

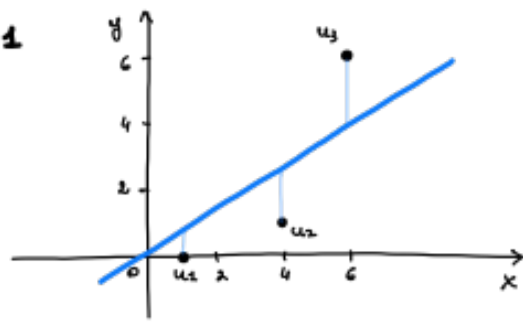
GEOMETRIC INTERPRETATION

Est'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 .

	x_i	y_i
u_1	1	0
u_2	4	1
u_3	6	6

$n=3$ $p=1$

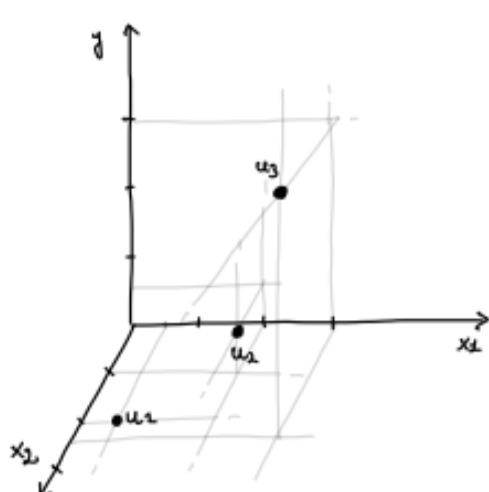


Our problem up to now was:

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}

	x_{i1}	x_{i2}	y_i
u_1	1	4	0
u_2	4	3	1
u_3	6	5	6



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

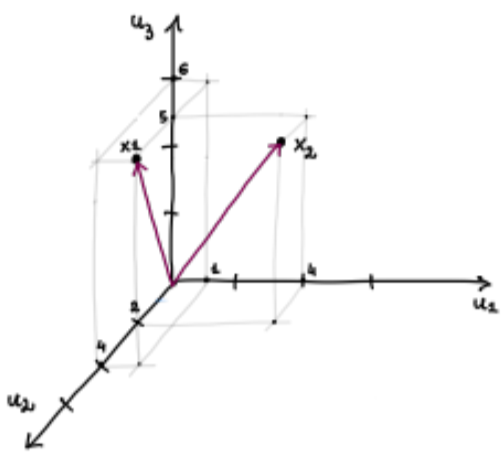
where the coordinates of each point are the values assumed by the p covariates and the response.

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

where $X = [x_1 \ x_2 \ \dots \ x_p]$, and the columns are p n -dimensional vectors

\rightarrow 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 x_1, \dots, x_p
 $C(X) = \{ \underline{\mu} \in \mathbb{R}^n : \underline{\mu} = X\underline{\beta} = \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p, \underline{\beta} \in \mathbb{R}^p \}$

In particular, $C(X)$ is the SUBSPACE of \mathbb{R}^n generated by (x_1, \dots, x_p) .

\hookrightarrow p linearly indep. vectors
 $\Rightarrow C(X)$ has dimension p

In our example, the 2 vectors identify a plane (2-dim space)
 \rightarrow any linear combination of x_1 and x_2 will lie on this plane

If we call $X = [x_1 \ x_2]$, $(n \times p) = (3 \times 2)$ matrix

$C(X) = \beta_1 x_1 + \beta_2 x_2$ the column space of X

$C(X)$ is a subspace of \mathbb{R}^3 of dimension 2

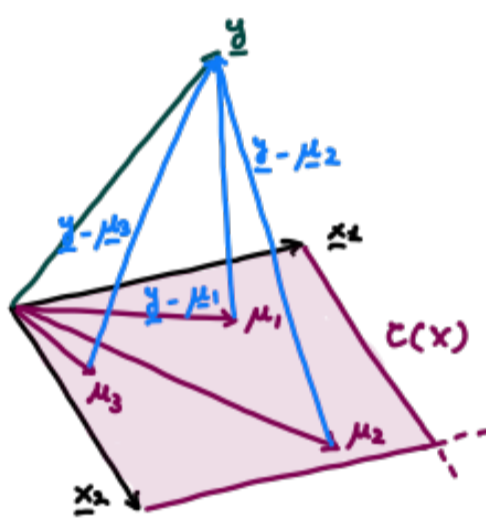
\Rightarrow any $\underline{\mu} = \beta_1 x_1 + \beta_2 x_2$ will lie on $C(X)$

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

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

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

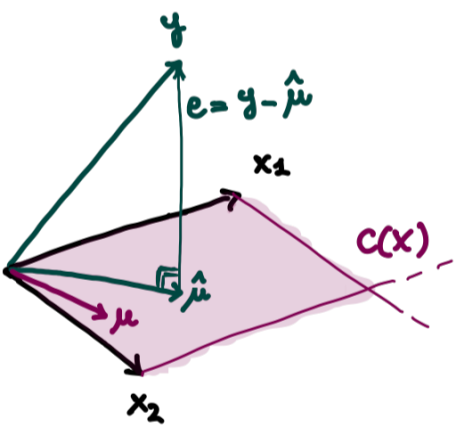
$\underline{y} - X\underline{\beta}$ is the difference between the response and that vector of $C(X)$.



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

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

\Rightarrow we want $\underline{y} - X\underline{\hat{\beta}}$ to be orthogonal to $C(X)$ (hence $\underline{y} - X\underline{\hat{\beta}}$ is orthogonal to the columns x_1, \dots, x_p of X)



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

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

$\Rightarrow \underline{y} - X\underline{\hat{\beta}} \perp x_j$ for all $j=1, \dots, p$

* orthogonality: $\begin{cases} (\underline{y} - X\underline{\hat{\beta}})^T x_1 = 0 \\ \vdots \\ (\underline{y} - X\underline{\hat{\beta}})^T x_p = 0 \end{cases}$
 \downarrow
 normal equations

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

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

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

$(I_n - P)$ is also a projection matrix of rank $n-p$ (it projects on the space $\perp C(X)$)

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

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

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

SUM OF SQUARES DECOMPOSITION

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

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

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

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

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

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

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

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

$$\text{moreover, } 1_n^T \underline{e} = 1_n^T (\underline{y} - \hat{\underline{y}}) = \sum_{i=1}^n y_i - \sum_{i=1}^n \hat{y}_i = 0 \Rightarrow \sum_{i=1}^n y_i = \sum_{i=1}^n \hat{y}_i$$

$$\underline{y} = \hat{\underline{y}} + \underline{e} \Rightarrow \underline{y} - 1_n \bar{y} = \hat{\underline{y}} - 1_n \bar{y} + \underline{e}$$

$$\Rightarrow \|\underline{y} - 1_n \bar{y}\|^2 = \|\hat{\underline{y}} - 1_n \bar{y}\|^2 + \|\underline{e}\|^2$$

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

SST SSR SSE

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

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

Its interpretation does not change.