

Modello logistico (Modello di regressione logistica)

Prof. Giuseppe Verlato

Sezione di Epidemiologia e Statistica Medica,
Dipartimento di Sanità Pubblica e Medicina di
Comunità, Università degli Studi di Verona

E per le variabili qualitative NOMINALI ?

2 VARIABILI (entrambe qualitative):
test del chi-quadrato, test esatto di Fischer

3 VARIABILI qualitative (2 var. + 1 var. di
stratificazione): test di Mantel-Haenszel

MOLTE VARIABILI:

y dicotomica (malato/sano) → modello LOGISTICO

y politomica (fumatore, ex-fumatore, mai-fumatore)
→ modello MULTINOMIALE

MODELLO DI REGRESSIONE LOGISTICA

19 / (19+132)	prevalenze
0 / (0+9)	
11 / (11+52)	
6 / (6+97)	

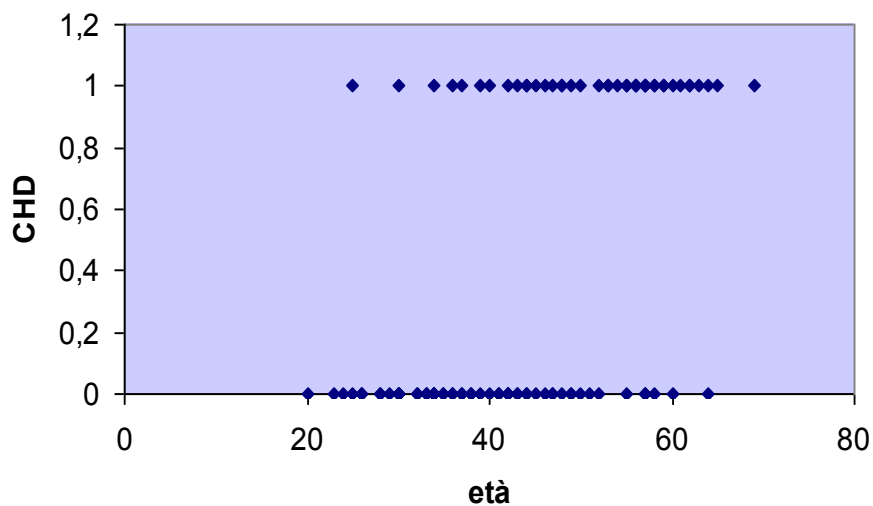
MODELLO LOG-LINEARE

19	132	conteggi
0	9	
11	52	
6	97	

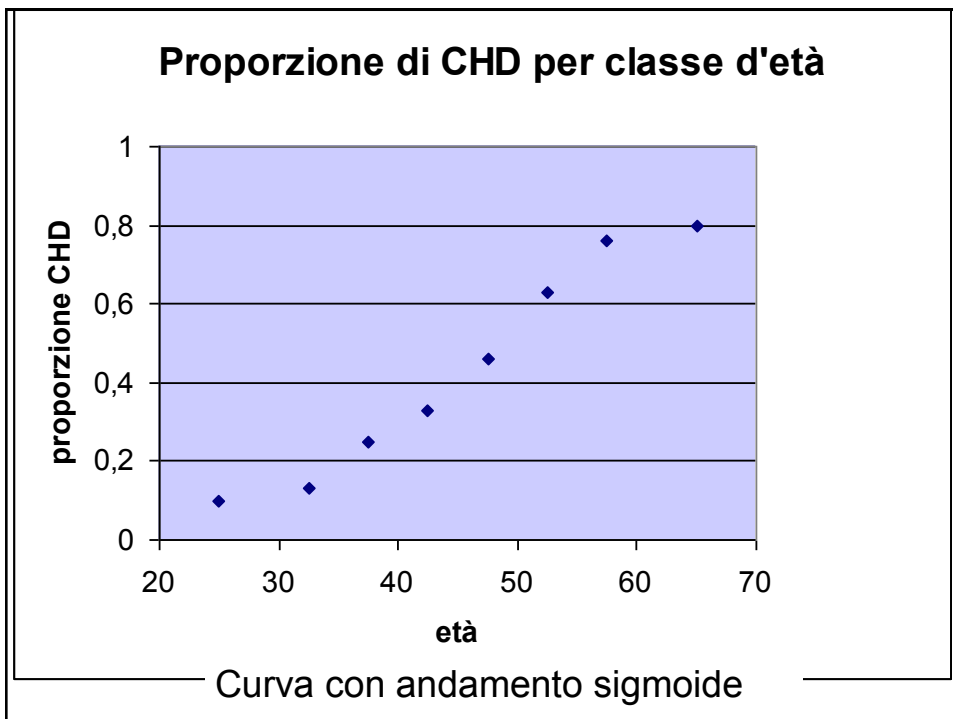
MODELLO DI POISSON

19 / 1510 persone-anno	incidenze
0 / 90 persone-anno	
11 / 630 persone-anno	
6 / 103 persone-anno	

Presenza di CHD in base all'età



CHD				
<i>classe d'età</i>	<i>N</i>	<i>assente</i>	<i>presente</i>	<i>proporzione</i>
20-29	10	9	1	0,1
30-34	15	13	2	0,13
35-39	12	9	3	0,25
40-44	15	10	5	0,33
45-49	13	7	6	0,46
50-54	8	3	5	0,63
55-59	17	4	13	0,76
60-69	10	2	8	0,8
totale	100	57	43	0,43



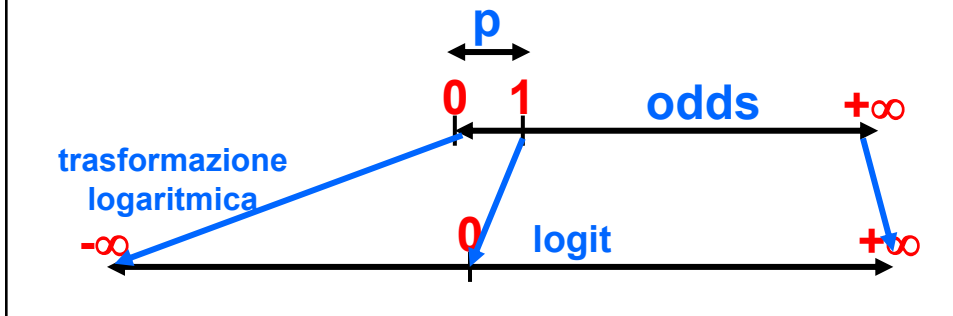
Nella regressione lineare multipla la Y varia tra $-\infty$ e $+\infty$

Nella regressione logistica

p(malattia) varia tra 0 e 1

odds(malattia) = $p/(1-p)$ varia tra 0 e $+\infty$

Logit = $\ln [p/(1-p)]$ varia tra $-\infty$ e $+\infty$



I MODELLI LINEARI GENERALIZZATI si differenziano per la **distribuzione dell'errore (error function)** e per la **funzione legame (link function)**

REGRESSIONE LINEARE MULTIPLA

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \beta_{13} x_1 x_3 + \varepsilon$$

La funzione legame (link-function) è l'IDENTITÀ

L'errore segue la distribuzione NORMALE

MODELLO DI REGRESSIONE LOGISTICA

$$\ln [y/(1-y)] = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \beta_{13} x_1 x_3 + \varepsilon$$

La funzione legame (link-function) è il LOGIT [LOG(ODDS)]

L'errore segue la distribuzione BERNOULLIANA

MODELLO LOG-LINEARE

$$\ln(y) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \beta_{13} x_1 x_3 + \varepsilon$$

La funzione legame (link-function) è il LOGARITMO

L'errore segue la distribuzione di POISSON

MODELLO DI REGRESSIONE LOGISTICA

Predittore lineare

$$\text{Ln} [\pi/(1-\pi)] = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \beta_{13} x_1 x_3$$

Logit

Var.qualitative e/o quantitative Termine d'interazione

$$\pi/(1-\pi) = \exp(\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \beta_{13} x_1 x_3)$$

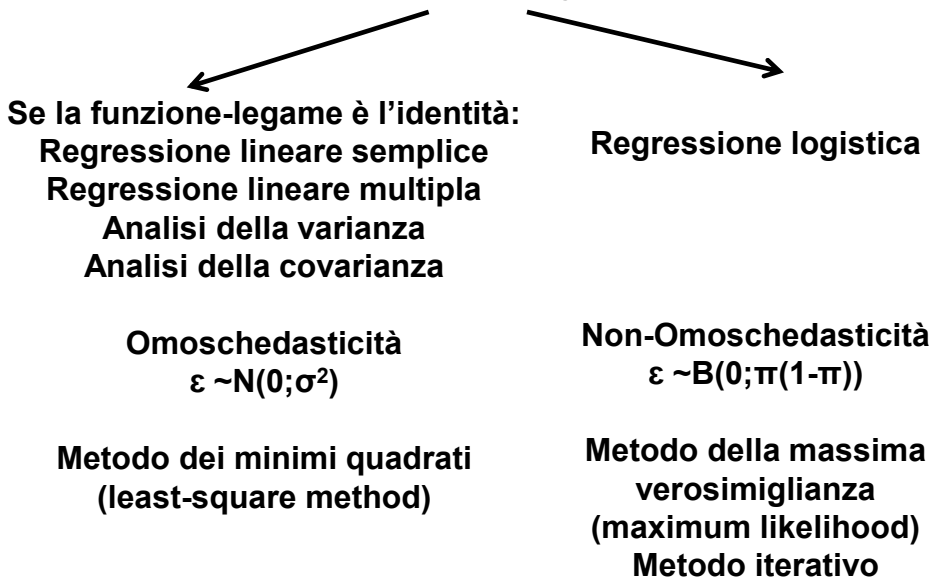
odds

$$\pi = \frac{\exp(\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \beta_{13} x_1 x_3)}{1 + \exp(\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \beta_{13} x_1 x_3)}$$

prevalenza

Metodi di ottimizzazione

per trovare il modello che meglio si adatta ai dati



30 FIGLI MASCHI SU 40 NASCITE

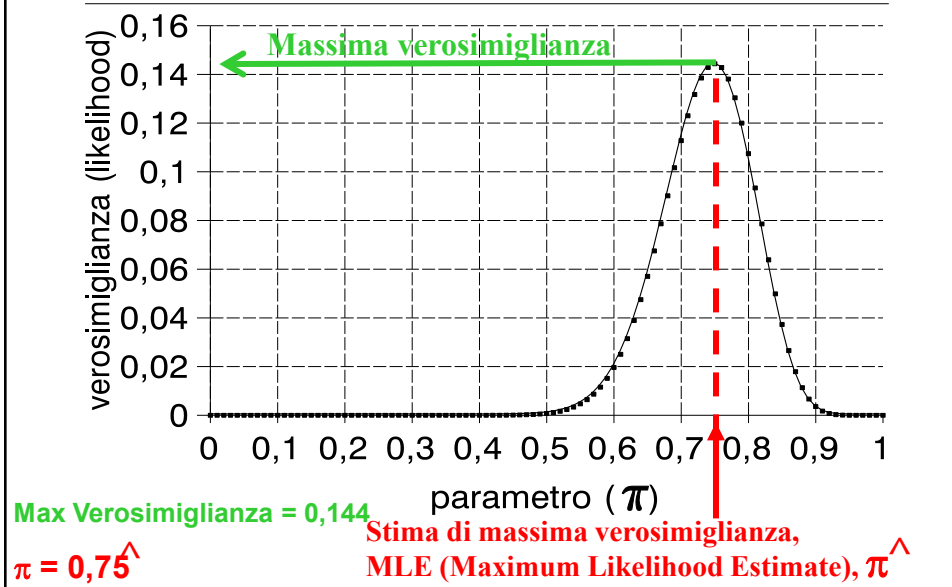


Table 4. Risk factors for e-cigarette ever use.

Risk factors	OR of e-cigarette ever use (95%IC)	p-value
Smoking Habits		
<i>Never Smokers</i>	1	
<i>Ex-Smokers</i>	5.97 (3.46 – 10.33)	<0.001
<i>Occasional Smokers</i>	7.75 (5.07 – 11.85)	<0.001
<i>Regular Smoker</i>	24.6 (16.62 – 36.50)	<0.001
Sex		
<i>Male</i>	1	
<i>Female</i>	0.45 (0.32 – 0.63)	<0.001
Family history of smoking		
<i>None</i>	1	
<i>Relative who smoked e-cig</i>	0.68 (0.37 – 1.25)	0.219
<i>Relative who smoked tobacco</i>	1.19 (0.92 – 1.53)	0.178
Housemates currently smoking		
<i>None</i>	1	
<i>Housemates smoking e-cig</i>	2.54 (1.15 – 5.59)	0.021
<i>Housemates smoking tobacco</i>	1.16 (0.94 – 1.43)	0.163
Centre		
<i>Verona</i>	1	
<i>Vicenza</i>	1.12 (1.01 – 1.24)	0.027
<i>Legnago</i>	0.89 (0.85 – 0.94)	<0.001
<i>Trento</i>	0.82 (0.73 – 0.91)	<0.001
<i>Bolzano</i>	0.66 (0.63 – 0.69)	<0.001
University class		
<i>1st year</i>	1	
<i>2nd year</i>	1.01 (0.73 – 1.41)	0.939
<i>3rd year</i>	0.79 (0.56 – 1.11)	0.173

OR (Odds Ratio), 95% IC, *P*-values were computed by logistic regression model. Significant results were highlighted in bold.

Canzan et al, BMC Public Health 2019

Logistic regression

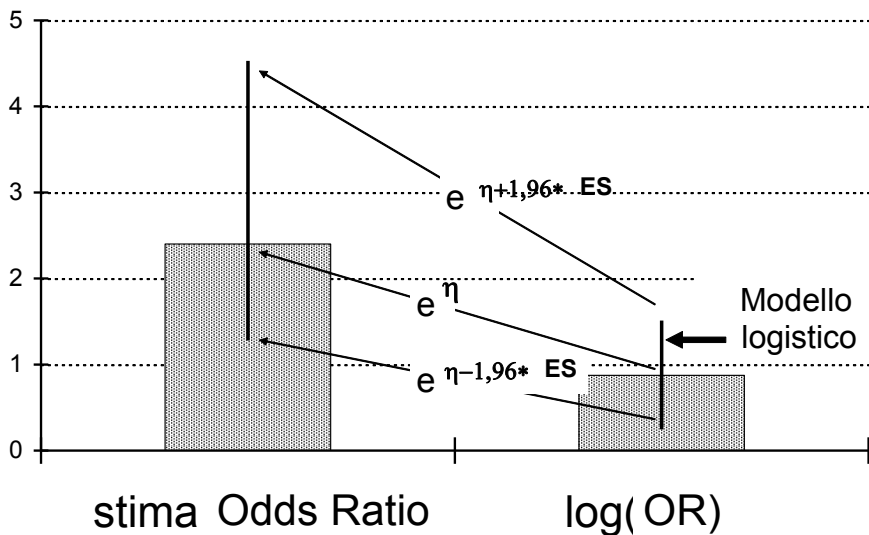
Log likelihood = -597.07448

Number of obs = 1,354
 LR chi2(14) = 476.34
 Prob > chi2 = 0.0000
 Pseudo R2 = 0.2852

provatoecg	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]	
_IsmokHab_1	1.787651	.2388358	7.48	0.000	1.319541	2.25576
_IsmokHab_2	2.047776	.2235388	9.16	0.000	1.609648	2.485904
_IsmokHab_3	3.203935	.1918879	16.70	0.000	2.827841	3.580028
_Icesso_2	-.8049849	.1684709	-4.78	0.000	-1.135182	-.4747881
_Ifamiliari_1	-.3819719	.2876365	-1.33	0.184	-.9457291	.1817854
_Ifamiliari_2	.1720737	.1825327	0.94	0.346	-.1856837	.5298312
_Iconvivent_1	.9307554	.3162551	2.94	0.003	.3109068	1.550604
_Iconvivent_2	.1494138	.1920308	0.78	0.437	-.2269596	.5257873
_Isede_2	.1159116	.2095446	0.55	0.580	-.2947882	.5266115
_Isede_3	-.1109284	.2274942	-0.49	0.626	-.5568088	.334952
_Isede_4	-.2020076	.2462749	-0.82	0.412	-.6846976	.2806824
_Isede_5	-.4146958	.2101358	-1.97	0.048	-.8265544	-.0028371
_Iannodisor_2	.0129026	.1907115	0.07	0.946	-.360885	.3866903
_Iannodisor_3	-.2393992	.1827322	-1.31	0.190	-.5975477	.1187494
_cons	-1.872191	.2466523	-7.59	0.000	-2.355621	-1.388762

CONSEGUENZE della TRASFORMAZIONE LOGARITMICA:

L'intervallo di confidenza diventa asimmetrico



Quando abbiamo una sola variabile indipendente, possiamo verificare che il parametro β stimato dalla regressione logistica corrisponde al $\ln(\text{OR})$ calcolato dalla tabellina corrispondente.

Es. Età e CHD: dividiamo l'età in 2 categorie $<55(x=0)$ e $\geq 55(x=1)$

	Coefficiente Stimato	Errore Standard	Coeff./ES	OR
AGE	2.094	0.529	3.96	8.1
Constant	-0.841	0.255	-3.30	

Costruendo ora la tabellina 2X2:

CHD (y)	AGE(x)		
	$\geq 55(1)$	<55	
Presente (1)	21	22	43
Assente (0)	6	51	57
Totale	27	73	100

$$\text{OR} = 21 \cdot 51 / (6 \cdot 22) = 8.11$$

che quindi corrisponde a quanto trovato con la regressione logistica