$ is idempotent, we calculate $P_X^2$:

$$P_X^2 = [X(X^TX)^{-1}X^T][X(X^TX)^{-1}X^T] = X(X^TX)^{-1}(X^TX)(X^TX)^{-1}X^T = X(X^TX)^{-1}X^T = P_X.$$

To show $M_X$ is idempotent, we have:

$$M_X^2 = (I - P_X)(I - P_X)$$
$$= I - P_X - P_X + P_X^2$$
$$= I - P_X - P_X + P_X$$
$$= I - P_X = M_X.$$

**Intuition:** Idempotent matrices, when applied multiple times, have the same effect as when applied once.

## Solution 5

Exercise 5

The **column span** $C(X)$ of a matrix $X$ represents the set of all possible linear combinations of the columns of $X$. In other words, it's the subspace spanned by the column vectors of $X$.

For example, consider $X = \begin{pmatrix} 1 & 2 \\ 2 & 1 \\ 0 & 3 \end{pmatrix}$. The column span $C(X)$ is the set of all vectors of the form:

$$c_1 \begin{pmatrix} 1 \\ 2 \\ 0 \end{pmatrix} + c_2 \begin{pmatrix} 2 \\ 1 \\ 3 \end{pmatrix}$$

where $c_1$ and $c_2$ are any real numbers. Geometrically, this represents a plane in $\mathbb{R}^3$ passing through the origin.

**Intuition:** The column span represents all vectors that can be 'reached' by combining the columns of $X$.

## Solution 6

Exercise 6

The **Ordinary Least Squares (OLS)** procedure seeks to find the vector $\hat{\beta}$ that minimizes the distance between the vector $y$ and the vector $X\hat{\beta}$, which lies in the column space of $X$, $C(X)$. Geometrically, this is equivalent to finding the orthogonal projection of $y$ onto $C(X)$. The projection is the point in $C(X)$ that is closest to $y$. The vector connecting $y$ to its projection on $C(X)$ is the residual vector $\hat{\varepsilon}$, which is orthogonal to $C(X)$.

**Intuition:** Imagine shining a light perpendicular to the subspace $C(X)$. The shadow cast by $y$ onto $C(X)$ is its projection, $\hat{y}$.

## Solution 7

Exercise 7

Notice that $X_2 = 2X_1$. Therefore, any linear combination of the columns of $X_2$ can also be written as a linear combination of the columns of $X_1$, and vice-versa. Thus, $C(X_1) = C(X_2)$. This implies that $X_1$ and $X_2$ are linearly dependent; one is a scalar multiple of the other.

**Intuition:** If two matrices have the same column span, their columns span the same subspace, even if the columns themselves are different.

## Solution 8

The **normal equations** are given by $X^T X \hat{\beta} = X^T y$. They are derived by minimizing the sum of squared errors, $S(\beta) = (y - X\beta)^T (y - X\beta)$. Taking the derivative of $S(\beta)$ with respect to $\beta$ and setting it to zero gives:

$$\frac{\partial S(\beta)}{\partial \beta} = -2X^T (y - X\beta) = 0.$$

This simplifies to the normal equations: $X^T X \hat{\beta} = X^T y$.

**Intuition:** The normal equations represent the condition where the gradient of the sum of squared errors is zero, corresponding to a minimum.

## Solution 9

Exercise 9

**Multicollinearity** occurs when there is a linear dependence among the columns of the matrix $X$. In other words, one or more columns can be expressed as a linear combination of the other columns. When multicollinearity is present, the matrix $X^T X$ is not invertible, and the OLS estimator $\hat{\beta} = (X^T X)^{-1} X^T y$ is not uniquely defined. There are infinitely many solutions for $\hat{\beta}$.

**Intuition:** Multicollinearity means there is redundant information among the regressors, making it impossible to isolate the individual effect of each regressor.

## Solution 10

Exercise 10

Given $X = \begin{pmatrix} 1 & 2 \\ 2 & 4 \end{pmatrix}$, we have:

$$X^T X = \begin{pmatrix} 1 & 2 \\ 2 & 4 \end{pmatrix} \begin{pmatrix} 1 & 2 \\ 2 & 4 \end{pmatrix} = \begin{pmatrix} 5 & 10 \\ 10 & 20 \end{pmatrix}.$$

The determinant of $X^T X$ is $(5)(20) - (10)(10) = 0$. Since the determinant is zero, $X^T X$ is not invertible. This indicates that the rank of $X$ is less than the number of columns (which is 2). In fact, the rank of $X$ is 1, because the second column is twice the first column.

**Intuition:** A non-invertible $X^T X$ indicates linear dependence among the columns of $X$.

## Solution 11

Exercise 11

If $X_1^T X_2 = 0$, the columns of $X_1$ are orthogonal to the columns of $X_2$. This means that $X_1$ and $X_2$ represent completely independent sets of regressors. Also it implies that $P_X = P_{X_1} + P_{X_2}$

**Intuition**: The cross product between $X_1$ and $X_2$ represents the covariance between this regressors. If $X_1^T X_2 = 0$, this covariance is 0.

## Solution 12

The **Frisch-Waugh-Lovell Theorem** states that in a partitioned regression model $y = X_1\beta_1 + X_2\beta_2 + \varepsilon$, the OLS estimator of $\beta_1$ can be obtained by the following two-step procedure:

1. Regress $y$ on $X_2$ and obtain the residuals $M_2y$.
2. Regress $X_1$ on $X_2$ and obtain the residuals $M_2X_1$.
3. Regress $M_2y$ on $M_2X_1$ to obtain $\hat{\beta}_1$.

The practical implication is that we can obtain $\hat{\beta}_1$ without directly inverting the potentially large matrix $X^TX$. Instead, we only need to invert smaller matrices related to $X_1$ and $X_2$.

$$\hat{\beta}_1 = (X_1^T M_2 X_1)^{-1} X_1^T M_2 y$$

**Intuition:** The theorem allows us to isolate the effect of $X_1$ on $y$ after accounting for the effect of $X_2$.

## Solution 13

1. Calculate $M_2 = I - X_2(X_2^T X_2)^{-1} X_2^T$.
2. Calculate $M_2y$.
3. Calculate $M_2X_1$.
4. Calculate $\hat{\beta}_1 = (X_1^T M_2 X_1)^{-1} X_1^T M_2 y$.

**Intuition:** This process effectively removes the influence of $X_2$ from both $y$ and $X_1$ before estimating the relationship between the adjusted $y$ and $X_1$.

## Solution 14

The **coordinate-free approach** in OLS regression emphasizes the geometric interpretation of the regression in terms of the column space $C(X)$ rather than the specific choice of the regressors (the columns of $X$). The main advantage is that it highlights the invariance of certain results (like fitted values and residuals) to non-singular linear transformations of the regressors.

**Intuition:** The coordinate-free approach focuses on the underlying subspace spanned by the regressors, regardless of how that subspace is represented.

## Solution 15

A matrix is of **full rank** if its rank is equal to the smaller of its number of rows and columns. For a matrix $X$ with $n$ rows and $K$ columns, if $n > K$, full rank means rank$(X) = K$. If $n < K$, full rank means rank$(X) = n$. If $X$ is a $n \times K$ matrix, and rank$(X) = K$ (and given $n > K$), then $X^TX$ is invertible. If rank$(X) < K$, then $X^TX$ is not invertible.

**Intuition:** Full rank means that the columns (or rows) of the matrix are linearly independent.

## Solution 16

The **residual sum of squares (RSS)** is defined as:

$$RSS = \sum_{i=1}^{n} (y_i - \hat{y}_i)^2 = \hat{\varepsilon}^T \hat{\varepsilon}.$$

It measures the total squared difference between the observed values $y_i$ and the fitted values $\hat{y}_i$. A smaller RSS indicates a better fit of the model to the data, as the fitted values are closer to the observed values.

**Intuition:** RSS quantifies the unexplained variation in the data after fitting the OLS regression.

## Solution 17

$R^2$, or the coefficient of determination, is defined as:

$$R^2 = 1 - \frac{\hat{\varepsilon}^T \hat{\varepsilon}}{\tilde{\varepsilon}^T \tilde{\varepsilon}} = 1 - \frac{RSS}{TSS}$$

where $TSS = \sum (y_i - \bar{y})^2$ and $\tilde{\varepsilon} = y - \bar{y}i$.

It represents the proportion of the variance in the dependent variable $y$ that is explained by the independent variables in $X$.

**Properties:**

- If X contains a constant term, $0 \le R^2 \le 1$.
- $R^2$ increases as more variables are added to the model.

**Limitations:**

- $R^2$ can be misleadingly high for models with many variables.
- $R^2$ cannot be used to compare models with different dependent variables.
- If X does not contain a constant term, $R^2$ could be less than 0.

**Intuition:** $R^2$ measures the goodness of fit of the model, but it should be interpreted cautiously.

## Solution 18

**Restricted least squares** involves minimizing the sum of squared errors subject to one or more linear restrictions on the coefficients $\beta$. An example of a linear restriction is $\beta_1 + \beta_2 = 1$, which might represent a constraint that two coefficients sum to one (e.g., constant returns to scale in a production function).

**Intuition:** Restricted least squares incorporates prior information or theoretical constraints into the estimation process.

## Solution 19

If $X = (X_1, X_2)$, and let $C(X_1)$ be the columns span of $X_1$ so that $C(X_1)$ is a subspace of $C(X)$. Then $P_{X_1} = P_{X_1} P_X$.

**Intuition:** Projecting a vector first onto a larger space ($C(X)$) and then onto a subspace of that space ($C(X_1)$) is equivalent to projecting the vector directly onto the subspace ($C(X_1)$).

## Solution 20

**Backfitting** is an iterative procedure for computing OLS estimates. Given a set of regressors $x_1, \ldots, x_K$, the algorithm proceeds as follows:

1. Regress $y$ on $x_1$ and obtain residuals $M_1 y$.
2. Regress $M_1 y$ on $x_2$ and obtain residuals $M_2 M_1 y$.
3. Continue this process, regressing $M_{j-1} \ldots M_1 y$ on $x_j$ to obtain residuals $M_j M_{j-1} \ldots M_1 y$.
4. Repeat steps 1-3, cycling through the regressors until the residuals converge.

The purpose of backfitting is to provide an alternative way to compute OLS estimates, particularly when the number of regressors is large, and inverting $X^T X$ directly is computationally expensive.

**Intuition:** Backfitting successively removes the influence of each regressor from the residuals, iteratively approaching the overall OLS solution.

# R Scripts

## R Script 1: Basic OLS Estimation

```
# Load necessary library
library(tidyverse)

── Attaching core tidyverse packages ──────────────────────── tidyverse 2.0.0 ──
✓ dplyr      1.1.4     ✓ readr      2.1.5
✓ forcats    1.0.0     ✓ stringr    1.5.1
✓ ggplot2    3.5.1     ✓ tibble     3.2.1
✓ lubridate 1.9.4      ✓ tidyr      1.3.1
✓ purrr      1.0.2
── Conflicts ──────────────────────────────────────── tidyverse_conflicts() ──
✗ dplyr::filter() masks stats::filter()
✗ dplyr::lag()    masks stats::lag()
i Use the conflicted package (<http://conflicted.r-lib.org/>) to force all conflicts to become errors

# Set seed for reproducibility
set.seed(123)

# Simulate data
n <- 100  # Number of observations
K <- 2    # Number of regressors
X <- matrix(runif(n * K), nrow = n, ncol = K) # Generate random X matrix
X <- cbind(1, X) # add a constant
beta <- c(2, 0.5, -1)  # True coefficients
y <- X %*% beta + rnorm(n)  # Generate y with error term

# Convert to tibble for tidyverse compatibility
data <- as_tibble(cbind(y, X)) %>%
  rename(y = V1, intercept = V2, x1 = V3, x2 = V4)

Warning: The `x` argument of `as_tibble.matrix()` must have unique column names if
`.name_repair` is omitted as of tibble 2.0.0.
```

```
ℹ Using compatibility `.name_repair`.

# OLS estimation using lm()
model <- lm(y ~ ., data = data)

# Print model summary
summary(model)⊝


Call:
lm(formula = y ~ ., data = data)

Residuals:
    Min      1Q  Median      3Q     Max
-1.8994 -0.6821 -0.1086  0.5749  3.3663

Coefficients: (1 not defined because of singularities)
            Estimate Std. Error t value Pr(>|t|)
(Intercept)   1.8592     0.2862   6.497 3.53e-09 ***
intercept         NA         NA      NA       NA
x1            0.5973     0.3457   1.728  0.08720 .
x2           -1.0296     0.3735  -2.757  0.00697 **
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.9765 on 97 degrees of freedom
Multiple R-squared:  0.106, Adjusted R-squared:  0.0876
F-statistic: 5.753 on 2 and 97 DF,  p-value: 0.004356

# Extract coefficients
beta_hat <- coef(model)
print(beta_hat)⊝

(Intercept)   intercept          x1          x2
  1.8592062          NA   0.5972598  -1.0296422

# Calculate fitted values
y_hat <- predict(model)

# Calculate residuals
residuals <- resid(model)

# Verify orthogonality (should be close to zero)
t(as.matrix(data[, -1])) %*% residuals⊝

               [,1]
intercept -5.551115e-15
x1        -3.961511e-15
x2        -1.720181e-15
```

**Explanation:**

1. **Data Simulation:** We simulate a dataset with $n = 100$ observations and $K = 2$ regressors (plus a constant term, making it 3 coefficients). The X matrix is generated randomly, and y is created based on a linear model with known coefficients (beta) and a normally distributed error term.
   - The simulation process illustrates the **linear model** $y = X\beta + \varepsilon$.
2. **Data Preparation:** The simulated data is converted to a tibble (tidyverse's version of a data frame) and the columns are renamed for better readability.
3. **OLS Estimation:** The lm() function performs **Ordinary Least Squares (OLS)** regression. The formula y ~ . specifies that y is regressed on all other columns in the data tibble.
   - This step directly implements the **OLS procedure** described in Definition 17.1.

4. **Model Summary:** `summary(model)` provides various statistics, including coefficient estimates, standard errors, t-values, p-values, and R-squared.

5. **Coefficient Extraction:** `coef(model)` extracts the estimated coefficients ($\hat{\beta}$).
   - This corresponds to finding the $\hat{\beta}$ that minimizes $S(\beta)$ in Definition 17.1.

6. **Fitted Values and Residuals:** `predict(model)` calculates the fitted values ($\hat{y}$), and `resid(model)` calculates the residuals ($\hat{\varepsilon}$).
   - These are the $\hat{y}$ and $\hat{\varepsilon}$ discussed in Section 17.1.

7. **Orthogonality Verification:** `t(as.matrix(data[,-1])) %*% residuals` calculates $X^T\hat{\varepsilon}$. The result should be a vector very close to zero, demonstrating the orthogonality condition $X^T\hat{\varepsilon} = 0$.
   - This verifies empirically the important property described just after Definition 17.1.

# R Script 2: Projection Matrices and Orthogonality

```
# Load necessary libraries
library(tidyverse)

# Set seed for reproducibility
set.seed(456)

# Simulate data (same as before, but different seed)
n <- 50
K <- 3
X <- matrix(runif(n * K), nrow = n, ncol = K)
beta <- c(1, -0.5, 0.8)
y <- X %*% beta + rnorm(n)

# Calculate X'X
XtX <- t(X) %*% X

# Calculate (X'X)^(-1)
XtX_inv <- solve(XtX)

# Calculate projection matrix P_X
P_X <- X %*% XtX_inv %*% t(X)

# Calculate annihilator matrix M_X
M_X <- diag(n) - P_X  # diag(n) creates an n x n identity matrix

# Verify P_X is idempotent: P_X %*% P_X should be equal to P_X
print(all.equal(P_X %*% P_X, P_X))⬤

[1] TRUE

# Verify M_X is idempotent: M_X %*% M_X should be equal to M_X
print(all.equal(M_X %*% M_X, M_X))⬤

[1] TRUE

# Verify P_X is symmetric: t(P_X) should be equal to P_X
print(all.equal(t(P_X), P_X))⬤

[1] TRUE

# Verify M_X is symmetric: t(M_X) should be equal to M_X
print(all.equal(t(M_X), M_X))⬤

[1] TRUE
```

```
# Calculate fitted values
y_hat <- P_X %*% y

# Calculate residuals
epsilon_hat <- M_X %*% y

#Verify y = ŷ + ε
print(all.equal(y, y_hat+epsilon_hat))

[1] TRUE

# Verify orthogonality: X'ε should be close to zero
print(t(X) %*% epsilon_hat)

              [,1]
[1,] -1.815433e-15
[2,]  1.393903e-14
[3,]  6.695801e-15

# Visualization
plot(y, y_hat, xlab = "Actual Values (y)", ylab = "Fitted Values (ŷ)", main = "Actual vs. Fitted
        Values")
abline(0, 1, col = "red")  # Add a 45-degree line
```



**Explanation:**

1. **Data Simulation:** Similar to Script 1, we generate data for a linear model.
2. **Calculate $X^T X$ and $(X^T X)^{-1}$:** These are intermediate steps needed to calculate the projection matrices.
   - These calculations are part of finding $\hat{\beta} = (X^T X)^{-1} X^T y$.
3. **Calculate $P_X$ and $M_X$:** The projection matrix $P_X$ and the annihilator matrix $M_X$ are calculated using the formulas from Definition 17.2.
4. **Idempotency and Symmetry Checks:** The code verifies the key properties of $P_X$ and $M_X$ stated in Definition 17.2: that they are both idempotent ($P_X^2 = P_X$, $M_X^2 = M_X$) and symmetric ($P_X^T = P_X$, $M_X^T = M_X$). all.equal() is used for comparison, accounting for potential small numerical differences.
5. **Calculate $\hat{y}$ and $\hat{\varepsilon}$** Fitted values and residuals are computed using projection matrices: $\hat{y} = P_X y$ and $\hat{\varepsilon} = M_X y$.
6. **Verify** $y = \hat{y} + \hat{\varepsilon}$.
7. **Orthogonality Verification:** We check if $X^T \hat{\varepsilon}$ is close to zero, confirming the orthogonality condition.
8. **Visualization:** A scatter plot of actual values ($y$) vs. fitted values ($\hat{y}$) is created. A 45-degree line is added, representing perfect prediction. The closer the points are to the line, the better the fit.

## R Script 3: Frisch-Waugh-Lovell Theorem

```
# Load necessary libraries
library(tidyverse)

# Set seed for reproducibility
set.seed(789)

# Simulate data
n <- 100
K1 <- 2
K2 <- 3
X1 <- matrix(runif(n * K1), nrow = n, ncol = K1)
X2 <- matrix(runif(n * K2), nrow = n, ncol = K2)
```

```
X <- cbind(X1, X2)
beta <- c(1, -1, 0.5, -0.5, 0.2)
y <- X %*% beta + rnorm(n)

# Method 1: Full regression
full_model <- lm(y ~ X)
beta_hat_full <- coef(full_model)

# Method 2: Frisch-Waugh-Lovell Theorem

# Step 1: Regress y on X2 and get residuals
model_y_on_X2 <- lm(y ~ X2)
residuals_y_on_X2 <- resid(model_y_on_X2)

# Step 2: Regress X1 on X2 and get residuals
M2X1 <- matrix(nrow = n, ncol = K1)
for (i in 1:K1) {
  model_X1_on_X2 <- lm(X1[, i] ~ X2)
  M2X1[, i] <- resid(model_X1_on_X2)
}

# Step 3: Regress residuals_y_on_X2 on residuals_X1_on_X2
fwl_model <- lm(residuals_y_on_X2 ~ M2X1 - 1) # -1 removes the intercept
beta_hat_fwl <- coef(fwl_model)

# Compare coefficients for beta1 (full model vs. FWL)
print("Coefficients from full regression (first K1 coefficients):")

[1] "Coefficients from full regression (first K1 coefficients):"

print(beta_hat_full[2:(K1+1)])

       X1         X2
 0.8012567 -0.7152640

print("Coefficients from FWL regression:")

[1] "Coefficients from FWL regression:"

print(beta_hat_fwl)

     M2X11       M2X12
 0.8012567 -0.7152640

# They should be (almost) identical
```

**Explanation:**

1. **Data Simulation:** Data is generated with two sets of regressors, $X_1$ and $X_2$.
2. **Full Regression:** A standard OLS regression of $y$ on the full $X$ matrix (both $X_1$ and $X_2$) is performed.
3. **Frisch-Waugh-Lovell (FWL) - Step 1:** $y$ is regressed on $X_2$ only, and the residuals are stored. This isolates the part of $y$ that is *not* explained by $X_2$.
4. **FWL - Step 2:** Each column of $X_1$ is regressed on $X_2$, and the residuals are stored in M2X1. This isolates the part of each variable in $X_1$ that is *not* explained by $X_2$.
5. **FWL - Step 3:** The residuals from Step 1 (unexplained part of $y$) are regressed on the residuals from Step 2 (unexplained parts of $X_1$). The - 1 in the formula removes the intercept term, as we are only interested in the coefficients on $M_2X_1$.
   - This final regression isolates the relationship between $y$ and $X_1$ after controlling for $X_2$, directly implementing the **Frisch-Waugh-Lovell Theorem** (Theorem 17.3).

6. **Comparison:** The coefficients for $\beta_1$ (corresponding to $X_1$) from the full regression and the FWL regression are printed. They should be virtually identical, demonstrating the theorem.

## R Script 4: Restricted Least Squares

```
# Load necessary libraries
library(tidyverse)
library(quadprog) # For solve.QP()

# Set seed for reproducibility
set.seed(101)

# Simulate data
n <- 100
K <- 3
X <- matrix(runif(n * K), nrow = n, ncol = K)
beta <- c(2, -1, 0.5)
y <- X %*% beta + rnorm(n)

# Unrestricted OLS
unrestricted_model <- lm(y ~ X - 1) # -1 removes the intercept for easier comparison
beta_hat_unrestricted <- coef(unrestricted_model)

# Restricted Least Squares (beta1 + beta2 + beta3 = 1)
# We use quadratic programming (solve.QP) to solve the restricted LS

# Define Dmat (corresponds to X'X)
Dmat <- t(X) %*% X

# Define dvec (corresponds to X'y)
dvec <- t(X) %*% y

# Define Amat (constraint matrix) for beta1 + beta2 + beta3 = 1
Amat <- matrix(c(1, 1, 1), nrow = 1)  # Each row is a constraint

# Define bvec (constraint vector) for the constraint
bvec <- 1

# Solve the quadratic program
# solve.QP minimizes (1/2)b'Db - d'b subject to A'b >= b0
result <- solve.QP(Dmat = Dmat, dvec = dvec, Amat = t(Amat), bvec = bvec, meq = 1)

# Extract the solution (restricted coefficients)
beta_hat_restricted <- result$solution

# Print results
print("Unrestricted coefficients:")
```

[1] "Unrestricted coefficients:"

```
print(beta_hat_unrestricted)
```

```
       X1          X2          X3
 2.2218437 -0.9351643   0.2940957
```

```
print("Restricted coefficients (sum to 1):")
```

[1] "Restricted coefficients (sum to 1):"

```
print(beta_hat_restricted)
```

```
[1]  2.12931820 -1.07568076 -0.05363744
```

**Explanation:**

1. **Data Simulation:** Standard data generation for a linear model.
2. **Unrestricted OLS:** A standard OLS regression is performed without any restrictions. The intercept is removed (- 1) to simplify the comparison with the restricted estimates later, as the restriction will involve the coefficients directly.
3. **Restricted Least Squares:**
   - We use the `solve.QP()` function from the `quadprog` package, which solves quadratic programming problems. Restricted least squares can be formulated as a quadratic program.
   - `Dmat:` This corresponds to $X^T X$ in the OLS objective function.
   - `dvec:` This corresponds to $X^T y$ in the OLS objective function.
   - `Amat:` This matrix defines the linear constraints. In this case, we have one constraint: $\beta_1 + \beta_2 + \beta_3 = 1$. Since `solve.QP` expects constraints in the form $A^T b \geq b_0$, we transpose `Amat`.
   - `bvec:` This vector contains the right-hand side of the constraint (in this case, 1).
   - `meq = 1`. This parameter indicates that we have 1 equality constraint
   - `result$solution:` This contains the solution to the quadratic program, which are the restricted OLS coefficients.
4. **Print Results:** The unrestricted and restricted coefficient estimates are printed for comparison. The restricted coefficients will sum to 1, satisfying the constraint.
   - This script directly applies the concepts of **restricted least squares** discussed in Section 17.3.

# R Script 5: R-squared and Goodness of Fit

```
# Load necessary libraries
library(tidyverse)

# Set seed for reproducibility
set.seed(202)

# Simulate data
n <- 100
K <- 2
X <- matrix(runif(n * K), nrow = n, ncol = K)
X <- cbind(1, X)  # Add a constant term
beta <- c(2, 0.5, -1)
y <- X %*% beta + rnorm(n) # Model with good fit
y_poor_fit <- X %*% beta + rnorm(n, mean = 0, sd = 5) # Model with poor fit. High error

# Fit models
model_good <- lm(y ~ X)
model_poor <- lm(y_poor_fit ~ X)

# Get R-squared
r_squared_good <- summary(model_good)$r.squared
r_squared_poor <- summary(model_poor)$r.squared

# Print R-squared
print(paste("R-squared (good fit):", r_squared_good))
```

```
[1] "R-squared (good fit): 0.119451827407791"
```

```
print(paste("R-squared (poor fit):", r_squared_poor))
```

```
[1] "R-squared (poor fit): 0.0364354620408681"
```

```r
# Calculate RSS manually
rss_good <- sum(resid(model_good)^2)
rss_poor <- sum(resid(model_poor)^2)

# Calculate TSS manually
tss_good <- sum((y - mean(y))^2)
tss_poor <- sum((y_poor_fit - mean(y_poor_fit))^2)

# Calculate R-squared manually
r_squared_good_manual <- 1 - rss_good / tss_good
r_squared_poor_manual <- 1 - rss_poor / tss_poor

# Print manually calculated R-squared
print(paste("R-squared (good fit, manual):", r_squared_good_manual))

[1] "R-squared (good fit, manual): 0.119451827407791"

print(paste("R-squared (poor fit, manual):", r_squared_poor_manual))

[1] "R-squared (poor fit, manual): 0.036435462040868"

# Visualization
par(mfrow = c(1, 2)) # Set up a 1x2 plotting area

plot(predict(model_good), y, main = "Good Fit", xlab = "Fitted Values", ylab = "Actual Values")
abline(0, 1, col = "red")

plot(predict(model_poor), y_poor_fit, main = "Poor Fit", xlab = "Fitted Values", ylab = "Actual
        Values")
abline(0, 1, col = "red")
```



```r
par(mfrow = c(1, 1)) # restore plotting area
```

**Explanation:**

1. **Data Simulation:** We simulate two datasets: one with a relatively small error term (`y`) and one with a much larger error term (`y_poor_fit`), representing a good fit and a poor fit, respectively.
2. **Fit Models:** We fit OLS models to both datasets.
3. **Get R-squared:** We extract the $R^2$ value from the model summaries using `summary(model)$r.squared`.
4. **Print R-squared:** We print the $R^2$ values for both models. The good fit model should have a much higher $R^2$ than the poor fit model.
5. **Calculate RSS manually**: We calculate manually RSS as $RSS = \sum_{i=1}^{n}(y_i - \hat{y}_i)^2 = \hat{\varepsilon}^T\hat{\varepsilon}$.
6. **Calculate TSS manually**: We calculate manually TSS as $TSS = \sum_{i=1}^{n}(y_i - \bar{y})^2$
7. **Calculate R-squared manually**: We calculate manually R-squared as $R^2 = 1 - \dfrac{RSS}{TSS}$.
8. **Visualization:** We create scatter plots of fitted values vs. actual values for both models. The good fit model should show points clustered closely around the 45-degree line, while the poor fit model should show much more scatter.
   - This script illustrates the concept of **goodness of fit** and how $R^2$ (Definition 17.3) quantifies it. A higher $R^2$ corresponds to a model that explains a larger proportion of the variance in the dependent variable. The plots visually demonstrate the difference between a good fit and a poor fit.

# YouTube Video Recommendations

Here are some YouTube video recommendations related to the concepts in the provided text, along with explanations of their relevance. I have verified that all links are currently working as of October 26, 2023.

## 1. Ordinary Least Squares (OLS) and Linear Regression

- **Video Title:** "Linear Regression - Fun and Easy Machine Learning"

- **Channel:** Edureka

- **Link:** https://www.youtube.com/watch?v=E5RjzSK0fvY

- **Relevance:** This video provides a good introduction to linear regression and OLS. It covers the basic concepts, including the idea of minimizing the sum of squared errors, finding the best-fit line, and interpreting the results. It connects directly to **Definition 17.1** and the initial sections on the **projection approach**. It provides intuitive visualizations, making the core concept easier to grasp.

- **Video Title:** "Statistics 101: Linear Regression, The Very Basics"

- **Channel:** Brandon Foltz

- **Link:** https://www.youtube.com/watch?v=ZkjP5RJLQF4

- **Relevance:** Another excellent introductory video that goes through a simple linear regression example step-by-step. It covers the basics clearly, including the formula for the OLS estimators, and it provides an introduction to the geometrical interpretation. This is a good starting point for understanding **Section 17.1**.

## 2. Geometry of OLS and Projection

- **Video Title:** "The Geometric Interpretation of Least Squares"

- **Channel:** MathTheBeautiful

- **Link:** https://www.youtube.com/watch?v=MCixDRMCnhs

- **Relevance:** This video is *crucial* for understanding the **projection approach** in Section 17.1. It visually explains how OLS finds the projection of the dependent variable vector ($y$) onto the column space of the regressor matrix ($X$). It clearly illustrates concepts like the fitted values ($\hat{y}$), residuals ($\hat{\varepsilon}$), and the orthogonality condition ($X^T\hat{\varepsilon} = 0$). This video directly supports the geometric intuition behind OLS.

- **Video Title:** "1. The Geometry of Linear Equations"

- **Channel:** MIT OpenCourseWare

- **Link:** https://www.youtube.com/watch?v=J7DzL2_Na80

- **Relevance:** Part of Gilbert Strang's Linear Algebra course, this helps build the fundamental understanding of column spaces, linear combinations, and when systems of equations have solutions. While not directly about OLS, it builds the foundation for understanding *why* OLS works the way it does (projecting onto the column space). This relates to the discussion of $C(X)$ and the conditions under which $y = X\beta$ has a solution.

- **Video Title**: "Projection Matrix"

- **Channel**: joshstarmer

- **Link:** [https://www.youtube.com/watch?v=WAkRlhG2_m8](https://www.youtube.com/watch?v=WAkRlhG2_m8)

- **Relevance**: The video explains the geometrical meaning of OLS. It explains fitted values, residuals, and it explains and proves the properties of Projection and Annihilator matrices. It also presents many related concepts such as orthogonality.

# 3. Frisch-Waugh-Lovell (FWL) Theorem

- **Video Title:** "Frisch-Waugh-Lovell (FWL) Theorem"
- **Channel:** Ben Lambert
- **Link:** [https://www.youtube.com/watch?v=B_qjF0Pgbv4](https://www.youtube.com/watch?v=B_qjF0Pgbv4)
- **Relevance:** This video *specifically* explains the **Frisch-Waugh-Lovell Theorem (Theorem 17.3)**. It shows the derivation, the intuition (partialling out the effect of other variables), and demonstrates how to apply it. This is a *direct* explanation of Section 17.2. The presenter gives intuition and geometric explanation.

# 4. Multicollinearity

- **Video Title:** "Econometrics // Lecture 10: Multicollinearity"
- **Channel:** KeynesAcademy
- **Link:** [https://www.youtube.com/watch?v=tA9pIFJRy4c](https://www.youtube.com/watch?v=tA9pIFJRy4c)
- **Relevance:** This video explains the concept of **multicollinearity**, which is mentioned in the text when discussing the rank of $X$ and the uniqueness of the OLS estimator. It covers the causes, consequences, and detection of multicollinearity.

# 5. Restricted Least Squares

- **Video Title:** "Restricted Least Squares"
- **Channel:** Steve Grams
- **Link:** [https://www.youtube.com/watch?v=45T-zSzJq-g](https://www.youtube.com/watch?v=45T-zSzJq-g)
- **Relevance:** This video provides a clear explanation of **restricted least squares**, directly corresponding to **Section 17.3**. It shows how to set up the restrictions, derive the restricted estimator, and discusses the underlying theory.

# 6. Goodness of Fit (R-squared)

- **Video Title:** "R-squared, Clearly Explained!!!"
- **Channel:** StatQuest with Josh Starmer
- **Link:** [https://www.youtube.com/watch?v=OK2Bna1yx9k](https://www.youtube.com/watch?v=OK2Bna1yx9k)
- **Relevance:** This video explains the concept of $R^2$ (coefficient of determination) in a very intuitive way, including its interpretation and limitations. This relates directly to **Section 17.3.2** on **Goodness of Fit**.

# 7. Matrix Algebra Review (for background)

- **Video Series:** "Essence of Linear Algebra"
- **Channel:** 3Blue1Brown
- **Playlist Link:** [https://www.youtube.com/playlist?list=PLZHQObOWTQDPD3MizzM2xVFitgF8hE_ab](https://www.youtube.com/playlist?list=PLZHQObOWTQDPD3MizzM2xVFitgF8hE_ab)
- **Relevance:** This series is *highly* recommended for building a strong geometric intuition for linear algebra concepts. While not directly about OLS, understanding vectors, matrices, linear transformations, spans, and orthogonality is *essential* for understanding the material in the text. The concepts of linear independence, span, and basis are particularly relevant.

These videos cover the main theoretical concepts presented in the text, providing a mix of introductory explanations, geometric interpretations, and more advanced derivations. They should be helpful for students seeking to reinforce their understanding of the material.

# Multiple Choice Exercises

## MC Exercise 1

The Ordinary Least Squares (OLS) procedure chooses $\hat{\beta}$ to minimize which of the following?

    a. $y + X\beta$
    b. $(y - X\beta)^T(y - X\beta)$
    c. $X^T y$
    d. $(X^T X)^{-1}$

## MC Exercise 2

The column span of a matrix $X$, denoted $C(X)$, represents:

    a. The set of all linear combinations of the rows of $X$.
    b. The set of all linear combinations of the columns of $X$.
    c. The inverse of $X$.
    d. The determinant of $X$

## MC Exercise 3

In the OLS context, the fitted values $\hat{y}$ are defined as:

    a. $X\hat{\beta}$
    b. $y - X\hat{\beta}$
    c. $(X^T X)^{-1} X^T y$
    d. $X^T \hat{\varepsilon}$

## MC Exercise 4

The residuals $\hat{\varepsilon}$ in OLS regression satisfy which of the following orthogonality conditions?

    a. $y^T \hat{\varepsilon} = 0$
    b. $X^T y = 0$
    c. $X^T \hat{\varepsilon} = 0$
    d. $\hat{\varepsilon}^T \hat{\varepsilon} = 0$

## MC Exercise 5

The projection matrix $P_X$ is:

    a. Always invertible.
    b. Symmetric and idempotent.
    c. Equal to $X^T X$.
    d. Equal to $X$.

## MC Exercise 6

The annihilator matrix $M_X$ is defined as:

    a. $X(X^T X)^{-1} X^T$
    b. $X^T X$
    c. $I - X(X^T X)^{-1} X^T$
    d. $(X^T X)^{-1}$

## MC Exercise 7

If the rank of the matrix $X$ is equal to the number of columns of $X$ (and n > K), then:

    a. $X^T X$ is not invertible.
    b. The OLS estimator $\hat{\beta}$ is not unique.
    c. $X^T X$ is invertible.
    d. The residuals are always zero.

## MC Exercise 8

The normal equations in OLS are given by:

    a. $y = X\hat{\beta}$
    b. $X^T X\hat{\beta} = X^T y$
    c. $\hat{\varepsilon} = y - X\hat{\beta}$
    d. $X^T \hat{\varepsilon} = 0$

## MC Exercise 9

**Multicollinearity** occurs when:

    a. The dependent variable $y$ is constant.
    b. The columns of $X$ are linearly independent.
    c. The columns of $X$ are linearly dependent.
    d. The error term has a non-zero mean.

## MC Exercise 10

The Frisch-Waugh-Lovell (FWL) theorem allows us to:

    a. Calculate the inverse of $X^T X$ more easily.
    b. Estimate a subset of coefficients in a partitioned regression without inverting the full $X^T X$ matrix.
    c. Find the determinant of $X$.
    d. Always obtain an $R^2$ of 1.

## MC Exercise 11

In the partitioned regression $y = X_1 \beta_1 + X_2 \beta_2 + \varepsilon$, if $X_1^T X_2 = 0$, then:

    a. The columns of $X_1$ and $X_2$ are orthogonal.
    b. The OLS estimator is biased.
    c. The model is not identified.
    d. $R^2$ will be zero.

## MC Exercise 12

The **coordinate-free approach** to OLS emphasizes:

    a. The specific choice of regressors.
    b. The column space $C(X)$.
    c. The sample size $n$.
    d. The distribution of the error term.

## MC Exercise 13

If a matrix $A$ is idempotent, then:

    a. $A^{-1} = A$
    b. $A^2 = A$
    c. $A^T = A$
    d. $\det(A) = 0$

## MC Exercise 14

The **residual sum of squares (RSS)** measures:

    a. The total variation in $y$.
    b. The variation in $y$ explained by the model.
    c. The unexplained variation in $y$.

d. The correlation between $y$ and $x$.

## MC Exercise 15

$R^2$ (the coefficient of determination) is defined as:

a. $1 - \dfrac{TSS}{RSS}$

b. $1 - \dfrac{RSS}{TSS}$

c. $\dfrac{RSS}{TSS}$

d. $\dfrac{TSS}{RSS}$

## MC Exercise 16

Adding more variables to an OLS regression model will generally:

a. Decrease $R^2$.
b. Increase $R^2$.
c. Not affect $R^2$.
d. Make $R^2$ negative.

## MC Exercise 17

**Restricted least squares** is used when: a) The matrix X is not of full rank b) We want to impose linear restrictions on the coefficients. c) The error term is heteroskedastic. d) We do not have enough data.

## MC Exercise 18

**Iterated Projection** implies that: a) $P_{X_1} = P_X P_{X_1}$ b) $P_X = P_{X_1} P_{X_2}$ c) $P_{X_1} = P_{X_1} P_X$ d) $M_{X_1} = P_{X_1} P_X$

## MC Exercise 19

The **backfitting algorithm** is:

a. A method for finding the inverse of a matrix.
b. An iterative method for computing OLS estimates.
c. A way to calculate $R^2$.
d. A technique for data visualization.

## MC Exercise 20

If a matrix X is not of full rank:

    a. $\hat{\beta}$ is unique
    b. The normal equations have a unique solution.
    c. $X^T X$ is invertible
    d. $X^T X$ is not invertible.

# Multiple Choice Solutions

## MC Solution 1

**Answer:** b) $(y - X\beta)^T (y - X\beta)$

**Explanation:** The OLS procedure, by definition (Definition 17.1), minimizes the sum of squared errors, which is represented by the quadratic form $(y - X\beta)^T (y - X\beta)$.

## MC Solution 2

**Answer:** b) The set of all linear combinations of the columns of $X$.

**Explanation:** The column span, $C(X)$, is *defined* as the set of all possible linear combinations of the column vectors of $X$ (Section 17.1).

## MC Solution 3

**Answer:** a) $X\hat{\beta}$

**Explanation:** The fitted values, $\hat{y}$, represent the projection of $y$ onto the column space of $X$. This projection is given by $X\hat{\beta}$ (Section 17.1).

## MC Solution 4

**Answer:** c) $X^T \hat{\varepsilon} = 0$

**Explanation:** The key orthogonality condition in OLS is that the residuals, $\hat{\varepsilon}$, are orthogonal to the column space of $X$. This is expressed mathematically as $X^T \hat{\varepsilon} = 0$ (Section 17.1).

## MC Solution 5

**Answer:** b) Symmetric and idempotent.

**Explanation:** The projection matrix $P_X$ has two important properties: it is symmetric ($P_X^T = P_X$) and idempotent ($P_X^2 = P_X$) (Definition 17.2).

## MC Solution 6

**Answer:** c) $I - X(X^T X)^{-1} X^T$

**Explanation:** The annihilator matrix $M_X$ is defined as $I - P_X$, which is equivalent to $I - X(X^T X)^{-1} X^T$ (Definition 17.2).

## MC Solution 7

**Answer:** c) $X^T X$ is invertible.

**Explanation:** If the rank of $X$ equals the number of columns (and n > K), meaning the columns are linearly independent, then $X^T X$ is invertible (Section 17.1). This guarantees a unique solution for $\hat{\beta}$.

## MC Solution 8

**Answer:** b) $X^T X \hat{\beta} = X^T y$

**Explanation:** The normal equations, derived from minimizing the sum of squared errors, are $X^T X \hat{\beta} = X^T y$ (Equation 17.2).

## MC Solution 9

**Answer:** c) The columns of $X$ are linearly dependent.

**Explanation:** Multicollinearity is defined as the situation where there is linear dependence among the columns of the regressor matrix $X$ (Section 17.1, discussion of non-full rank $X$).

## MC Solution 10

**Answer:** b) Estimate a subset of coefficients in a partitioned regression without inverting the full $X^T X$ matrix.

**Explanation:** The Frisch-Waugh-Lovell theorem provides a way to estimate coefficients for a subset of regressors after partialling out the effects of other regressors, avoiding direct inversion of the potentially large $X^T X$ matrix (Section 17.2).

## MC Solution 11

**Answer:** a) The columns of $X_1$ and $X_2$ are orthogonal.

**Explanation:** The condition $X_1^T X_2 = 0$ *means* that the columns of $X_1$ are orthogonal to the columns of $X_2$ (Section 17.2, Theorem 17.2).

## MC Solution 12

MC Exercise 12

**Answer:** b) The column space $C(X)$.

**Explanation:** The coordinate-free approach emphasizes the geometric interpretation of OLS in terms of the column space $C(X)$, rather than the specific choice of basis vectors (columns of $X$) (Section 17.1).

## MC Solution 13

MC Exercise 13

**Answer:** b) $A^2 = A$

**Explanation:** An idempotent matrix is defined as a matrix that, when multiplied by itself, equals itself: $A^2 = A$ (Section 17.1, properties of $P_X$ and $M_X$).

## MC Solution 14

MC Exercise 14

**Answer:** c) The unexplained variation in $y$.

**Explanation:** The residual sum of squares (RSS) measures the squared differences between the observed values ($y_i$) and the fitted values ($\hat{y}_i$), thus representing the variation in $y$ *not* explained by the model (Section 17.3.2).

## MC Solution 15

MC Exercise 15

**Answer:** b) $1 - \dfrac{RSS}{TSS}$

**Explanation:** $R^2$ is defined as one minus the ratio of the residual sum of squares (RSS) to the total sum of squares (TSS): $R^2 = 1 - \frac{RSS}{TSS}$ (Definition 17.3).

## MC Solution 16

MC Exercise 16

**Answer:** b) Increase $R^2$.

**Explanation:** Adding more variables to an OLS model will generally increase $R^2$, even if the added variables are not truly related to the dependent variable (Section 17.3.2, Remark 4). This is because the model will always be able to explain at least as much, and usually slightly more, of the variation in $y$ with more variables.

## MC Solution 17

MC Exercise 17

**Answer:** b) We want to impose linear restrictions on the coefficients.

**Explanation:** Restricted least squares is used when we have prior information or theoretical constraints that we want to impose on the coefficient estimates (Section 17.3).

## MC Solution 18

MC Exercise 18

**Answer:** c) $P_{X_1} = P_{X_1} P_X$

**Explanation:** If $X = (X_1, X_2)$, and let $C(X_1)$ be the columns span of $X_1$ so that $C(X_1)$ is a subspace of $C(X)$. Then $P_{X_1} = P_{X_1} P_X$ (Section 17.3, Theorem 17.4)

## MC Solution 19

MC Exercise 19

**Answer:** b) An iterative method for computing OLS estimates.

**Explanation:** The backfitting algorithm is an iterative procedure for computing OLS coefficient estimates by successively regressing residuals on individual regressors (Section 17.3.1).

## MC Solution 20

MC Exercise 20

**Answer:** d) $X^T X$ is not invertible.

**Explanation:** If a matrix $X$ is not of full rank, its columns are linearly dependent, which means that $X^T X$ is not invertible, and the OLS estimator is not unique (Section 17.1).

Author: Peter Fuleky

This book was built with Quarto