Concepts	Description
	Correlation and Regression
Covariance	Covariance: measure the linear relationship between 2 random variables
	$Cov_{X,Y} = \frac{\sum_{i=1}^{n} (X_i - \bar{X}) \times (Y_i - \bar{Y})}{n-1}$
	n-1
Correlation coefficient	Correlation coefficient: strength of the linear relationship
	$Cov_{X,Y}$
	$r_{XY} = rac{Cov_{X,Y}}{\sigma_X imes \sigma_Y}$
	$-1 \le r_{XY} \le 1$
Scatter plot	Scatter plot : collection of point on the graph , each represents the value of 2 variables (X and Y)
	- Upward scatter plot : positive correlation
	- Downward scatter plot : negative correlation
Limitation to correlation analysis	1. Outliers: extreme values for sample observations → statistical evidence that significant relationship exist when there is none, or
	→ no relationship when there is
	2. Spurious Correlation : may appear to have a relationship when there is none
	3. Correlation only measure linear relationship, but not non-linear relationship
Test the correlation between 2	Methodology: use t-test to test whether the correlation between 2 variables = 0
variables	H_o : $\rho = 0$; H_a : $\rho \neq 0$
	$t = \frac{r \times \sqrt{n-2}}{\sqrt{1-r^2}}$ $df = n-2$
	$\sqrt{1-r^2}$ $df=n=2$
	$\omega - \kappa$
	Reject H_0 if $t > +t_{critical}$ or $t < -t_{critical}$
Dependent / Independent variables	Dependent variables: variable whose variation is explained by independent variables Independent variables: variable is used to explain the variation of dependent variables
variables	independent variables . Variable is used to explain the variation of dependent variables
Assumptions of linear regression	Independent variables: variable is used to explain the variation of dependent variables 1. Linear regression exists between dependent and independent variables 2. Independent variable: uncorrelated with residuals 3. Expected value of residual term $E(\varepsilon) = 0$ 4. Variance of the residual term is constant for all observations $E(\varepsilon_i^{\;2}) = \sigma_\varepsilon^{\;2}$ 5. Residual term is independently distributed (residual for observation A is not correlated with residual for custs varion B) 6. Residual term is normally distributed $Y_i = b_0 + b_1 \times X_i + \varepsilon_i$ In which: $Y = dependent \ variable$ $X = independent \ variable$ $X = independent \ variable$ $b_0 = reares sign \ short or fixed in the bound of the properties of the prope$
	2. Independent variable : uncorrelated with residuals
	3. Expected value of residual term E(ε) = 0
	4. Variance of the residual term is constant for all observations $E(\varepsilon_t^2) = \sigma_g^2$
	5. Residual term is independently distributed (residual for observation A is not correlated with religible term is open-publication by the design of the control of the con
Linear regression model	6. Residual term is normany distributed V. – b. + b. × Y. + c.
	$I_i = \nu_0 + \nu_1 \wedge \lambda_i + \varepsilon_i$
	In which:
	Y = depedent variable
	A = maepenaent partable - h = repression principal de la
	$b_0 = regress $ or slave for futuent
	DAG DAG
Linear equation for regression line	$\hat{V} = \hat{v}_1 + \hat{v}_2 \times \hat{v}_3$
'	1 - 20 - 61 - 52
	III WIIICII .
	$\widehat{b_1} = \text{estimated slope coefficient} = \frac{\text{Cov}_{XY}}{\sigma_v^2}$
	σ_{χ^2}
	$\widehat{b_0} = \text{estimated intercept term} = \overline{Y} - \widehat{b_1} \times \overline{X}$
	$\vec{Y} = mean \ of \ Y$
	$\bar{X} = mean of X$
Confidence interval for regression	Confidence interval of regression slope coefficient:
slope coefficient (Range of b1)	$\overline{b_1} \pm (t_c imes s_{\overline{b_1}})$
	$ \begin{array}{c} \text{or} \\ \overline{b_1} - (t_c \times s_{\overline{b_1}}) < b_1 < \overline{b_1} + (t_c \times s_{\overline{b_1}}) \end{array} $
	$v_1 - (\iota_c \times s_{\overline{b_i}}) < v_1 < v_1 + (\overline{\iota}_c \times s_{\overline{b_i}})$
	In which:
	t_c = critical 2 – tailed t value for the selected confidence level , with $df = n - 2$
	$s_{\widehat{b_1}} = standard\ of\ error$
Test hypothesis that slope	From the confidence interval for slope coefficient → use t-test to test the hypothesis that slope coefficient = hypothesised value
coefficient = hypothesized value	
	$t_{b_1} = rac{\widehat{b_1} - b_1}{s_{\widehat{b_1}}}$
	"1 S _b -
	Reject H_0 if $t > +t_{critical}$ or $t < -t_{critical}$