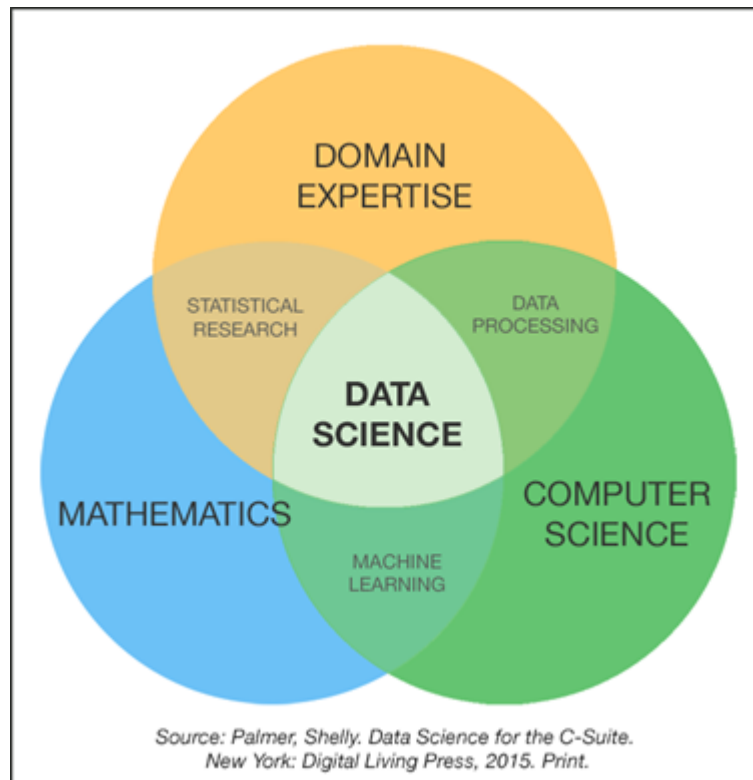




# **CIENCIA DE DATOS: APRENDE LOS FUNDAMENTOS DE MANERA PRÁCTICA**



## **SESION 03 APRENDIZAJE SUPERVISADO EDA II**

**Juan Antonio Chipoco Vidal**  
jchipoco@gmail.com



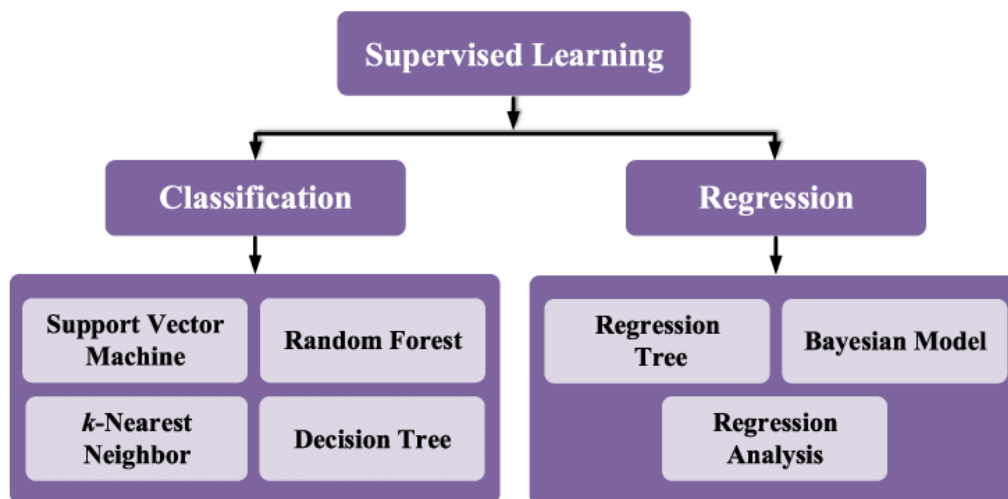
# ÍNDICE

OBJETIVO .....	4
MEDIDAS DE DISPERSION .....	5
MEDIDAS DE DISPERSION: CORRELACION .....	6
MEDIDAS DE DISPERSION: CORRELACION .....	7
MEDIDAS DE DISPERSION: COEFICIENTE DE CORRELACION .....	8
MEDIDAS DE DISPERSION: COEFICIENTE DE DETERMINACION $R^2$ .....	9
MEDIDAS DE DISPERSION: COEFICIENTE DE DETERMINACION $R^2$ .....	10
MEDIDAS DE DISPERSION: COEFICIENTE DE DETERMINACION $R^2$ .....	11
COVARIANZA Y EL COEFICIENTE DE CORRELACION .....	12
COVARIANZA Y EL COEFICIENTE DE CORRELACION .....	13

## Objetivo

El objetivo de esta sesión es profundizar en el análisis de la correlación lineal de dos variables, la cual cuantifica que tan relacionadas están las mismas. Esta técnica está estrechamente relacionada con la **regresión lineal** la cual da lugar a una ecuación que describe dicha relación en términos matemáticos.

En la práctica de esta sesión, continuación de la práctica de la sesión anterior, finalizaremos el análisis exploratorio de datos para poder ya aplicar diversos algoritmos de **clasificación** para obtener la variable objetivo buscada, en este caso la supervivencia o no de un pasajero del Titanic.





## Medidas de Dispersion

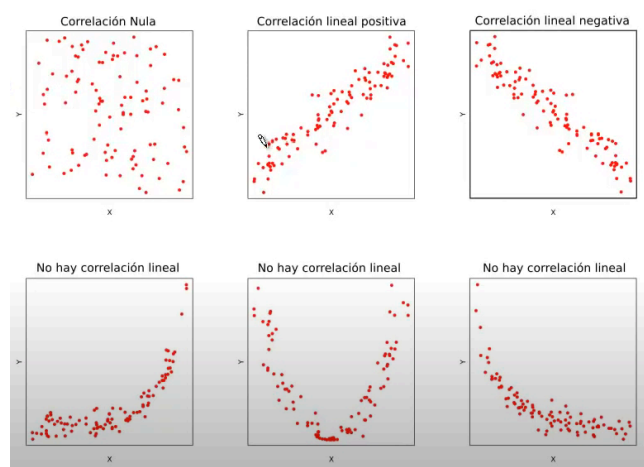
Las medidas de dispersión, se utiliza para describir la variabilidad en una muestra o población. Por lo general, se usa junto con una medida de tendencia central, como la media o la mediana, para proporcionar una descripción general de un conjunto de datos.

## Medidas de Dispersion: Correlacion

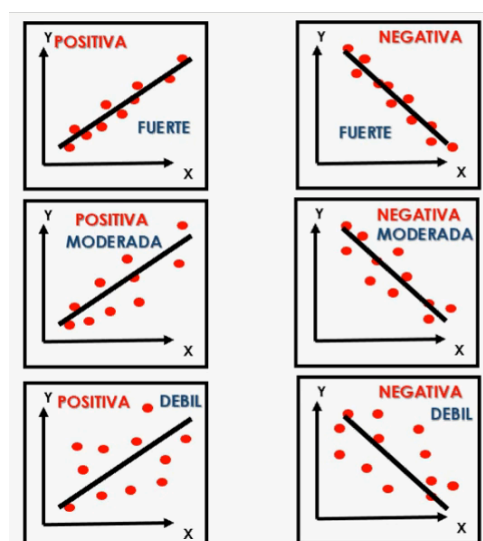
La correlacion sirve para medir la relacion que existe entre dos o mas variables.

La correlacion contesta preguntas como las siguientes:

La practica de algun deporte esta relacionada con una vida mas longeva?  
 Existe una relacion entre la cantidad de carne ingerida diariamente y el cancer?  
 Mayor estudio implica mejores notas en un examen?



Si la correlacion es lineal su direccion puede ser positiva o negativa. Su fuerza varia entre perfecta y nula.



## Medidas de Dispersion: Correlacion



## Medidas de Dispersion: Coeficiente de Correlacion

Para cuantificar las relaciones anteriores tenemos el Coeficiente de Correlacion al cual se le asignara un valor entre -1 y 1.

Este coeficiente nos da una medida de la fuerza y el sentido de una relacion lineal entre variables cuantitativas.

Cuando el signo es positivo la asociacion lineal es positiva lo que implica que cuando el valor de una variable x aumenta tambien aumenta el valor de la otra variable y.

Cuando el signo es negativo la asociacion lineal es negativa lo que implica que cuando el valor de una variable x aumenta el valor de la otra variable y disminuye.

$\pm 0.96$  ,  $\pm 1.0$  PERFECTA

$\pm 0.85$  ,  $\pm 0.95$  FUERTE

$\pm 0.70$  ,  $\pm 0.84$  SIGNIFICATIVA

$\pm 0.50$  ,  $\pm 0.69$  MODERADA

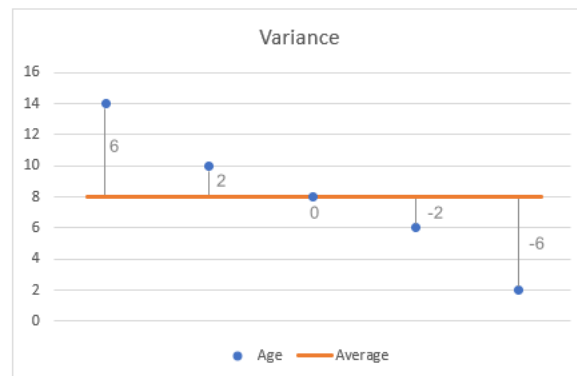
$\pm 0.20$  ,  $\pm 0.49$  DÉBIL

$\pm 0.10$  ,  $\pm 0.19$  MUY DÉBIL

$\pm 0.09$  ,  $\pm 0.0$  NULA

## Medidas de Dispersion: Coeficiente de Determinacion $R^2$

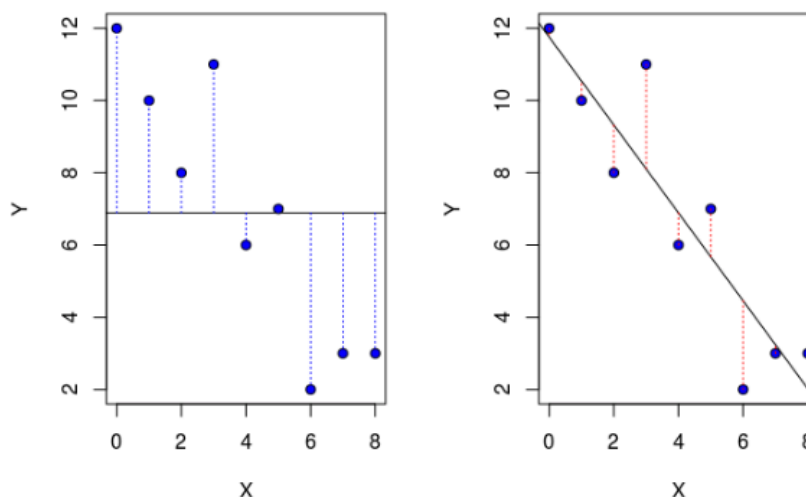
Recordemos que la varianza es la medida de la variabilidad de un conjunto de datos que indica hasta qué punto se distribuyen los diferentes valores. Matemáticamente, se define como la suma de los cuadrados de las diferencias entre una variable y su media, dividido entre el numero de datos.



$$Mean = \frac{14 + 10 + 8 + 6 + 2}{5} = 8$$

$$Variance = \frac{6^2 + 2^2 + 0^2 + (-2)^2 + (-6)^2}{5} = 16$$

El 16 nos da una idea de la dispersion de los datos. Un valor de 0 indica que no hay variabilidad, mayor el valor, mayor la dispersion de los datos.



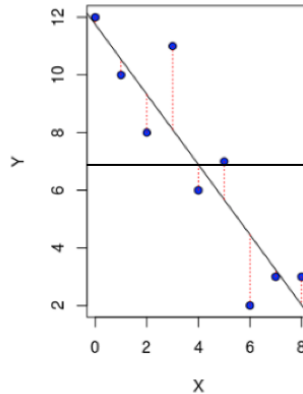
$$Var(mean) = \sum (y_i - \bar{y})^2$$

$$Var(line) = \sum (y_i - (mx_i + b))^2$$

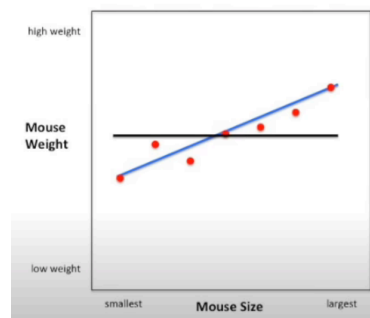
En la grafica anterior trataremos de averiguar que tan bien se ajusta la recta del lado derecho al conjunto de datos. ¿Cual es la bondad del ajuste?.



## Medidas de Dispersion: Coeficiente de Determinacion $R^2$



¿Es este ajuste mejor que el ajuste con la media? Si es así, ¿Qué tan mejor es? ¿Cómo cuantificamos esta diferencia?



Variables correlacionadas

La suma total de cuadrados de los residuos de la imagen anterior  $\text{Var}(\text{line})$  representa la variación del modelo ajustado, o variación no explicada por el modelo (recta de regresión).

Supongamos que  $\text{Var}(\text{mean}) = 32$  y  $\text{Var}(\text{line}) = 6$

Por lo que  $\text{Var}(\text{line})/\text{Var}(\text{mean})$  nos indicara que porcentaje de la variación total en y (peso del ratón) no está explicada por la variación en x (tamaño del ratón).

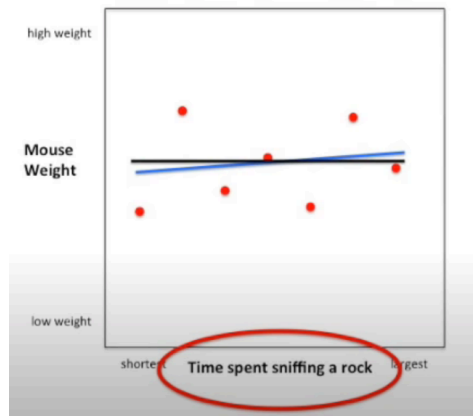
$$\text{Var}(\text{line})/\text{Var}(\text{mean}) = 6/32 = 19\%$$

Así pues para saber que porcentaje de la variación total en y (peso del ratón) está explicada por la variación en x (tamaño del ratón) usamos  $1 - \text{Var}(\text{line})/\text{Var}(\text{mean}) = 81\%$

En otras palabras la relación entre las dos variables explica el 81% de la variación de los datos. Esta relación es significativa.

A este último resultado se le conoce como coeficiente de determinación  $R^2$

## Medidas de Dispersion: Coeficiente de Determinacion $R^2$



Variables no correlacionadas

$$\text{Var}(\text{mean}) = 32 \text{ y } \text{Var}(\text{line}) = 30$$

$\text{Var}(\text{line})/\text{Var}(\text{mean})$  nos indicara que porcentaje de la variacion total en y (pero del raton) no esta explicada por la variacion en x (tiempo oliendo una roca).

$$\text{Var}(\text{line})/\text{Var}(\text{mean}) = 30/32 = 94\%$$

Asi pues para saber que porcentaje de la variacion total en y (peso del raton) esta explicada por la variacion en x (tiempo oliendo una roca) usamos  $1 - \text{Var}(\text{line})/\text{Var}(\text{mean}) = 6\%$

En otras palabras la relacion entre las dos variables explica el 6% de la variacion de los datos. Esta relacion no es significativa.

Si el coeficiente de correlacion  $R = 0.9$  entonces el coeficiente de determinacion  $R^2 = 0.81$ , la relacion entre las dos variables explica el 81% de la variacion de los datos.

$R^2$  es mas facil de interpretar, por ejemplo que tan mejor es  $R = 0.7$  que  $R = 0.5$

$$R^2 = 0.7^2 = 0.49$$

$$R^2 = 0.5^2 = 0.25$$

Con  $R^2$  es facil ver que la primera correlacion es el doble mejor que la segunda correlacion.

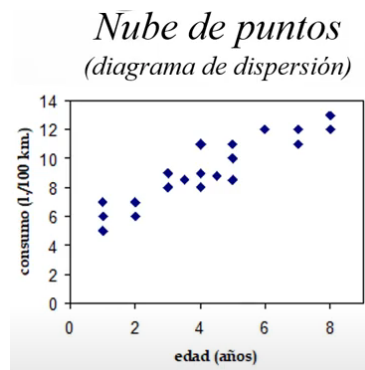
## Covarianza y el coeficiente de correlacion

En el siguiente grafico tenemos informacion de la edad de 20 automoviles asi como el consumo de gasolina en litros por cada 100 km según la edad del automovil.

*Edad y consumo de gasolina de 20 automóviles*

edad auto (años)	Consumo (l/100km)
7	11
5	10
3	8
2	7
7	12
8	12
5	11
4	11
4	8
8	13
1	7
6	12
1	6
3	9
2	6
3,5	8,5
1	5
4,5	8,75
5	8,5
4	9

Al graficar el diagrama de dispersion podemos ver que hay una relacion lineal positiva o directa entre ambas variables. Nos da informacion sobre la covariacion (variacion conjunta) y sus características, si es lineal, su signo y su intensidad.



Ahora calculemos la covarianza y el coeficiente de correlacion para estos datos:

edad auto	consumo				
$x_i$	$y_i$	$x_i y_i$	$x_i^2$	$y_i^2$	
7	11	77	49	121	
5	10	50	25	100	
3	8	24	9	64	
2	7	14	4	49	
7	12	84	49	144	
8	12	96	64	144	
5	11	55	25	121	
4	11	44	16	121	
4	8	32	16	64	
8	13	104	64	169	
1	7	7	1	49	
6	12	72	36	144	
1	6	6	1	36	
3	9	27	9	81	
2	6	12	4	36	
3,5	8,5	29,75	12,25	72,25	
1	5	5	1	25	
4,5	8,5	39,375	20,25	76,5625	
5	8,5	42,5	25	72,25	
4	9	36	16	81	
Totales	84	182,75	856,625	446,5	1770,0625

## Covarianza y el coeficiente de correlacion

De la formula de covarianza tenemos:

$$\sigma_{XY} = \sum_{i=1}^N \frac{(X_i - \bar{X})(Y_i - \bar{Y})}{N}$$

$$\sigma_{XY} = 4.4548$$

$$S_{XY} \begin{cases} > 0 \Rightarrow \text{covariación lineal directa (positiva)} \\ < 0 \Rightarrow \text{covariación lineal inversa (negativa)} \\ = 0 \Rightarrow \text{no hay covariación lineal} \end{cases}$$

De la formula del coeficiente de correlacion:

$$\rho_{XY} = \sigma_{XY} / (\sigma_X \sigma_Y)$$

$$\rho_{XY} = 0.9194$$

Se trata entonces de una relación directa o positiva y muy fuerte.