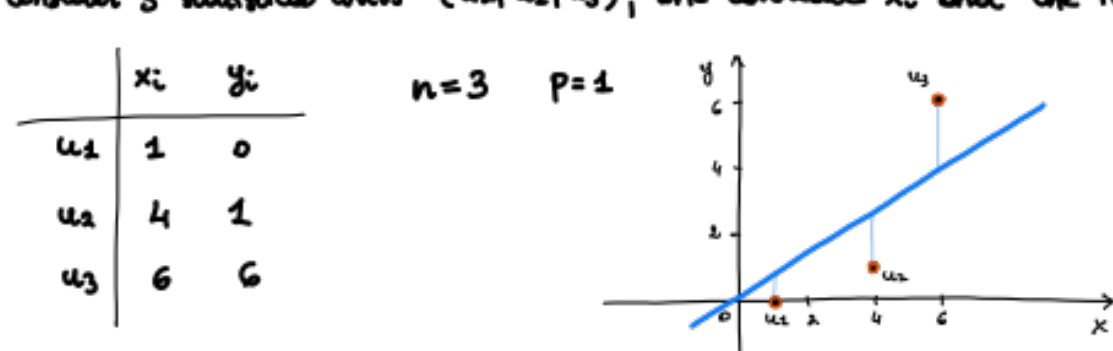


GEOMETRIC INTERPRETATION

Let's start with a simple example.

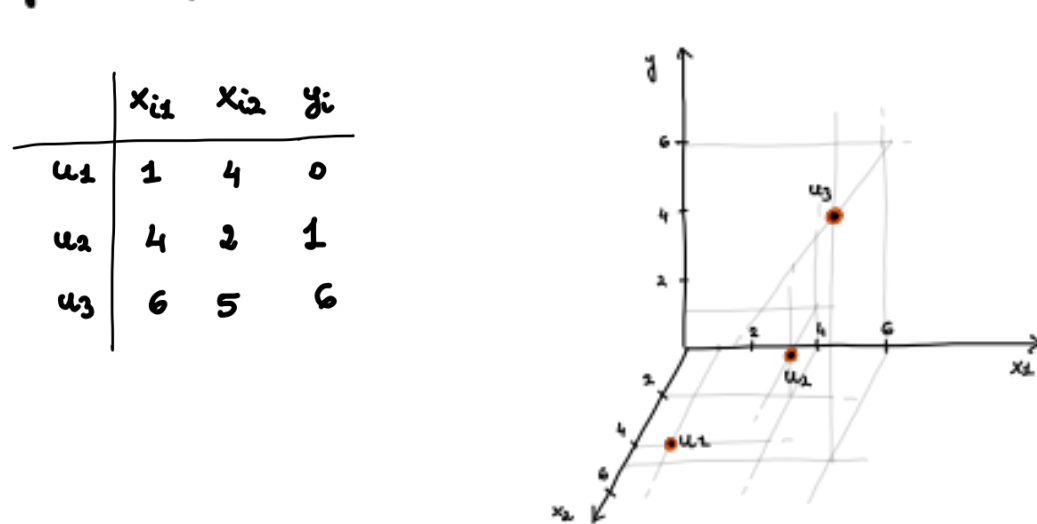
consider 3 statistical units (u_1, u_2, u_3) , one covariate x_i and the response y_i



Our problem up to now was:

1 look for the line that minimizes the "vertical distances".

If we consider now the same units, but 2 covariates x_{i1} and x_{i2}



We represent n points in a $(p+1)$ -dimensional space : n points in \mathbb{R}^{p+1}

units # covariates + 1 (y)

where the coordinates of each point are the values assumed by the p covariates and the response.

In the multiple linear model we have $\underline{Y} = \underline{X}\underline{\beta} + \underline{\varepsilon}$

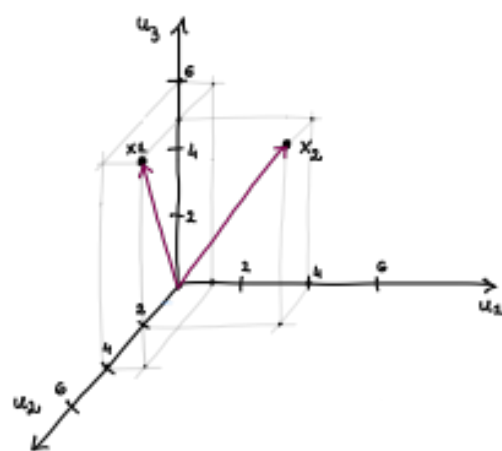
where $\underline{X} = [\underline{x}_1 \ \underline{x}_2 \ \dots \ \underline{x}_p]$, and the columns are p n -dimensional vectors

→ we can change perspective on the data: now UNITS ARE THE AXES

VARIABLES ARE VECTORS

We represent p vectors in a n -dimensional space : p n -dimensional vectors in \mathbb{R}^n

The coordinates of each vector are the observations of that variable on the n units



$p=2$ n -dimensional linearly independent vectors in an n -dimensional space.

On this space, we can define the set of all possible LINEAR COMBINATIONS of $\underline{x}_1, \dots, \underline{x}_p$

calligraphic " \mathcal{C} " $\mathcal{C}(X) = \{ \underline{\mu} \in \mathbb{R}^n : \underline{\mu} = \underline{X}\underline{\beta} = \beta_1 \underline{x}_1 + \beta_2 \underline{x}_2 + \dots + \beta_p \underline{x}_p, \ \underline{\beta} \in \mathbb{R}^p \}$

In particular, $\mathcal{C}(X)$ is the SUBSPACE of \mathbb{R}^n generated by $(\underline{x}_1, \dots, \underline{x}_p)$.

↳ p linearly indep. vectors

⇒ $\mathcal{C}(X)$ has dimension p

In our example, the 2 vectors identify a plane (2-dim space)

→ any linear combination of \underline{x}_1 and \underline{x}_2 will lie on this plane

If we call $\underline{X} = [\underline{x}_1 \ \underline{x}_2]$, $(n \times p) = (3 \times 2)$ matrix

$\mathcal{C}(X) = \beta_1 \underline{x}_1 + \beta_2 \underline{x}_2$ the column space of \underline{X}

$\mathcal{C}(X)$ is a subspace of \mathbb{R}^3 of dimension 2 ⇒ any $\underline{\mu} = \beta_1 \underline{x}_1 + \beta_2 \underline{x}_2$ will lie on $\mathcal{C}(X)$.

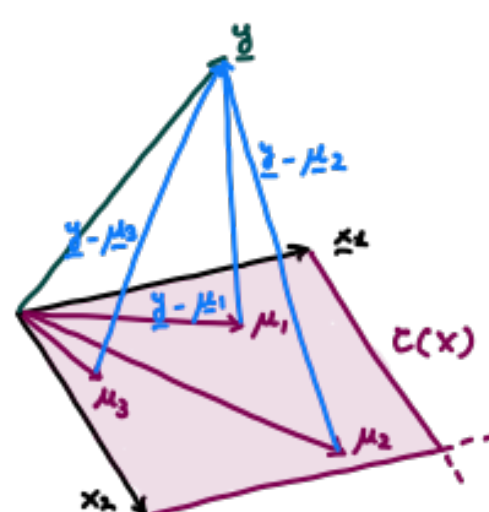
The assumption of LINEARITY is equivalent to asking that the mean of \underline{Y} belongs to $\mathcal{C}(X)$.

For a given $(\beta_1, \beta_2) = \underline{\beta}$, $\underline{\mu} = \underline{X}\underline{\beta}$ is a vector in the subspace

When we introduce \underline{y} , in general it will not lie on $\mathcal{C}(X)$

Now, consider \underline{y} and a generic vector of $\mathcal{C}(X)$ $\underline{\mu} = \underline{X}\underline{\beta}$.

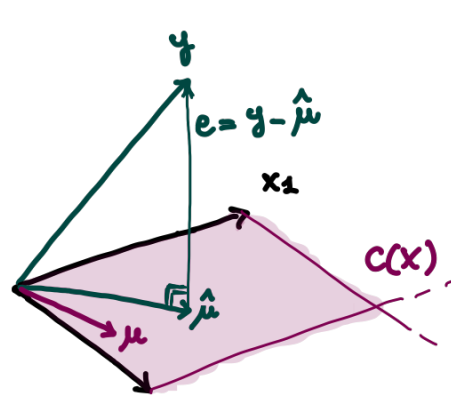
$\underline{y} - \underline{X}\underline{\beta}$ is the difference between the response and that vector of $\mathcal{C}(X)$.



$(\underline{y} - \underline{X}\underline{\beta})^T (\underline{y} - \underline{X}\underline{\beta}) = S(\underline{\beta})$ is the squared length of the difference

⇒ minimizing $S(\underline{\beta})$ means finding, in $\mathcal{C}(X)$, the vector $\underline{X}\underline{\hat{\beta}}$ so that $\underline{y} - \underline{X}\underline{\hat{\beta}}$ has minimum length.

⇒ we want $\underline{y} - \underline{X}\underline{\hat{\beta}}$ to be orthogonal to $\mathcal{C}(X)$ (hence $\underline{y} - \underline{X}\underline{\hat{\beta}}$ is orthogonal to the columns $\underline{x}_1, \dots, \underline{x}_p$ of \underline{X})



Indeed, $\underline{\hat{y}} = \underline{\hat{\mu}} = \underline{X}\underline{\hat{\beta}}$ is the ORTHOGONAL PROJECTION of \underline{y} onto $\mathcal{C}(X)$

⇒ $\underline{y} - \underline{X}\underline{\hat{\beta}} \perp \mathcal{C}(X)$

⇒ $\underline{y} - \underline{X}\underline{\hat{\beta}} \perp \underline{x}_j$ for all $j=1, \dots, p$ *

* orthogonality: $\begin{cases} (\underline{y} - \underline{X}\underline{\hat{\beta}})^T \underline{x}_1 = 0 \\ \vdots \\ (\underline{y} - \underline{X}\underline{\hat{\beta}})^T \underline{x}_p = 0 \end{cases}$
↓
normal equations

$\underline{\hat{\mu}} = \underline{\hat{y}} = \underline{X}\underline{\hat{\beta}} = \underline{X}(\underline{X}^T \underline{X})^{-1} \underline{X}^T \underline{y} = \underline{P}\underline{y}$ and $\underline{P} = \underline{X}(\underline{X}^T \underline{X})^{-1} \underline{X}^T$ is the projection matrix
($n \times n$), symmetric, idempotent, with rank= p
 $\underline{P}^T = \underline{P}$ $\underline{P}^2 = \underline{P}$

The vector of residuals $\underline{e} = \underline{y} - \underline{\hat{y}} = \underline{y} - \underline{P}\underline{y} = (\underline{I}_n - \underline{P})\underline{y}$ is also a projection of \underline{y} .

\underline{e} is the projection of \underline{y} on the subspace of \mathbb{R}^n perpendicular to $\mathcal{C}(X)$: $\underline{e} \perp \mathcal{C}(X)$.

$(\underline{I}_n - \underline{P})$ is also a projection matrix of rank $n-p$ (it projects on the space $\perp \mathcal{C}(X)$)

⇒ the vector of fitted values $\underline{\hat{y}}$ and the vector of residuals \underline{e} are orthogonal: $\underline{e}^T \underline{\hat{y}} = 0$

the vector \underline{e} and \underline{X} are orthogonal: $\underline{e}^T \underline{X} = 0 \Leftrightarrow \underline{X}^T \underline{e} = 0$

$\underline{X}^T (\underline{y} - \underline{X}\underline{\hat{\beta}}) = 0 \rightarrow$ the normal equation

SUM OF SQUARES DECOMPOSITION

the least squares estimate decomposes the response vector into two orthogonal components

$$\underline{y} = \underline{\hat{\mu}} + \underline{e} = \underline{\hat{y}} + \underline{e} = \underline{P}\underline{y} + (\underline{I}_n - \underline{P})\underline{y}$$

thanks to the orthogonality between \underline{e} and $\underline{\hat{y}}$ we can write

$$\begin{aligned} \underline{y}^T \underline{y} &= \underline{y}^T (\underline{P} + \underline{I}_n - \underline{P}) \underline{y} = \\ &= \underline{y}^T \underline{P} \underline{y} + \underline{y}^T (\underline{I}_n - \underline{P}) \underline{y} = \\ &= \underline{y}^T \underline{P}^T \underline{P} \underline{y} + \underline{y}^T (\underline{I}_n - \underline{P})^T (\underline{I}_n - \underline{P}) \underline{y} = \quad (P \text{ and } (I_n - P) \text{ are symmetric and idempotent}) \\ &= \underline{\hat{y}}^T \underline{\hat{y}} + \underline{e}^T \underline{e} \quad \Rightarrow \underline{P} = \underline{P}^2, \underline{P}\underline{P} = \underline{P}^T \underline{P} \end{aligned}$$

$$\Rightarrow \underline{y}^T \underline{y} = \underline{\hat{y}}^T \underline{\hat{y}} + \underline{e}^T \underline{e}$$

or, equivalently $\|\underline{y}\|^2 = \|\underline{\hat{y}}\|^2 + \|\underline{e}\|^2$

Consider a model that includes the intercept: $\underline{X} = [\underline{1}_n \ \underline{x}_1 \ \dots \ \underline{x}_p]$, then $\underline{1}_n \in \mathcal{C}(X)$

and for the normal equations: $\underline{1}_n^T \underline{e} = 0 \Rightarrow \sum_{i=1}^n e_i = 0$

$$\begin{aligned} \text{moreover, } \underline{1}_n^T \underline{e} &= \underline{1}_n^T (\underline{y} - \underline{\hat{y}}) = \underline{1}_n^T \underline{y} - \underline{1}_n^T \underline{\hat{y}} = 0 \\ &= n\bar{y} - n\bar{\hat{y}} \Rightarrow \bar{y} = \bar{\hat{y}} \end{aligned}$$

Let's consider

$$\begin{aligned} (\underline{y} - \underline{1}\bar{y})^T (\underline{y} - \underline{1}\bar{y}) &= (\underline{y} - \underline{1}\bar{y})^T (\underline{I}_n - \underline{P} + \underline{P}) (\underline{y} - \underline{1}\bar{y}) \\ &= \underbrace{(\underline{y} - \underline{1}\bar{y})^T (\underline{I}_n - \underline{P}) (\underline{y} - \underline{1}\bar{y})}_{(a)} + \underbrace{(\underline{y} - \underline{1}\bar{y})^T \underline{P} (\underline{y} - \underline{1}\bar{y})}_{(b)} \end{aligned}$$

$$\begin{aligned} (a) \ (\underline{y} - \underline{1}\bar{y})^T (\underline{I}_n - \underline{P}) (\underline{y} - \underline{1}\bar{y}) &= \underline{y}^T (\underline{I}_n - \underline{P}) \underline{y} - \underline{y}^T (\underline{I}_n - \underline{P}) \underline{1}\bar{y} - \bar{y} \underline{1}^T (\underline{I}_n - \underline{P}) \underline{y} + \bar{y} \underline{1}^T (\underline{I}_n - \underline{P}) \underline{1}\bar{y} \\ &= \underline{e}^T \underline{e} - \underline{e}^T \underline{1}\bar{y} - \bar{y} \underline{1}^T \underline{e} + \bar{y} \underline{1}^T \underline{1}\bar{y} - \bar{y} \underline{1}^T \underline{P} \underline{1}\bar{y} \end{aligned}$$

$$\begin{aligned} (b) \ (\underline{y} - \underline{1}\bar{y})^T \underline{P} (\underline{y} - \underline{1}\bar{y}) &= \underline{y}^T \underline{P} \underline{y} - \underline{y}^T \underline{P} \underline{1}\bar{y} - \bar{y} \underline{1}^T \underline{P} \underline{y} + \bar{y} \underline{1}^T \underline{P} \underline{1}\bar{y} \\ &= \underline{\hat{y}}^T \underline{\hat{y}} - \underline{\hat{y}}^T \underline{1}\bar{y} - \bar{y} \underline{1}^T \underline{\hat{y}} + \bar{y} \underline{1}^T \underline{P} \underline{1}\bar{y} \end{aligned}$$

$$\begin{aligned} \Rightarrow (\underline{y} - \underline{1}\bar{y})^T (\underline{y} - \underline{1}\bar{y}) &= \underline{e}^T \underline{e} + \bar{y} \underline{1}^T \underline{1}\bar{y} - \bar{y} \underline{1}^T \underline{P} \underline{1}\bar{y} + \underline{\hat{y}}^T \underline{\hat{y}} - \underline{\hat{y}}^T \underline{1}\bar{y} - \bar{y} \underline{1}^T \underline{\hat{y}} + \bar{y} \underline{1}^T \underline{P} \underline{1}\bar{y} \\ &= \underline{e}^T \underline{e} + (\underline{\hat{y}} - \underline{1}\bar{y})^T (\underline{\hat{y}} - \underline{1}\bar{y}) \end{aligned}$$

$$\Rightarrow (\underline{y} - \underline{1}\bar{y})^T (\underline{y} - \underline{1}\bar{y}) = (\underline{\hat{y}} - \underline{1}\bar{y})^T (\underline{\hat{y}} - \underline{1}\bar{y}) + (\underline{y} - \underline{\hat{y}})^T (\underline{y} - \underline{\hat{y}})$$

$$\sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 + \sum_{i=1}^n (y_i - \hat{y}_i)^2 \Rightarrow \text{DEVIANCE decomposition}$$

SST SSR SSE

This is the same decomposition that we found in the simple lin.

Also in this case, we can define the coefficient of determination $R^2 = \frac{SSR}{SST}$.

Its interpretation does not change.

Notice that we derived the decomposition using the fact that $\underline{1}_n^T \underline{e} = 0$ i.e. $\bar{e} = 0$, which

holds only if $\underline{1}_n \in \mathcal{C}(X)$ ⇒ if the model includes the intercept.

If we don't include the intercept, in general $SST \neq SSR + SSE$, and R^2 loses its interpretation.