



# Statistics – Bachelor in Computer Engineering - DEIOAC

Sample	<b>Parameters</b>
	Sample

Average:

Median or C<sub>2</sub>:

$$\overline{x} = \frac{\sum_{i=1}^{N} x_i}{N}$$

If N is odd  $\Rightarrow$  value at the position (N+1)/2

If N is even  $\Rightarrow$  average of values at positions N/2 and (N/2+1)

### **Quartiles:**

<b>C</b> <sub>1</sub> (Lower quartile) is the first quartile if:
No. data $\leq C_1$ is higher or equal to N/4

 $C_3$  (*Upper quartile*) is the third quartile if: No. data  $\leq C_3$  is higher or equal to 3N/4

No. data  $\geq C_1$  is higher or equal to 3N/4

No. data  $\geq C_3$  is higher or equal to N/4

Variance: 
$$S^2 = \sum_{i=1}^{N} \frac{(x_i - \overline{X})^2}{N-1}$$

Standard deviation:  $S = \sqrt{S^2}$ 

Range:

Interquartile range:

Coefficient of variation:

$$R = X_{max} - X_{min}$$

$$RI = C_3 - C_1$$

$$CV = \frac{S}{\overline{X}}$$

**Skewness coefficient:** 

Standardized skewness coefficient (SSC):

CA = 
$$\frac{\sum_{i=1}^{N} (x_i - \overline{x})^3 / (N - 1)}{S^3}$$

If SSC < -2  $\Rightarrow$  negative skew If SSC  $\in$  [-2, 2]  $\Rightarrow$  symmetric distribution

If SSC  $\in$  [-2, 2]  $\Rightarrow$  symmetric distribution If SSC > 2  $\Rightarrow$  positive skew

Coefficient of kurtosis:

Standardized kurtosis coefficient (SKC):

 $CK = \frac{\sum_{i=1}^{N} (x_i - \overline{x})^4 / N - 1}{S^4} - 3$ 

If SKC < -2  $\Rightarrow$  platykurtic data If SKC  $\in$  [-2, 2]  $\Rightarrow$  mesokurtic data ("normal") If SKC > 2  $\Rightarrow$  leptokurtic data

Covariance:

Correlation coefficient:

$$cov_{xy} = \frac{\sum_{i=1}^{N} (x_i - \overline{x})(y_i - \overline{y})}{N-1}$$

 $r_{xy} = \frac{\text{cov}_{xy}}{\text{s}_x \cdot \text{s}_y}$   $r_{xy} \in [-1,+1]$ 

# Probability Properties: $0 \le P(A) \le 1$ If **A** and **B** are **exclusive** $\Rightarrow$ $P(A \cup B) = P(A) + P(B)$ $\Rightarrow$ $P(A \cap B) = 0$ Complementary event: $P(\overline{A}) = 1 - P(A)$ $P(\emptyset) = 0$ being $\phi$ the empty set Rule of Laplace: De Morgan's laws:

$$P(A) = \frac{\text{favorable cases}}{\text{possible cases}}$$

$$\overline{A \cup B} = \overline{A} \cap \overline{B}$$

 $\overline{A \cap B} = \overline{A} \bigcup \overline{B}$ 





## Statistics – Bachelor in Computer Engineering - DEIOAC

#### Union of events:

 $P(A \cup B)=P(A)+P(B)-P(A \cap B)$ 

 $P(A \cup B \cup C) = P(A) + P(B) + P(C) - P(A \cap B) - P(A \cap C) - P(B \cap C) + P(A \cap B \cap C)$ 

In general:

$$P(A_{_{1}} \cup ... \cup A_{_{n}}) = \sum (P(A_{_{i}})) - \sum (P(A_{_{i}} \cap A_{_{j}})) + \sum (P(A_{_{i}} \cap A_{_{j}} \cap A_{_{k}})) + ... + (-1)^{n+1} \left(\sum (P(A_{_{1}} ... A_{_{n}})) + ... + (-1)^{n+1} \left(\sum (P(A_$$

Conditional probability:

$$P(A/B) = \frac{P(A \cap B)}{P(B)}$$

Intersection of events:

$$P(A \cap B) = P(A).P(B/A)$$
  
 $P(A \cap B) = P(B).P(A/B)$ 

If **A** and **B** are independent ⇒

$$P(A \cap B) = P(A).P(B)$$

Total probability theorem:

$$P(B) = \sum_{j=1}^{n} P(A_{j}) P(B/A_{j}) = P(A_{1}) P(B/A_{1}) + ... + P(A_{n}) P(B/A_{n})$$

$$P(A_{j}/B) = \frac{P(A_{j} \cap B)}{P(B)} = \frac{P(A_{j}) P(B/A_{j})}{\sum_{j=1}^{n} P(A_{j}) P(B/A_{j})}$$

Bayes' Theorem

$$P(A_i/B) = \frac{P(A_i \cap B)}{P(B)} = \frac{P(A_i)P(B/A_i)}{\sum_{j=1}^{n} P(A_j)P(B/A_j)}$$

Probability distributions				
Distribution function: $F(x) = P(X \le x)$	<b>Property:</b> $P(a < X \le b) = F(b) - F(a)$			
Discrete random variables	Continuous random variables			
<b>Probability function:</b> $P(X = x_i)$	Density function: $f(x) = \frac{dF(x)}{dx}$			

## Mathematical expectation and population parameters

Average: m=E(X) Discrete distributions:  $E(X) = \sum_{i=1}^{n} x_i \cdot P(X = x_i)$  Continuous distr.:  $E(X) = \int x \cdot f(x) \cdot dx$ 

Variance:  $\sigma^2 = E(X - m)^2$  Discrete distr.:  $\sigma^2 = \sum_{i=1}^n (x_i - m)^2 \cdot P(X = x_i)$  Standard deviation:  $\sigma = \sqrt{\sigma^2}$ 

#### Properties of the average:

If  $Y = a_0 \pm a_1 \cdot X_1 \pm a_2 \cdot X_2 \pm ... \pm a_n \cdot X_n \implies m_Y = a_0 \pm a_1 \cdot m_{X_1} \pm a_2 \cdot m_{X_2} \pm ... \pm a_n \cdot m_n$ 

Particular cases:

$$\begin{array}{ll} \text{If} \quad Y=a+b\cdot X \ \Rightarrow \ m_Y=a+b\cdot m_X \\ \\ \text{If} \quad Y=X_1+X_2 \Rightarrow \ m_Y=m_{X1} \ + \ m_{X2} \\ \end{array} \qquad \begin{array}{ll} \text{If} \quad Y=a-b\cdot X \ \Rightarrow \ m_Y=a-b\cdot m_X \\ \\ \text{If} \quad Y=X_1-X_2 \Rightarrow \ m_Y=m_{X1} \ - \ m_{X2} \\ \end{array}$$

Properties of the variance: If  $Y = a_0 \pm a_1 \cdot X_1 \pm a_2 \cdot X_2 \implies \sigma_Y^2 = a_1^2 \cdot \sigma_{X_1}^2 + a_2^2 \cdot \sigma_{X_2}^2 \pm 2 \cdot a_1 \cdot a_2 \cdot \text{cov}_{X_1 X_2}$ 

Particular cases:

Variation coefficient:  $CV_x = \sigma_x/m_x$  Interquartile range:  $C_3 - C_1$  Range:  $C_{\text{max}} - C_{\text{min}}$ 



 $X_N \approx N(m_{x_N}, \sigma_{x_N}^2)$ 

# **TABLE OF FORMULAS**



# Statistics - Bachelor in Computer Engineering - DEIOAC

Most important distributions						
<b>Binomial:</b> $X \sim B(n, p)$ (X = 0, 1,, n)						
Probability function:		Distribution	function:	Average		Variance:
$P(X = x) = \binom{n}{x} \cdot p^{x} \cdot (1 - p)^{n - x} = \frac{n! \cdot p^{x} \cdot (1 - p)^{n - x}}{x! \cdot (n - x)!}$		$P(X \le x) = \sum_{x_i=0}^{x} P(X=x_i)$		$m_X = E(X) = n.p$		$\sigma_{X}^{2} = \text{n.p.}(1-p)$
	$X_1 \approx B(n_1, p)$					
Properties:	$  \Rightarrow Y$	$= X_1 + + X_N$	$\approx B(n_1 + \dots$	$+n_N,p)$		
	$X_{N} \approx B(n_{N}, p)$					
Poisson: $X \sim Ps(\lambda)$ (X =	0, 1,, ∞)					
Probability function:	$P(X \le x) = \sum_{x=0}^{x} P$	(X-x )	Average:		Variance:	
$P(X = x) = e^{-\lambda} \cdot \frac{\lambda^{2}}{x!}$	$x_i = 0$	( <b>/\-</b> \ <sub>i</sub> ) ⇒	$m_X = E(X)$	$)=\lambda$	$\sigma^2 = \lambda$	
	Abacus of F	Poisson	^ ( )		$O_X - \lambda$	
	$X_1 \approx Ps(\lambda_1)$	7 37 37	D (0	•		
Properties:	⇒ Y	$\mathbf{Y} = \mathbf{X}_1 + \dots + \mathbf{X}_n$	$_{\rm N} \approx {\rm Ps}(\lambda_1 + \lambda_2)$	$\dots + \lambda_{N}$ )		
Uniform: X ~U (a,b)	$\frac{X_{N} \approx Ps(\lambda_{N})}{(a < X < b)}$					
	(a < x < b)		A		Variana	
Density function: $f(x) = \frac{1}{b-a}$ $a \le x \le b$	$P(X < x) - \frac{x - a}{x}$	l - a < x < h	Average:			ce:
$f(x) = \frac{1}{b-a}$ $a \le x \le b$	b-a	1	$m_X = E(X)$	$=\frac{1}{2}$	$\sigma_{X}^{2} = \frac{100}{100}$	12
Exponential: X ~ Exp (α)	$(0 \le X \le \infty)$					
Density function:			Average:		Variand	ce:
$f(x) = \alpha e^{-\alpha x}  x \ge 0$	$P(X \le x) = 1 - e^{-c}$	$x^{-\alpha x}$ $x \ge 0$ $\Rightarrow$ $m_x = E(X)$		$=\frac{1}{\alpha}$	$\sigma_{X}^{2} = \frac{1}{\alpha^{2}}$	
Normal: $X \sim N(m, \sigma)$	$(-\infty < X < \infty)$					
Density function:			Average:		Variand	e:
$f(x) = \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{(x-m)^2}{2\sigma^2}}$	Z ~ N(0;1) stand		$m_X = E(X)$	) = m	$\sigma_{\rm X}^2 = \sigma$	-2
$\sigma\sqrt{2\pi}$	$P(Z \ge z) =$	⇒ Table				
Standard Normal distribution: $Z \sim N(m_z = 0, \sigma_z = 1)$ If $X \approx N(m, \sigma) \Rightarrow \frac{X - m}{\sigma} \approx N(0; 1)$						
$P[N(m,\sigma) > x_i] = P\left[N(0; 1) > \frac{x_i - m}{\sigma}\right] = P[Z > z] \Rightarrow Table$						
$X_1 \approx N(m_{x_1}, \sigma_{x_1}^2)$						
<b>Properties:</b> $\Rightarrow Y = X_1 + + X_N \approx N(m_Y = m_{x_1} + + m_{x_N}; \sigma_Y^2 = \sigma_{X_1}^2 + + \sigma_{X_N}^2)$						
$\mathbf{X} = \mathbf{N}(\dots -2)$						

Standard deviation of Y:  $\sigma_{\rm Y} = \sqrt{\sigma_{\rm X_1}^2 + \cdots + \sigma_{\rm X_N}^2}$ 





# Statistics - Bachelor in Computer Engineering - DEIOAC

บว	rtici	ular	case	c
гα	าแบเ	JIAI	Last	o

If 
$$Y = a + b \cdot X \Rightarrow Y \sim N(m_Y = a + b.m_X; \sigma_Y^2 = b^2.\sigma_X^2)$$
 If  $X \sim N(m_X, \sigma_X)$  and  $Y \sim N(m_Y, \sigma_Y)$  independent  $Z = X \mp Y \sim N \ (m_Z = m_X \mp m_Y, \sigma_Z^2 = \sigma_X^2 + \sigma_Y^2)$ 

If 
$$\mathbf{X} \sim \mathsf{N}(\mathbf{m_x} \ , \ \sigma_{\mathbf{x}} \ ) \Rightarrow \begin{cases} 68,26\% \ \text{of X values} \ \in \ [\text{m-$\sigma$, m+$\sigma}] \\ 95,44\% \ \text{of X values} \ \in \ [\text{m-$2\sigma$, m+$2\sigma}] \\ 99,73\% \ \text{of X values} \ \in \ [\text{m-$3\sigma$, m+$3\sigma}] \end{cases}$$

## **Normal approximations**

#### **Central Limit Theorem:**

$$\begin{split} X_{1} &\approx g_{1}\left(m_{X_{1}}, \sigma_{X_{1}}^{2}\right) \\ &\qquad \dots \\ X_{N} &\approx g_{1}\left(m_{X_{N}}, \sigma_{X_{N}}^{2}\right) \end{split} \Rightarrow Y = X_{1} + \dots + X_{N} \approx N\left(m_{Y} = m_{X_{1}} + \dots + m_{X_{N}}; \sigma_{Y}^{2} = \sigma_{X_{1}}^{2} + \dots + \sigma_{X_{N}}^{2}\right) \end{split}$$

Being:

**g** → any distribution (Binomial, Poisson, etc.)

 $N \to \infty$  (N very high)

Previous teaching material from R. Alcover (DEIOAC - UPV)

This work is under a license Recognizing-Non commercial – sharing under the same license 2.5 Spain of Creative Commons. To see one copy of this license, visit <a href="http://creativecommons.org/licenses/by-nc-sa/2.5/es/">http://creativecommons.org/licenses/by-nc-sa/2.5/es/</a>







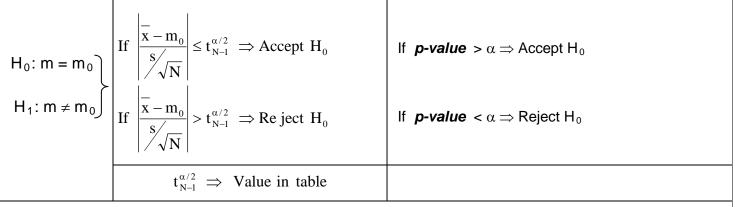
# Statistics - Bachelor in Computer Engineering - DEIOAC

Distributions in the sampling of normal populations	
$X \sim N(m,\sigma)$ and $\overset{-}{x}$ is the average of a sample with size N	$\frac{\overline{x} - m}{\sigma / \sqrt{N}} \sim N(0, 1)$
$X \sim N(m,\sigma)$ and $s^2$ is the variance of a sample with size N	$(N-1)\frac{s^2}{\sigma^2} \sim \chi_{N-1}^2$
$X \sim N(m,\sigma)$ and $x = 0$ are the mean and variance of a sample with size N	$t = \frac{\overline{x} - m}{s / \sqrt{N}} \sim t_{N-1}$
$X_1 \sim N(m_1, \sigma_1)$ , $X_2 \sim N(m_2, \sigma_2)$ independent. $S_1^2$ and $S_2^2$ are the sample variances of $X_1$ and $X_2$ (sizes $N_1$ and $N_2$ )	$\frac{s_1^2/\sigma_1^2}{s_2^2/\sigma_2^2} \sim F_{N_1-1,N_2-1}$

# Inference in normal populations

 $\alpha$ : significance level (type-I risk); **1-** $\alpha$ : confidence level;  $\beta$ : type-II risk; **1-** $\beta$ : power

#### Hypothesis test for the mean (t - test)



#### Confidence interval for the mean (Cl<sub>m</sub>)

$$CI_{m} \Rightarrow \left[ \overline{x} - t_{N-1}^{\alpha/2} \frac{s}{\sqrt{N}}, \ \overline{x} + t_{N-1}^{\alpha/2} \frac{s}{\sqrt{N}} \right] \\ \qquad \qquad t_{N-1}^{\alpha/2} \ \Rightarrow \ \text{Value in t table}$$

## Hypothesis test for the mean using the confidence interval

$$\begin{array}{ll} H_0: \mathbf{m} = \mathbf{m}_0 \\ H_1: \mathbf{m} \neq \mathbf{m}_0 \end{array} \qquad \begin{array}{ll} If \ m_0 \in CI_m \ \Rightarrow Accept \ H_0 \\ If \ m_0 \notin CI_m \ \Rightarrow \mathrm{Re} \ ject \ H_0 \end{array}$$





# Statistics - Bachelor in Computer Engineering - DEIOAC

## Confidence interval for the variance ( $\sigma^2$ )

$$IC_{\sigma^2} \Rightarrow \left\lceil \frac{(N-1)S^2}{g_2}, \frac{(N-1)S^2}{g_1} \right\rceil$$

 $g_1$  and  $g_2 \Rightarrow$  values in Chi<sup>2</sup> table so that:

$$P\!\left(\!\chi_{\scriptscriptstyle N-1}^{\scriptscriptstyle 2}>g_{\scriptscriptstyle 1}\right)\!=\!1\!-\!\frac{\alpha}{2}\quad\text{ and }P\!\left(\!\chi_{\scriptscriptstyle N-1}^{\scriptscriptstyle 2}>g_{\scriptscriptstyle 2}\right)\!=\!\frac{\alpha}{2}$$

## Hypothesis test for the variance ( $\sigma^2$ ) using the confidence interval

$$H_0: \sigma^2 = \sigma_0^2$$

$$H_1: \sigma^2 \neq \sigma_0^2$$

If 
$$\sigma_0^2 \in CI_{\sigma^2} \Rightarrow Accept H_0$$

If 
$$\sigma_0^2 \notin CI_{\sigma^2} \Rightarrow \text{Re ject } H_0$$

# Hypothesis test for the comparison of means

$$\mathbf{H}_0: \mathbf{m}_1 = \mathbf{m}_2$$

$$\mathbf{H}_1: \mathbf{m}_1 \neq \mathbf{m}_2$$

$$\begin{array}{c|c} H_0: m_1 = m_2 \\ H_1: m_1 \neq m_2 \end{array} \qquad \qquad If \ \, \frac{ \begin{vmatrix} - & - \\ \bar{x}_1 - \bar{x}_2 \end{vmatrix} }{s_{(\bar{x}_1 - \bar{x}_2)}} \le t_{N_1 + N_2 - 2}^{\alpha/2} \ \, \Rightarrow Accept \ \, H_0 \\ \end{array}$$

If ***p*-value** > 
$$\alpha \Rightarrow$$
 Accept H<sub>0</sub>

If 
$$\left| \frac{\overline{x_1 - x_2}}{\overline{x_{(\overline{x_1 - x_2})}}} \right| > t_{N_1 + N_2 - 2}^{\alpha/2} \implies \text{Re ject } H_0$$

If **p-value**  $< \alpha \Rightarrow \text{Reject H}_0$ 

 $t_{N_1+N_2-1}^{\alpha/2} \implies \text{Value in table}$ 

 $\alpha$  = Significance level

$$S_{(\bar{X}_1 - \bar{X}_2)} = S \sqrt{\frac{1}{N_1} + \frac{1}{N_2}}$$

$$S = \sqrt{\frac{\left(N_1 - 1\right) S_1^2 + \left(N_2 - 1\right) S_2^2}{N_1 + N_2 - 2}}$$

## Confidence interval for the difference of means $(m_1 - m_2)$

$$IC_{m_1-m_2} \Longrightarrow (\overline{x}_1 - \overline{x}_2) \pm t_{N_1+N_2-2}^{\alpha/2} \ S_{(\overline{x}_1-\overline{x}_2)} \qquad \qquad t_{N_1+N_2-1}^{\alpha/2} \ \Longrightarrow \ Value \ in \ t \ table$$

$$t_{N_1+N_2-l}^{\alpha/2} \Rightarrow \text{Value in t table}$$

# Hypothesis test for the comparison of means using the confidence interval

$$H_0: m_1 = m_2$$

$$H_1: m_1 \neq m_2$$

If 
$$0 \in CI_{m_1-m_2} \Rightarrow Accept H_0$$

If 
$$0 \notin CI_{m_1-m_2} \Rightarrow \text{Re ject } H_0$$

# Confidence interval for the ratio of variances $(\sigma^2_1 / \sigma^2_2)$

$$\mathsf{IC}_{\sigma_1^2/\sigma_2^2} \Rightarrow \left(\frac{\mathsf{S}_1^2}{\mathsf{S}_2^2 \, \mathsf{f}_2}, \frac{\mathsf{S}_1^2}{\mathsf{S}_2^2 \, \mathsf{f}_1}\right)$$

 $f_1$ ,  $f_2 \Rightarrow$  Values in F table so that:

$$P\!\!\left(\!F_{_{(N_{1}-1),(N_{2}-1)}} > f_{_{1}}\right) = 1 - \frac{\alpha}{2} \quad \text{and} \quad P\!\!\left(\!F_{_{(N_{1}-1),(N_{2}-1)}} > f_{_{2}}\right) = \frac{\alpha}{2}$$

#### Hypothesis test for the comparison of variances using the confidence interval

$$H_0: \sigma_1^2 = \sigma_2^2$$

$$H_1: \sigma_1^2 \neq \sigma_2^2$$

$$\begin{array}{ll} \textit{If} & 1 \in CI_{\sigma_{1}^{2}/\sigma_{2}^{2}} \text{ [or if } s_{1}^{2} \big/ s_{2}^{2} < F_{n_{1}-1,n_{2}-1}^{\alpha} \big] \Rightarrow \textit{Accept } H_{0} \\ \textit{If } & 1 \notin CI_{\sigma_{1}^{2}/\sigma_{2}^{2}} \text{ [or if } s_{1}^{2} \big/ s_{2}^{2} > F_{n_{1}-1,n_{2}-1}^{\alpha} \big] \Rightarrow \text{Re ject } H_{0} \end{array} \quad being \ s_{1}^{2} > s_{2}^{2} \end{array}$$





#### Statistics – Bachelor in Computer Engineering - DEIOAC

Analysis of Variance (ANOVA)						
	Nomenclature					
$\mathbf{N}$ = total number of observations $\mathbf{F}_i$ = factor i $\mathbf{F}_j$ = factor j $\mathbf{F}_i$ x $\mathbf{F}_j$ = interaction of $\mathbf{F}_j$ with $\mathbf{F}_i$ .			<ul><li>I = number of levels in factor F<sub>i</sub></li><li>J = number of levels in factor F<sub>j</sub></li></ul>		SS = Sum of Squares df = degrees of freedom	
SS <sub>TOTAL</sub> = total sum of squares			$df_{Tot}$ = total degrees of freedom $df_{Fi}$ = degrees of freedom associated to the SS of factor $F_i$			
$\mathbf{SS}_{Fi} = \mathbf{Sum}$ of squares of factor i $\mathbf{SS}_{Fj} = \mathbf{Sum}$ of squares of factor j $\mathbf{SS}_{FixFj} = \mathbf{SS}$ of the interaction $\mathbf{F}_i \times \mathbf{F}_j$ $\mathbf{SS}_{res} = residual$ sum of squares $\mathbf{df}_{Fi} = \mathbf{deg}$ . freeedom associated to the SS of factor $\mathbf{df}_{Fi} = \mathbf{deg}$ . fr. associated to the SS of the interaction $\mathbf{df}_{Fi} = \mathbf{deg}$ . fr. associated to the SS of the interaction $\mathbf{df}_{Fi} = \mathbf{deg}$ . fr. associated to the SS of factor $\mathbf{df}_{Fi} = \mathbf{deg}$ .			so the SS of factor $\mathbf{F}_{\mathbf{j}}$			
$df_{Tot} \rightarrow (N-1)$	$df_{Fi} \rightarrow (I-1)$	$df_{Fj} \rightarrow (J-1)$	$\mathbf{df}_{\mathbf{F}\mathbf{j}\mathbf{x}\mathbf{F}\mathbf{j}} \to (\mathbf{I}\mathbf{-}1)\mathbf{x}(\mathbf{J}\mathbf{-}1)$	$df_{res} = df_{total} - \left( \sum_{\forall Fact} df_{factors} + \sum_{\forall Fact} df_{int  eractions} \right)$		
		Fundame	ental equation of AN	OVA		
$\mathbf{SS}_{\mathrm{total}} = \sum_{\forall \mathrm{Fact}} \mathbf{SS}_{\mathrm{factors}} + \sum_{\forall \mathrm{int}} \mathbf{SS}_{\mathrm{int eractions}} + \sum_{\mathrm{resid}} \mathbf{SS}_{\mathrm{resid}}$						
		Hypothesi	is test for ANOVA (F	test)		
		$\frac{ctor}{ctor} pprox F_{df_F, df_{resid}}$	MS <sub>Factor</sub> = Mean square associated to the effect of one factor or interaction MS <sub>resid</sub> = mean square of residuals			





#### Statistics – Bachelor in Computer Engineering - DEIOAC

Int	roduction to Multi	iple Linear Regressior	n models		
	Ne	omenclature			
			$X_1=x_{1t}$ ) is the average of the conditional $X_1=x_{1t}$ when $X_1=x_{1t}$ $X_1=x_{1t}$		
<b>N</b> = Total number of observation	$\beta_i = \text{model parame}$ $= \text{Total number of observations}  b_i = \text{estimated mod}$ $I = \text{number of expli}$		SS	SS <sub>total</sub> = SS <sub>expl</sub> + SS <sub>resid</sub>	
\$\$\text{SS}_{total}\$ = total Sum of Squares\$\$\$\$S_{expl}\$ = explained \$\$\$\$S_{resid}\$ = residual \$\$\$\$\$\$	es $df_{tot}$ = total degrees of freed $df_{exp}$ = d.f. associated to SS $df_{res}$ = residual degrees of fr		M MS	$S_{\text{expl}} = SS_{\text{expl}} / df_{\text{exp}}$ $S_{\text{resid}} = SS_{\text{resid}} / df_{\text{res}}$	
	Global significand	e test of the model (AN	OVA)		
$H_0: \beta_1 = \beta_2 = =$	$H_0$ : $\beta_1 = \beta_2 = = \beta_1 = 0$ Coefficient of determine		ermination:	Residual variance:	
$H_1$ : at least one if $\frac{MS_{\mathrm{exp}l}}{MS_{resid}} > F_{I,N-\mathrm{l-}I}^{\alpha} \Longrightarrow reject$		$R^2 = \frac{\exp(-100)}{100}$		$\sigma_{\rm resid}^2 = { m MS}_{ m resid}$	
Sig	nificance test for th	e effect of one X <sub>1</sub> varial	ble ( <i>t</i> test)		
$ \begin{aligned} \mathbf{H}_{0}: \boldsymbol{\beta}_{\mathrm{i}} &= 0 \\ \mathbf{H}_{1}: \boldsymbol{\beta}_{\mathrm{i}} &\neq 0 \end{aligned} \qquad \mathbf{t}_{\mathrm{statistic}} &= \frac{\mathbf{b}_{\mathrm{i}}}{\mathbf{s}_{\mathbf{b}_{\mathrm{i}}}} \approx \mathbf{t}_{\mathrm{N-1-I}} \qquad \begin{vmatrix} if & \left  b_{i} \middle/ s_{b_{i}} \right  < \mathbf{b}_{\mathrm{i}} \\ if & \left  b_{i} \middle/ s_{b_{i}} \right  > \mathbf{b}_{\mathrm{i}} \end{aligned} $		$\left  \langle t_{N-1-I}^{\alpha/2} \Rightarrow accept \ H \right $ $\left  \langle t_{N-1-I}^{\alpha/2} \Rightarrow accept \ H \right $	If <b>p-value</b> > $\alpha \Rightarrow \text{Accept H}_0$ If <b>p-value</b> < $\alpha \Rightarrow \text{Reject H}_0$		
1		Predictions			
(Y/X <sub>1</sub> =x <sub>1t</sub> ,	,X <sub>I</sub> =x <sub>It</sub> ) ~ Norma	al [ m = E(Y/X <sub>1</sub> =x <sub>1t</sub> ,	$(X_1=X_{lt}); \sigma=0$	$\sigma_{res}$ ]	
	Simple Line	ar Regression models	<b>,</b>		
Model $E(Y/X=x_t) = \beta_0 + \beta_1 \cdot x_t$	E(Y/X=x <sub>t</sub> ) is the a	verage of the condition	al distribution	of Y when X=x <sub>t</sub>	
Estimated regression line:		Slope		Intercept	
$Y = a + b \cdot X$	b =	$b = r_{xy} \cdot \frac{s_y}{s_x} = \frac{cov}{s_x^2}$		$a = \overline{Y} - b \cdot \overline{X}$	
Coefficient of determination:	Residual var	Residual variance: Re		Residual:	
$R^2 = (r_{xy})^2 \cdot 100 = \left(\frac{\text{cov}_{xy}}{\text{s}_x \cdot \text{s}_y}\right)^2 \cdot 100$	$S^2_{residual} = S^2$	$S^2_{residual} = S^2_y \cdot (1 - r^2_{xy})$		+ b·x <sub>i</sub> )	