

Corso integrato di informatica, statistica e analisi dei dati sperimentali Esercitazione VII

Un breve richiamo sul test t-Student

Siano A^{exp} ($a_1, a_2 \dots a_n$) e B^{exp} ($b_1, b_2 \dots b_m$) due set di dati i cui valori medi sono:

$$\bar{A} = \frac{1}{n} \sum_i a_i \quad e \quad \bar{B} = \frac{1}{m} \sum_i b_i$$

e le varianze campionarie:

$$\sigma_a^2 = \frac{1}{n-1} \sum_i (a_i - \bar{A})^2 \quad e \quad \sigma_b^2 = \frac{1}{m-1} \sum_i (b_i - \bar{B})^2$$

Il test t- di Student serve per verificare quanto siano significative le differenze osservate tra due gruppi di dati.

Il rapporto:

$$t_{oss} = \frac{\text{differenza tra le medie campionarie}}{\text{errore standard sulla differenza tra le medie campionarie}}$$

cioè:

$$t_{oss} = \frac{\bar{A} - \bar{B}}{\sigma_{\bar{A}-\bar{B}}}$$

è una variabile aleatoria. Nell'ipotesi (H_0) che i due campioni provengano dalla stessa popolazione, cioè nell'ipotesi che il valore atteso e la varianza vera siano gli stessi per ambedue i campione e che le differenze osservate siano dovute al caso, la t_{oss} segue la distribuzione di probabilità di una variabile t-Student: t_ν con $\nu = n + m - 2$ gradi di libertà.

Il test si può applicare in due modi:

- si calcola la probabilità di osservare un valore t_ν maggiore o eguale a t_{oss} . Questa rappresenta la probabilità che, pur essendo valida l'ipotesi H_0 (distribuzioni eguali), abbiamo osservato per caso una certa differenza tra i valori medi. Questa probabilità può essere calcolata utilizzando la funzione Excel

DISTRIB.T(t_{oss} , ν , *code*)

con *code* = 1 calcola la probabilità per una coda della distribuzione:

$$P(t_\nu > t_{oss}) = \int_{t_{oss}}^{\infty} f(t_\nu) dt_\nu$$

Con *code* = 2 calcola la probabilità per una coda della distribuzione:

$$P(|t_\nu| > |t_{oss}|) = \int_{-\infty}^{-t_{oss}} f(t_\nu) dt_\nu + \int_{t_{oss}}^{\infty} f(t_\nu) dt_\nu$$

Si usa *code* = 1 quando, per qualche motivo, si conosce il segno della differenza: $\bar{A} - \bar{B}$

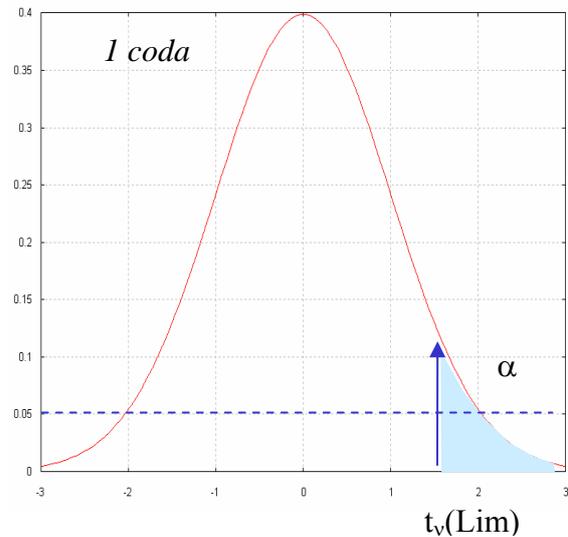
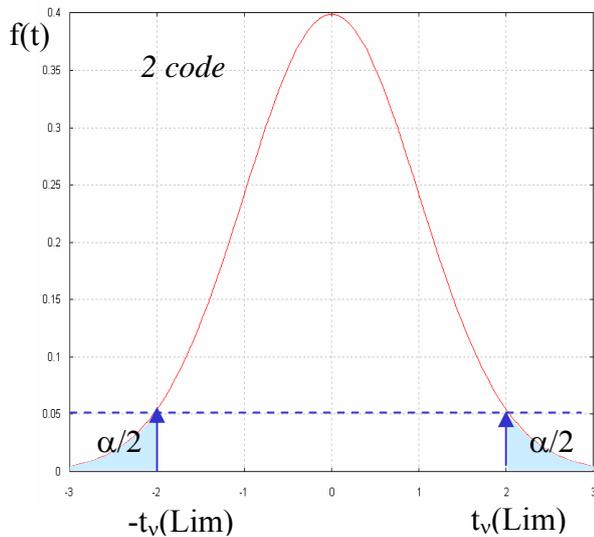
- si definisce un limite di rischio α (confidenza $1-\alpha$) alla probabilità di sbagliare affermando che le distribuzioni sono diverse e si calcola il valore limite $t_\nu(\mathbf{Lim})$ tale che:

$$\alpha = \int_{t_\nu(\mathbf{Lim})}^{\infty} f(t_\nu) dt_\nu$$

si confronta quindi il valore osservato con il valore limite, se $t_{oss} > t_\nu(\mathbf{Lim})$ posso rigettare l'ipotesi con un livello di confidenza maggiore di $1-\alpha$. I valori limite si possono trovare sulle tavole (dispense, internet, etc...) o si possono calcolare utilizzando la funzione Excel:

INV.T(α, ν)

Attenzione: $\text{INV.T}(\alpha, \nu)$ calcola la probabilità associata alla variabile t_ν a due code (quando non conosco il segno della differenza $\bar{A} - \bar{B}$). Per un test ad una coda utilizzare: $\text{INV.T}(2\alpha, \nu)$



Problema: come calcolare $\sigma_{\bar{A}-\bar{B}}$:

a) se i due campioni hanno la stessa numerosità, cioè $n = m$:

$$\sigma_{\bar{A}-\bar{B}} = \sqrt{\sigma_{\bar{A}}^2 + \sigma_{\bar{B}}^2} = \sqrt{\frac{\sigma_a^2}{n} + \frac{\sigma_b^2}{n}} = \sqrt{\frac{\sigma_a^2 + \sigma_b^2}{n}}$$

dove σ_a^2 è la varianza dei dati A e $\sigma_{\bar{A}}$ l'errore standard sulla media \bar{A} (lo stesso vale per i dati B), abbiamo usato il fatto che l'errore standard della media è:

$$\sigma_{\bar{X}} = \frac{\sigma}{\sqrt{n}}$$

b) se i due campioni non hanno la stessa numerosità $n \neq m$:

$$\sigma_{\bar{A}-\bar{B}} = \sqrt{\frac{(n-1)\sigma_a^2 + (m-1)\sigma_b^2}{n+m-2} \left(\frac{1}{n} + \frac{1}{m} \right)}$$

Che si riconduce alla precedente per $m=n$.

Il t-Test può essere usato per stabilire se, dato un set di dati $X^{\text{exp}} (x_1, x_2, \dots, x_n)$ il valor medio \bar{X} è significativamente diverso da un valore determinato μ . In questo caso la t_{oss} :

$$t_{\text{oss}} = \frac{\bar{X} - \mu}{\sigma_{\bar{X}}}$$

Deve essere confrontata con la statistica di una variabile t_ν con $\nu = n-1$.

Esercizio 1) (il file Excel: **student.xls** mostra possibili soluzioni)

a) Il file **T_test1.dat** riporta i risultati di due set di misure ($n_1 = n_2$). Stabilire con un livello di confidenza $1-\alpha=95\%$ se le differenze osservate sono significative

- 1: importare i dati in un foglio Excel,
- 2: calcolare valori medi, varianza, dev.st ed errore standard della media
- 3: si calcola il valore di t_{oss}
- 4: usando la funzione DISTRIB.T si calcola il valore di significatività (due code): se i dati provenissero da una stessa distribuzione la probabilità di osservare un tale valore di t per caso è 0.89%, molto minore del limite del 5%, quindi le differenze sono significative con un livello di confidenza maggiore del 95% (maggiore del 99% in questo caso)
- 4b: usando la funzione INV.T si calcola il valore limite per una probabilità $\alpha = 0.05$ ($\alpha=5\%$) per una variabile t_v . Il valore $t_{oss} = 2.74$ è maggiore del valore limite $t_v(Lim)=2.02$, quindi le differenze sono significative con un livello di confidenza maggiore del 95%

X	X.2		X.1	X.2
1.3	1.56	N	23	23
0.82	1.02	media	0.970	1.170
0.85	1.05	err. st. media	0.052	0.052
0.64	0.84	var.	0.061	0.061
0.6	0.8	dev. st	0.248	0.248
0.98	1.18			
0.84	1.04			
0.66	0.86			
0.65	0.85			
1.24	1.44	t_oss	2.74	
1.02	1.22	Significatività	0.0089	
0.96	1.16			
1.16	1.36			
1.1	1.3			
1	1.2			
1.46	1.66			
1.02	1.22			
0.93	1.13			
0.8	1			
1.07	1.27			
1.45	1.65			
0.97	1.17			
0.73	0.93			

	α	v
	0.05	44
t_n(Lim)		2.02

b) Il file **T_test2.dat** riporta i risultati di due set di misure (diversa numerosità). Stabilire con un livello di confidenza $1-\alpha=95\%$ se le differenze osservate sono significative

- 1: importare i dati in un foglio Excel,
- 2: calcolare valori medi, varianza, dev.st ed errore standard della media
- 3: si calcola il valore di t_{oss} .
- 3b: per non complicare troppo le formule è utile calcolare separatamente il valore:

$$\sigma_{\bar{A}-\bar{B}} = \sqrt{\frac{(n-1)\sigma_a^2 + (m-1)\sigma_b^2}{n+m-2} \left(\frac{1}{n} + \frac{1}{m}\right)}$$
- 4: usando la funzione DISTRIB.T si calcola il valore di significatività (due code): se i dati provenissero da una stessa distribuzione la probabilità di osservare un tale valore di t per caso è 9.5%, maggiore del limite del 5%, quindi le differenze non sono significative con un livello di confidenza maggiore del 95%.
- 4b: usando la funzione INV.T si calcola il valore limite per una probabilità $\alpha = 0.05$ ($\alpha=5\%$) per una variabile t_v . Il valore osservato $t_{oss} = 1.71$ è minore del valore $t_v(Lim)$ quindi le differenze non sono significative con un livello di confidenza maggiore del 95%.

X	X.2		X.1	X.2
11.5	12.44	N	29	18
13.7	14.93	media	13.933	14.989
12.5	13.53	err. st. media	0.351	0.546
11.57	12.45	var.	3.576	5.363
11.12	11.92	dev. st	1.891	2.316
13.77	15.02			
10.95	11.73			
15.5	17.03			
16.25	17.9			
15.23	16.72			
14.7	16.1			
16.08	17.71			
16.35	18.02			
12.05	13.01			
15.39	16.91			
15.07	16.53			
14.79	16.2			
10.89	11.65			
12.56				
11.05				
16.01				
14.95				
15.89				
15.88				
15.99				
13.89				
15.13				
12.29				
12.96				

	s^2
	0.382792

	α	v
	0.05	45
t_n(Lim)		2.01

Nota: le differenze sono significative con un livello di confidenza del: $1-0.0949 = 90.5\%$

Regressione lineare

Una retta di regressione permette di stimare i parametri che caratterizzano una relazione lineare tra due variabili: una variabile indipendente (X) e una variabile dipendente (Y).

Siano X^{exp} ($x_1, x_2 \dots x_n$) e Y^{exp} ($y_1, y_2 \dots y_n$) un set di misure sperimentali tra le quali si ipotizza esista una relazione lineare del tipo:

$$Y^{\text{th}} = aX + b$$

Nell'ipotesi che

- i) **l'errore sulle ascisse (X) sia nullo,**
- ii) **l'errore sulle variabili dipendenti sia costante e indipendente da X**

si può utilizzare il metodo dei minimi quadrati: i parametri **a** e **b** sono quelli che minimizzano la somma dei quadrati degli scarti tra valori sperimentali (Y^{exp}) e valori misurati (Y^{th}):

$$\text{Min} \sum_i (y_i^{\text{exp}} - y_i^{\text{th}})^2 = \text{Min} \sum_i (y_i^{\text{exp}} - ax_i^{\text{exp}} - b)^2$$

Si dimostra che il termine noto (o intercetta) è:

$$b = \frac{(\sum_i y_i)(\sum_i x_i^2) - (\sum_i x_i)(\sum_i x_i y_i)}{n(\sum_i x_i^2) - (\sum_i x_i)^2}$$

e il coefficiente angolare è:

$$a = \frac{n(\sum_i x_i y_i) - (\sum_i x_i)(\sum_i y_i)}{n(\sum_i x_i^2) - (\sum_i x_i)^2}$$

I parametri a e b possono essere calcolati in modo esplicito (calcolando i vari termini) ma può essere complesso e fonte di errori di battitura. Più semplice è utilizzare le funzioni Excel:

PENDENZA(dati_Y;Dati_X): restituisce il coefficiente angolare della retta: **a**

INTERCETTA(dati_Y;Dati_X): restituisce l'intercetta della retta: **b**

Il parametro s_{xy} quantifica la variabilità dei dati attorno alla retta di regressione:

$$s_{xy} = \sqrt{\frac{\sum (y_i^{\text{exp}} - (ax_i + b))^2}{n - 2}}$$

dove **n-2** sono i gradi di libertà per la stima della varianza. Ho dovuto ridurre di 2 i gradi di libertà poiché ho utilizzato i dati sperimentali per stimare il valore di due parametri: **a** e **b**. Per calcolare il parametro s_{xy} si può usare la funzione Excel:

ERR.STD.YX(dati_Y;Dati_X)

Attenzione all'ordine con cui vengono inseriti i dati: prima la colonna per la variabile dipendente Y e poi la colonna dei dati per la variabile indipendente X.

Il parametro s_{xy} serve per calcolare gli errori standard sui parametri **a** e **b** della retta di regressione:

$$s_a = \frac{s_{xy}}{\sqrt{\sum_i (x_i - \bar{x})^2}} \quad \text{e} \quad s_b = s_{xy} \sqrt{\frac{1}{n} + \frac{(\bar{x})^2}{\sum_i (x_i - \bar{x})^2}}$$

Per il calcolo si noti che: $\sum_i (x_i - \bar{x})^2 = (\mathbf{n-1}) \mathbf{var(x)}$

Per stabilire se **a** o **b** sono significativamente diverso da 0 (o da un valore dato) utilizziamo il test t-

Student: la variabile t è definita:

$$t_{\nu} = \frac{\text{stima del parametro} - \text{valore da testare}}{\text{err. standard sulla stima del parametro}}$$

con ν =numero di gradi di libertà = numero di punti sperimentali – numero di parametri stimati, nel caso di una retta di regressione i parametri stimati sono 2: a e b, quindi $\nu=n-2$. Si usano i valori s_a e s_b come errori standard sulle stime dei parametri.

Per applicare il test di Student si segue quanto fatto nell'esercizio precedente: o si usano le tabelle, o si usano le funzioni Excel.

Il coefficiente di correlazione lineare (coefficiente di correlazione di Pearson) permette di quantificare quanto la variabile dipendente dipende linearmente dalla variabile indipendente:

$$r = \frac{\sum_i (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_i (x_i - \bar{x})^2 \sum_i (y_i - \bar{y})^2}} = (n - 1) \frac{\sum_i (x_i - \bar{x})(y_i - \bar{y})}{s_x s_y}$$

dove s_x e s_y sono le deviazioni standard stimate sulla variabile X e sulla variabile Y.

La funzione Excel:

CORRELAZIONE(matriceY,matriceX)

calcola il coefficiente di correlazione di Pearson.

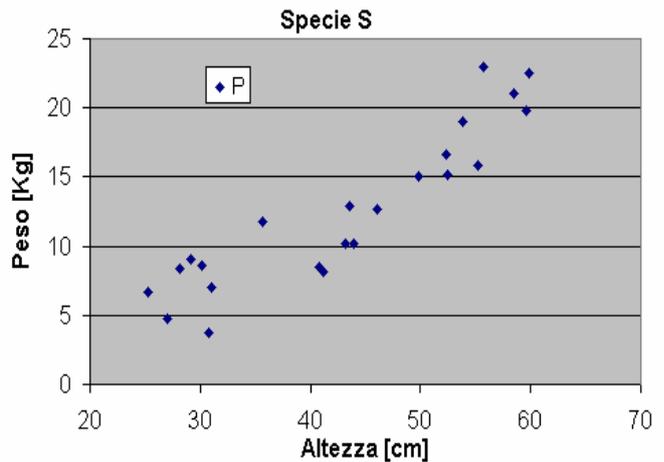
Il coefficiente di determinazione R^2 rappresenta la frazione di variabilità dei dati Y^{exp} che può essere spiegata tramite la retta di correlazione. Nel caso di regressione lineare semplice è uguale al quadrato del coefficiente di correlazione: $R^2 = r^2$

Esercizio 2) il file **regressione.xls** contiene suggerimenti utili per la soluzione.

Sono stati misurati l'altezza (H) e il peso (P) per un gruppo di n=23 animali della specie S e i dati sono riportati nel file **Specie1.dat**.

0. importare i dati in un foglio Excel
1. graficare i dati Peso in funzione dell'altezza,
2. calcolare i parametri che definiscono la retta di regressione con i loro errori standard (riportare i dati con il numero corretto di cifre significative)
3. riportare su grafico la retta di regressione $y = ax + b$
4. stabilire se **a** e **b** sono diversi da 0 con un livello di confidenza del 95 % ($\alpha = 0.05$)
5. calcolare il coefficiente di correlazione e il parametro R^2

1.: si usi l'opzione grafica di Excel: **dispersione (xy)** per graficare i dati riportando sugli assi l'unità di misura (cm per l'altezza e kg per il peso)



2/3.: per calcolare i parametri della retta di regressione possiamo seguire tre metodi. In modo complicato: si calcolano i termini x^2 e xy per ogni coppia (2) di dati e quindi le somme (3):

$$\sum_i x_i \quad \sum_i x_i^2 \quad \sum_i y_i$$

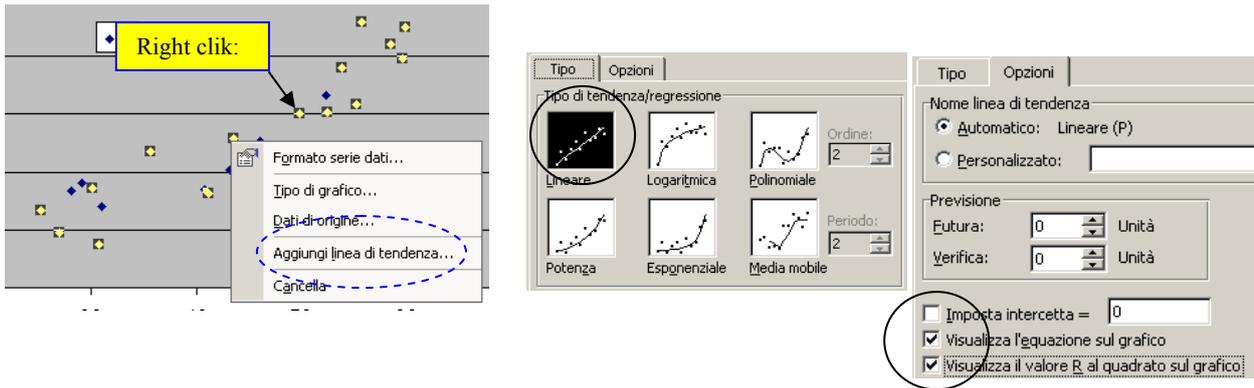
$$(\sum_i x_i)^2 \quad \sum_i x_i y_i$$

E, infine, si calcolano i termini a e:

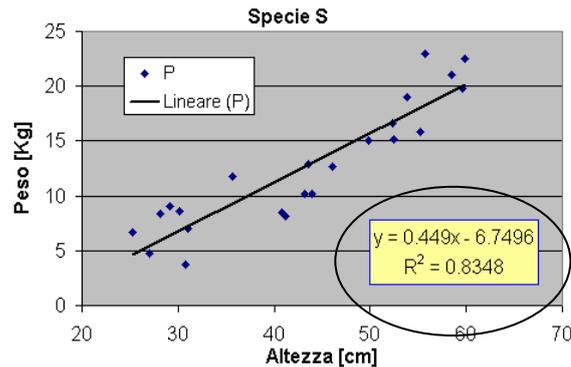
b=	-6.74956
a=	0.449047

1		2	
H	P	x^2	$x*y$
cm	kg	cm ²	cm kg
53.80	19.03	2894.44	1023.81
31.09	7.07	966.59	219.81
49.87	15.08	2487.02	752.04
35.64	11.82	1270.21	421.26
28.20	8.41	795.24	237.16
30.16	8.58	909.63	258.77
59.60	19.85	3552.16	1183.06
29.17	9.01	850.89	262.82
25.31	6.65	640.60	168.31
43.61	12.86	1901.83	560.82
46.06	12.66	2121.52	583.12
30.82	3.71	949.87	114.34
40.78	8.51	1663.01	347.04
27.00	4.78	729.00	129.06
52.42	15.14	2747.86	793.64
43.20	10.21	1866.24	441.07
55.66	23.01	3098.04	1280.74
58.46	21.00	3417.57	1227.66
43.88	10.23	1925.45	448.89
41.17	8.09	1694.97	333.07
55.18	15.85	3044.83	874.60
52.29	16.68	2734.24	872.20
59.89	22.55	3586.81	1350.52
$\sum_i x_i$	23	$\sum_i x_i^2$	$\sum_i x_i y_i$
sum(x)	993.26	sum(y)	290.78
$(\sum_i x_i)^2$	986565.4	sum(x ²)	45848.02
		sum(xy)	13883.82

Questo è un modo complesso. In modo più rapido si possono calcolare direttamente i parametri della retta di regressione dal grafico cliccando con il tasto destro del mouse su uno dei punti e scegliendo l'opzione: "aggiungi linea di tendenza". nella finestra "Opzioni" scegliere di visualizzare l'equazione e il valore di R^2 sul grafico:



Si ottiene:



La formula sul grafico mostra anche il parametro R^2 , dal quale si ricava il coefficiente di correlazione:

$$r = (R^2)^{0.5}$$

Un terzo modo per calcolare i parametri della retta di regressione consiste nell'utilizzare le funzioni Excel:

PENDENZA(dati_Y;Dati_X): restituisce il coefficiente angolare della retta: **a**
INTERCETTA(dati_Y;Dati_X): restituisce l'intercetta della retta: **b**

Si dovrebbero ottenere gli stessi risultati mostrati in precedenza.

Per calcolare gli errori sui parametri bisogna calcolare s_{xy} : si può usare la funzione Excel: **ERR.STD.YX(dati_Y;Dati_X)** e si ottiene:

$$s_{xy} = 2.37 \text{ [kg]}$$

Il parametro s_{xy} serve per calcolare gli errori standard sui parametri **a** e **b** della retta di regressione:

Funzioni Excel			
b=	-6.749582	+/-	1.945896
a=	0.449047	+/-	0.043584
s_{xy}	2.368748		

A questo punto è necessario riportare i risultati con il numero corretto di cifre significative: utilizzando due cifre significative per l'errore e accordando le cifre significative dei parametri si ottiene:

$$a = 0.449 \pm 0.044 \text{ [kg cm}^{-1}\text{]}$$

$$b = -6.7 \pm 1.9 \text{ [kg]}$$

4.: Per stabilire se **a** è significativamente diverso da 0 (lo stesso vale per **b**) utilizziamo il test t-Student

come descritto sopra. La t_{oss} :

$$t_{oss}(a) = \frac{a}{s_a}$$

quindi, per **a**:

$$t_{oss}(a) = 0.449 / 0.044 = 10.3$$

confrontiamo il valore ottenuto con il valore della variabile t-Student con $23-2=21$ gradi di libertà calcolate per un rischio $\alpha = 0.05$.

Per questo possiamo

- utilizzare le tabelle della variabile t-Student,
- utilizzare la funzione

$$\text{INV.T}(\alpha, \nu)$$

e otteniamo un valore limite per $t_{21}(\text{Lim}) = 2.08$. Dal momento che osserviamo un valore di t molto maggiore ($t_{oss}(a) = 10.3$) possiamo affermare che il valore trovato è diverso da 0 con un livello di confidenza maggiore del 95% (rischio minore del 5%).

- possiamo utilizzare la funzione

$$\text{DISTRIB.T}(t_{oss}(a), \nu, 2)$$

per calcolare la probabilità associata al valore di t osservato (Si usa un test a due code poiché siamo interessati a sapere se a è significativamente diverso da 0, se volessimo sapere se a è significativamente maggiore di 0 potremmo usare un test a una coda). Nel nostro caso abbiamo:

$$\text{DISTRIB.T}(10.03, 21, 2) = 1.1 \cdot 10^{-9}$$

cioè la probabilità di osservare per caso il valore di t osservato nell'ipotesi che il valore di **a** sia nullo è molto piccola, quindi il valore trovato è significativo

Un procedimento analogo si può seguire per il parametro b. Si ottiene la tabella seguente:

		t_limite
t(a)=	10.30	2.08
significatività	1.15E-09	
t(b)=	3.47	2.08
significatività	2.30E-03	

5.: per il calcolo del coefficiente di correlazione di Pearson usiamo la funzione Excel `CORRELAZIONE(datiX,datiY)`: otteniamo la tabella:

Correlazione	0.913699
R ²	0.834846

Ci sono diversi modi per effettuare il calcolo in modo più veloce utilizzando Excel, ne vediamo uno di

questi che utilizza la macro “regressione” del menu: ”componenti aggiuntivi-Analisi dati” :

STATISTICA DELLA REGRESSIONE	
R multiplo	0.914
R al quadrato	0.835
R al quadrato corretto	0.827
Errore standard	2.369
Osservazioni	23

ANALISI VARIANZA		
	ddl	SQ
Regressione	1	595.6261485
Residuo	21	117.8302949
Totale	22	713.4564435

	Coefficienti	Errore standard	Stat t	significatività	Inf. 95%	Sup. 95%	Inf. 90.0%	Sup. 90.0%
Intercetta	-6.7496	1.9459	-3.47	2.3E-03	-10.796	-2.703	-10.098	-3.401
Variabile X 1	0.4490	0.0436	10.30	1.1E-09	0.358	0.540	0.374	0.524

La macro regressione fornisce molte informazioni statistiche sulla retta di regressione. Le principali sono evidenziate sopra. Oltre a quelle già viste (a, b, s_{xy} , errori su a e b) ci sono:

La significatività dei parametri: rappresenta la probabilità associata a un t-test sul parametro: rappresenta la probabilità di osservare per caso il valore di a o b trovato, nell'ipotesi che il valore vero di a o b sia nullo. Se questa probabilità è alta (es. > 5%) il valore trovato è poco significativo rispetto a 0.

In questo caso i valori trovati e per i parametri a e b sono diversi da 0 con un livello di confidenza molto maggiore del 95%.

L'intervallo di confidenza: intervallo di valori per il parametro (a o b) entro il quale è stimato trovarsi il valore vero con una probabilità del 95%. Questo intervallo è dato per default. Si può definire un intervallo con un livello di confidenza definito.

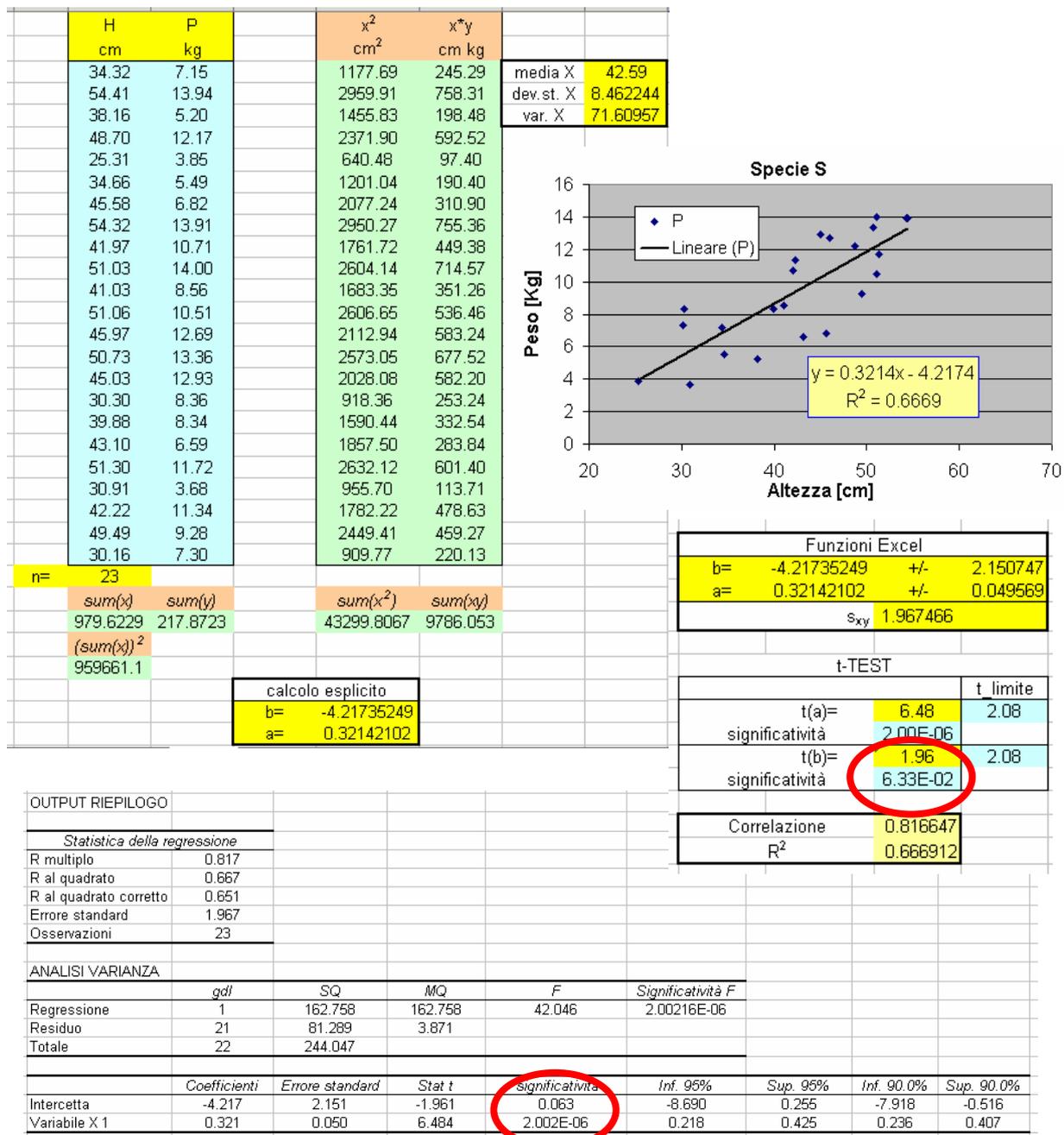
Esercizio 3)

Sono stati misurati l'altezza (H) e il peso (P) per un gruppo di n=23 animali della specie S e i dati sono riportati nel file **Specie2.dat**. Gli animali sono stati trattati con una dieta di tipo diverso rispetto al caso precedente. Si vuole sapere se la dieta ha influenza sui parametri di crescita. Per questo si vuole stabilire se le rette di regressione Peso(altezza) sono diverse per i due gruppi di animali.

- 1.: calcolare la retta di regressione per gli animali del secondo gruppo.
- 2.: utilizzare un test t per stabilire se i coefficienti delle rette di regressione per i due gruppi sono diverse con un livello di significatività migliore del 5%.

Per il punto 1 possiamo utilizzare una copia del foglio usato in precedenza per calcolare rapidamente i parametri della regressione e la loro significatività:

Si vede che nel secondo caso il termine noto (intercetta) non è significativamente diverso da 0.



Per valutare se due rette sono significativamente diverse utilizziamo un t-test. In particolare vogliamo verificare se le differenze riscontrate nei coefficienti angolari e nelle intercette sono diverse con un livello di confidenza $1-\alpha=95\%$

La variabile t-Student è definita:

$$t_{osservato} = \frac{|\text{differenza tra i parametri}|}{\text{err. standard sulla differenza tra i parametri}}$$

Quindi, per il coefficiente angolare (pendenza) si ha:

$$t_{oss} = \frac{|a_1 - a_2|}{s_{a_1 - a_2}}$$

Se il numero di osservazioni è lo stesso per ambedue le rette (è questo il caso) si ha:

$$t_{oss} = \frac{|a_1 - a_2|}{\sqrt{s_{a_1}^2 + s_{a_2}^2}}$$

L'analogo vale per l'intercetta.

Il valore della variabile t osservata deve essere confrontato con la statistica di una variabile t con v gradi di libertà, dove: $v = n_1 + n_2 - 4$ (abbiamo in totale n_1+n_2 osservazioni ma abbiamo stimato 4 parametri a_1, b_1, a_2, b_2 dai dati sperimentali)

Otteniamo quindi la tabella seguente:

	retta 1	retta 2	t_{oss}	significatività	t_{lim}
a	0.449	0.321	1.934	0.060	2.018
s_a	0.044	0.050			
b	-6.750	-4.217	0.873	0.388	
s_b	1.946	2.151			
n	23	23			

Da cui possiamo concludere che, con un livello di confidenza del 95% non possiamo affermare che il coefficiente angolare e l'intercetta delle due rette siano diverse.

Nota: se il numero di osservazioni fosse diverso l'errore standard sulla differenza tra i parametri è un po' più complicato e non verrà trattato.