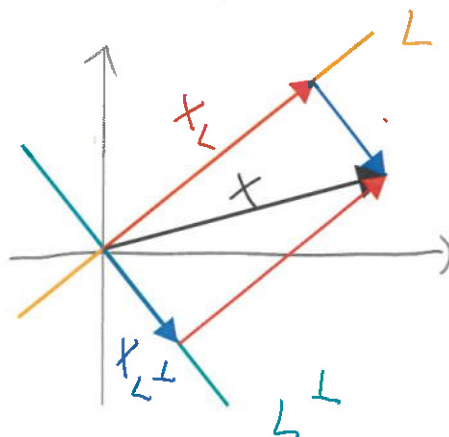Recall that any $\mathbf{x} \in \mathbb{R}^n$ can be uniquely written as $\mathbf{x} = \mathbf{x}_L + \mathbf{x}_{L^\perp}$, where $\mathbf{x}_L \in L$ and $\mathbf{x}_{L^\perp} \in L^\perp$.

$$\left.\begin{array}{l} P\mathbf{x}_L = \mathbf{x}_L \\ P\mathbf{x}_{L^\perp} = 0 \end{array}\right] \Rightarrow P\mathbf{x} = P(\mathbf{x}_L + \mathbf{x}_{L^\perp}) = P\mathbf{x}_L = \mathbf{x}_L$$

Let $\mathbf{x} \in \mathbb{R}^n$. Then $P^2\mathbf{x} = P(\underbrace{P\mathbf{x}}) = P\mathbf{x}_L = P\mathbf{x}$. Hence, $\underline{P^2 = P}$. $\qquad \mathbf{x}_L$

$\cdot \quad Y = Y_L + Y_{L^\perp} \qquad \text{(VECTOR)}$

$\cdot$ IF WE MULTIPLY AN ELEMENT OF $L$ WITH AN ~~œ~~ ELEMENT OF $L^\perp$

For all $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$: $\qquad\qquad\qquad$ WE GET $0$

$$\mathbf{x}^T P^T \mathbf{y} = (\underbrace{P\mathbf{x}}_{\in L})^T \mathbf{y} = (P\mathbf{x})^T \underbrace{\mathbf{y}_L}_{\mathbf{x}_L} = \mathbf{x}_L^T P\mathbf{y} = \mathbf{x}^T P\mathbf{y}$$

$\longrightarrow \mathbf{x}^T P y = (\mathbf{x}_L + \mathbf{x}_{L^\perp})^T P y = \mathbf{x}_L^T P y$

Hence, $P^T = P \qquad \longrightarrow P\mathbf{y} = \mathbf{y}_L$

$\longrightarrow (P\mathbf{x})^T \mathbf{y}_{L^\perp} = 0 \qquad$ BECAUSE

$\qquad\qquad\qquad\qquad\qquad\qquad \mathbf{x}_{L^\perp}^T P y = 0$

$\boxed{\Leftarrow:}$

Let $L$ be the space spanned by the columns of $P$. $\quad L = \text{Span}(P) = \{ P\mathbf{z} : \mathbf{z} \in \mathbb{R}^m \}$

- Let $\mathbf{x} \in L$. Then $\exists \mathbf{z} \in \mathbb{R}^n : \mathbf{x} = P\mathbf{z}$. Hence, $\underline{P\mathbf{x}} = P^2\mathbf{z} \overset{\text{idempot}}{=} P\mathbf{z} = \underline{\mathbf{x}}$.

- Let $\underline{\mathbf{x} \in L^\perp}$. Then for all $\mathbf{y} \in \mathbb{R}^n$: $(P\mathbf{x})^T \mathbf{y} = \mathbf{x}^T P^T \mathbf{y} \overset{\text{symm}}{=} \mathbf{x}^T \underbrace{P\mathbf{y}}_{\in L} = 0$. Hence
  
  $\underline{P\mathbf{x} = 0}$.

$\qquad\qquad\qquad\qquad\qquad\qquad e_i = (0, \dots, 0, \underset{\underset{i\text{-th}}{\uparrow}}{1}, 0, \dots 0)$

- The projection matrix is unique. Indeed, for each $i$, the vector $\mathbf{e}_i$ can be uniquely written as $\mathbf{e}_i = \mathbf{x} + \mathbf{y}$, where $\mathbf{x} \in L$ and $\mathbf{y} \in L^\perp$. Then the $i$th column of $P$ is $P\mathbf{e}_i = \mathbf{x}$.

- If $\mathbf{x}_1, \dots, \mathbf{x}_r$ are a basis of $L$ then the projection onto $L$ is given by

$$P = X(X^T X)^{-1} X^T,$$

where $X = (\mathbf{x}_1, \dots, \mathbf{x}_r)$. [prove this directly via the definition of the projection matrix or check $P^2 = P$, $P^T = P$, $\underbrace{\text{span}(P)}_{\text{space spanned by the columns of } P} = L$ or .] $\quad y \in \mathbb{R}^m$

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad Y = Y_L + Y_{L^\perp}$

- If $\mathbf{x}_1, \dots, \mathbf{x}_r$ are an *orthonormal* basis then $P = XX^T$. $\quad \rightarrow Y_L = X\mathbf{z}$, FOR SOME $\mathbf{z} \in \mathbb{R}^m$

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad X^T X = I \qquad P Y_L = X(X^T X)^{-1} X^T Y_L =$

- $I_n - P$ is the projection matrix onto $L^\perp$ $\qquad\qquad\qquad\qquad = X(X^T X)^{-1} X^T X\mathbf{z} =$

$\left[\begin{array}{l} Y = Y_L + Y_{L^\perp} \\ Y = I Y = (I - P)Y + PY = (I - P)Y + Y_L \end{array}\right. \qquad = X\mathbf{z} = Y_L$

$\longrightarrow (I - P)Y = Y_{L^\perp} \qquad\qquad\qquad\qquad P Y_{L^\perp} = X\underbrace{(X^T X)^{-1} X^T Y_{L^\perp}}_{=0}$

73

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad = 0$

(can be checked using original definition).

> **Example 54**
> $n = 3$
>
> If $L = \text{span}(\begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix})$ then $P = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{pmatrix}$.
>
> If $L = \text{span}(\begin{pmatrix} 1 \\ 1 \\ 0 \end{pmatrix}, \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix})$ then $P = \begin{pmatrix} 1/2 & 1/2 & 0 \\ 1/2 & 1/2 & 0 \\ 0 & 0 & 1 \end{pmatrix}$.
>
> If $L = \text{span}(\mathbf{x})$ for some $\mathbf{x} \in \mathbb{R}^n$ then $P = \frac{\mathbf{x}\mathbf{x}^T}{\mathbf{x}^T\mathbf{x}}$.

$$Py = \frac{\mathbf{x}\,\mathbf{x}^T y}{\mathbf{x}^T\mathbf{x}} = \frac{\mathbf{x}^T y\,\mathbf{x}}{\|\mathbf{x}\|^2} =$$

$$= \frac{\langle \mathbf{x}, y \rangle\,\mathbf{x}}{\|\mathbf{x}\|^2}$$

> **Lemma 12**
> If $A$ is an $n \times n$ projection matrix (i.e. $A = A^T$, $A^2 = A$) of rank $r$ then
>
> 1. $r$ of the eigenvalues of $A$ are 1 and $n - r$ are 0,
>
> 2. $\text{rank } A = \text{trace } A$,

IDEMPOTENT
↓

**Proof** Let $\mathbf{x}$ be an eigenvector of $A$, with eigenvalue $\lambda$. Then $\lambda \mathbf{x} = A\mathbf{x} = A^2\mathbf{x} = A\lambda\mathbf{x} = \lambda A\mathbf{x} = \lambda^2 \mathbf{x}$. $\overset{\mathbf{x}\neq 0}{\Longrightarrow} \lambda = \lambda^2 \implies \lambda \in \{0, 1\}$.

1. $A$ symmetric $\implies \exists P$ (orthogonal) s.t. $P^{-1}AP = D$, where $D$ is diagonal with 0s and 1s on the diagonal. Since $P$ is non-singular, $\text{rank } A = \text{rank } D$. Hence $D$ has $r$ ones down the diagonal.
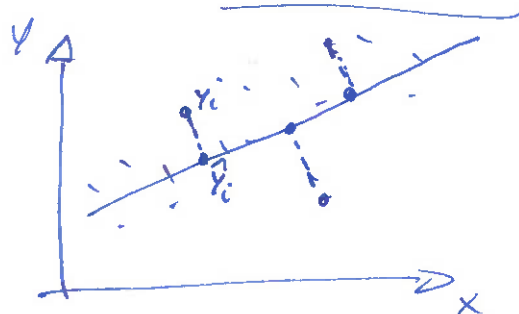
2. $\text{trace}(A) = \text{trace}(APP^{-1}) \overset{L46}{=} \text{trace}(P^{-1}AP) = \text{trace } D = \text{rank } A$

$$\text{trace}(AB) = \text{trace}(BA)$$

## 9.9 Residuals, Estimation of the variance

> **Definition 20**
> $\hat{\mathbf{Y}} = X\hat{\beta}$, where $\hat{\beta}$ is a least squares estimator, is called the *vector of fitted values*.

$$\hat{Y} = X\hat{\beta} \qquad \text{BY DEF OF } \hat{\beta}$$

In the full rank case, $\hat{\mathbf{Y}} = \underbrace{X(X^TX)^{-1}X^T}_{P}\mathbf{Y}$.

> **Lemma 13**
> $\hat{\mathbf{Y}}$ is unique and
> $$\hat{\mathbf{Y}} = P\mathbf{Y},$$
> where $P$ is the projection matrix onto the column space of $X$.

$$\text{span}(X) = \{X z : z \in \mathbb{R}^n\}$$

Because of this lemma, $P$ is sometimes called the *hat matrix* (it puts the hat on $\mathbf{Y}$, i.e. $\hat{\mathbf{Y}} = P\mathbf{Y}$).

**Proof** Suppose $\hat{\beta}$ is a LSE of $\beta$. We already know $P$ is unique, hence $P\mathbf{Y}$ is unique. Thus it suffices to show $\hat{\mathbf{Y}} = P\mathbf{Y}$. Since $P\mathbf{Y} \in \text{span}(X)$ there exists $\gamma$ s.t. $X\gamma = P\mathbf{Y}$. Then

$$S(\hat{\beta}) = \|Y - X\hat{\beta}\|^2$$

$$S(\gamma) =$$
$$= \|Y - X\gamma\|^2$$

$$\gamma \in \mathbb{R}^m$$

$$S(\hat{\beta}) = \|\mathbf{Y} - P\mathbf{Y} + P\mathbf{Y} - X\hat{\beta}\|^2$$

$$= \underbrace{\|\mathbf{Y} - P\mathbf{Y}\|^2}_{=S(\gamma)} + \|P\mathbf{Y} - X\hat{\beta}\|^2 + 2\underbrace{\underbrace{(\mathbf{Y} - P\mathbf{Y})^T}_{=\mathbf{Y}^T(I-P)} \underbrace{(P\mathbf{Y} - X\hat{\beta})}_{\in \text{span}(X)}}_{=0}$$

$$\geq S(\hat{\beta}) + \|P\mathbf{Y} - X\hat{\beta}\|^2,$$

since $\hat{\beta}$ minimises $S$. Thus $\|P\mathbf{Y} - X\hat{\beta}\| = 0$. Therefore, $P\mathbf{Y} = X\hat{\beta}. = \hat{Y}$

$$\uparrow$$
$$\text{BY DEF OF } \hat{Y}$$

> **Definition 21**
> $\mathbf{e} = \mathbf{Y} - \hat{\mathbf{Y}}$ is called the *vector of residuals*.

**Remark** Equivalent form: Using Lemma 13,

$$e = (I - P)Y = QY$$

$$\mathbf{e} = \mathbf{Y} - P\mathbf{Y} = Q\mathbf{Y},$$

where $Q = I - P$ is the projection matrix onto $\text{span}(X)^{\perp}$.

**Remark**

$$E[Y] = X\beta$$
$$\uparrow$$
$$E(\mathbf{e}) = E[Q\mathbf{Y}] = Q E\mathbf{Y} = \underbrace{QX}_{=0}\beta = \mathbf{0}$$
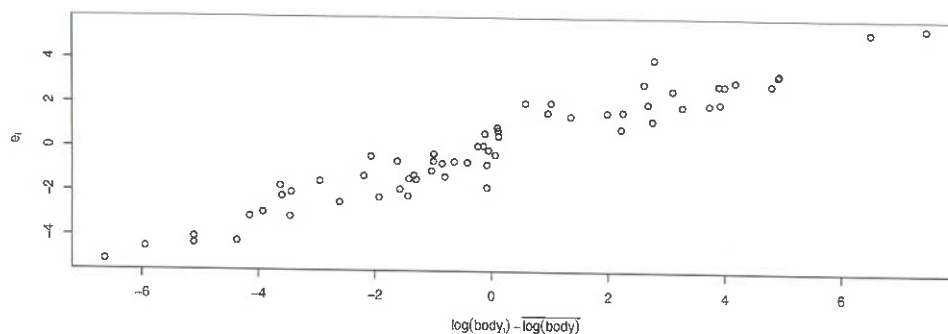
$\mathbf{e}$ can be used to see how well the model and the data agree and to see if certain observations are larger or smaller than predicted by the model.

Suppose we suspect $\mathbf{Z} = \mathbf{o}$ might be important. To investigate this, we plot $(Q\mathbf{o})_i = o_i - \bar{o}$ vs $e_i$. If the model (1) is true then the plot below should roughly look like the previous plot.



The fit of (1) does not seem to be good; however simply including $\mathbf{o}$ in the model does not seem to be reasonable because the above plot does not look like a linear relationship.

Let $z_j = \log(o_j)$. A plot of $(Q\mathbf{z})_i$ vs $e_i$:



This looks like a linear relationship with slope $\neq 0 \to$ could include $z_i \beta_2$ in model (1).

## Residual Sum of Squares

**Definition 22**
RSS $= \mathbf{e}^T \mathbf{e}$ is called the *residual sum of squares*.

RSS quantifies the departure of the data from the model. It is the minimum of $S(\beta)$.

77

**Remark** Other forms:

PY ↑

- RSS $= \sum_{i=1}^n e_i^2$

- RSS $= S(\hat{\beta}) = \|Y - X\hat{\beta}\|^2$ $= \|Y - \hat{Y}\|^2 = \|e\|^2$

- RSS $= (Q\mathbf{Y})^T Q\mathbf{Y} = \mathbf{Y}^T Q^T Q\mathbf{Y} = \mathbf{Y}^T Q\mathbf{Y}$

  ↳ $Q$ IS SYMMETRIC AND IDEMPOTENT
  
  (underbrace: $e$)

- RSS $= \mathbf{Y}^T\mathbf{Y} - \hat{\mathbf{Y}}^T\hat{\mathbf{Y}}$.

  Indeed, RSS $= (\mathbf{Y} - \hat{\mathbf{Y}})^T(\mathbf{Y} - \hat{\mathbf{Y}}) = \mathbf{Y}^T\mathbf{Y} - 2\hat{\mathbf{Y}}^T\mathbf{Y} + \hat{\mathbf{Y}}^T\hat{\mathbf{Y}} = \mathbf{Y}^T\mathbf{Y} - \hat{\mathbf{Y}}^T\hat{\mathbf{Y}}$.
  The last equality holds because $\hat{\mathbf{Y}}^T\mathbf{Y} = (P\mathbf{Y})^T\mathbf{Y} = (PP\mathbf{Y})^T\mathbf{Y} = \mathbf{Y}^T P^T P^T\mathbf{Y} =$
  $(P\mathbf{Y})^T P\mathbf{Y} = \hat{\mathbf{Y}}^T\hat{\mathbf{Y}}$.

  ↑ BY DEF $\hat{Y}$  ↳ $P^2 = P$

---

> **Theorem 8**
> $\hat{\sigma}^2 := \frac{\text{RSS}}{n-r}$ is an unbiased estimator of $\sigma^2$.

$I\sigma^2 = \text{cov}(\varepsilon)$

BY SOA

Recall: $r = \text{rank}(X) = \text{rank}(P)$

**Proof** Let $Q = I - P$. Since $P$ is a projection matrix, $Q$ is a projection matrix as
well. Hence, RSS $= \mathbf{Y}^T Q\mathbf{Y}$.

BECAUSE RSS IS A ~~CERT~~ ONE-DIMENSIONAL OBJECT

trace$(AB)$ = trace$(BA)$ BY LINEARITY OF TRACE

$E(\text{RSS}) = E\,\text{trace}\,\text{RSS} = E\,\text{trace}(\mathbf{Y}^T Q\mathbf{Y}) \overset{\text{Le }6}{=} E\,\text{trace}(Q\mathbf{Y}\mathbf{Y}^T) = \text{trace}(Q\,E(\mathbf{Y}\mathbf{Y}^T))$

$= \text{trace}(Q[\text{cov}\,\mathbf{Y} + E(\mathbf{Y})\,E(\mathbf{Y})^T]) = \text{trace}(Q\sigma^2) + \text{trace}(QX\beta(X\beta)^T)$

$\in \text{SPAN}(X)$  ↓  $= \sigma^2\,\text{trace}(I - P) + 0 = \sigma^2(n - \text{trace}(P))$

$QX\beta = 0$

$\overset{\text{Le }12}{=} \sigma^2(n - \text{rank}(P)) = \sigma^2(n - r)$.

$E[YY^T] = \text{Cov}(Y) + E[Y]E[Y]^T$

$\text{Cov}(Y) = \sigma^2 I$

**Remark** This is a generalisation of the result that the sample variance $s^2$ is an unbiased
estimator for $\sigma^2$ when $Y_1, \ldots, Y_n$ are i.i.d. with unknown mean $\mu$ and unknown variance
$\sigma^2$.

$\underset{\boxed{1}}{X}$

Indeed, we can write this iid setup as the linear model $\mathbf{Y} = \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} \mu + \epsilon$ with $E\,\epsilon = 0$ and

$\text{cov}\,\epsilon = \sigma^2 I$. Then $P = \mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T = \frac{1}{n}\mathbf{X}\mathbf{X}^T$ and thus $\mathbf{e} = \mathbf{Y} - \hat{\mathbf{Y}} = \mathbf{Y} - P\mathbf{Y} =$

$\underset{A}{\smile}$   $X^T X = n$

$= Y - \frac{1}{n}\begin{pmatrix} 1 & \cdots & 1 \\ \vdots & & \vdots \\ 1 & \cdots & 1 \end{pmatrix} Y$

$= Y - \begin{pmatrix} \bar{Y} \\ \vdots \\ \bar{Y} \end{pmatrix}$

78

$$e = \mathbf{Y} - \begin{pmatrix} \bar{Y} \\ \vdots \\ \bar{Y} \end{pmatrix}. \text{ Hence,}$$

$$e^\mathsf{T} e = \sum_{i=1}^{n} (Y_i - \bar{Y})^2$$

$$\frac{RSS}{n - r} = \underbrace{\frac{\sum (Y_i - \bar{Y})^2}{n - 1}}_{=s^2 = \text{sample variance}} = s^2$$

which we already know is unbiased for $\sigma^2$.

$$E[s^2] = \sigma^2$$

## Coefficient of Determination ($R^2$)

In the simplest model with only an intercept term, i.e. in

$$\mathbf{Y} = \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} \beta_1 + \epsilon, \quad \mathsf{E}\,\epsilon = 0$$

we have $RSS = \sum_{i=1}^{n}(Y_i - \bar{Y})^2$. Larger models, i.e. models with more columns in $X$ will only lead to smaller RSS.

For models containing an *intercept term*, (i.e. $X$ contains a column consisting of 1s (or any other constant)), a popular measure of the quality of a model is

$$R^2 = 1 - \frac{RSS}{\sum_{i=1}^{n}(Y_i - \bar{Y})^2},$$

called the *coefficient of determination* or simply $R^2$. A smaller RSS is "better", thus we want a large $R^2$. Note: $0 \leq R^2 \leq 1$ and $R^2 = 1$ for a "perfect" model.

**Remark (Intuitive interpretation)** $RSS / n$ is an estimator of $\sigma^2$. $\frac{1}{n}\sum_{i=1}^{n}(Y_i - \bar{Y})^2$ is an estimator of $\sigma^2$ in the model with only the intercept term (let us call this the "total variance").

Thus $\frac{RSS / n}{\frac{1}{n}\sum_{i=1}^{n}(Y_i - \bar{Y})^2} \approx \frac{\text{Variance in the model}}{\text{total variance}}$ and hence

$$R^2 \approx \frac{\text{total variance} - \text{variance in model}}{\text{total variance}}$$

Hence, $R^2 \approx$ fraction of the total variance of the data that "is explained" by the model.