

## Contents

<b>1</b>	<b>Generalizzazioni del modello lineare</b>	<b>5</b>
1.1	Un approccio un po' antico: trasformazioni di variabili; riduzione dell'eteroscedasticità . . . . .	5
1.2	Trasformazioni classiche . . . . .	11
<b>2</b>	<b>Regressione logistica</b>	<b>12</b>
2.1	Risposte singole . . . . .	12
2.2	Esempi visti in aula . . . . .	16
2.3	Modellizzazione della probabilità di un evento . . . . .	20
<b>3</b>	<b>Estensioni della regressione logistica semplice</b>	<b>34</b>
<b>4</b>	<b>Verosimiglianza ed inferenza nella regressione logistica</b>	<b>39</b>
<b>5</b>	<b>Conteggi: Regressione di Poisson</b>	<b>41</b>
<b>6</b>	<b>Verosimiglianza nella regressione di Poisson</b>	<b>45</b>

<b>7</b>	<b>Modelli lineari generalizzati</b>	<b>47</b>
7.1	La famiglia esponenziale . . . . .	47
7.2	Valore atteso e varianza . . . . .	51
<b>8</b>	<b>Generalized linear models (GLM)</b>	<b>53</b>
8.1	cenni sull'inferenza nei GLM . . . . .	55
8.2	proprietà asintotiche di $\mathbf{b}$ (link canonico) . . . . .	59
8.3	legame non canonico . . . . .	60
8.4	confronto fra modelli . . . . .	61

**List of Figures**

1	Un grafico di residui in funzione di valori stimati che mostra una possibile eteroscedasticità . . . . .	6
2	Un grafico location scale che mostra una possibile eteroscedasticità . . . . .	7
3	Un grafico deviazione standard vs medie calcolate in 4 gruppi . . . . .	8
4	una variabile binaria in funzione di una var. quantitativa . . . . .	14
5	una variabile binaria in funzione di una var. quantitativa . . . . .	15

6	Risultati deille partite di un campionato di calcio in funzione della quota data dagli allibratori per la vittoria della squadra di casa; $Y = 1$ indica la vittoria della squadra di casa, $Y = 0$ pareggio o sconfitta . . . . .	17
7	Dati raggruppati: frequenza di immatricolazione alla nostra facoltà nel 2009-2010 in funzione del risultato nel test d'ingresso . . . . .	18
8	Probabilità costante al variare di $x$ . . . . .	23
9	Massima dipendenza della probabilità da $x$ . . . . .	24
10	Probabilità funzione lineare a tratti di $x$ . . . . .	26
11	Funzione logistica e approssimazione lineare in $x_0$ . . . . .	30
12	Funzione logistica (2) per quattro diversi valori di $\beta$ : $\beta = 1; 2; 4e10$ . . . . .	32
13	Risultati deille partite di un campionato di calcio: gol segnati dalla squadra di casa in funzione della quota data dagli allibratori per la vittoria della squadra di casa . . . . .	42

(14 maggio 2020, versione molto provvisoria, considerate questo materiale un promemoria per alcune formule e dimostrazioni e un supporto alle lezioni, che quest'anno peraltro sono in buona parte registrate e restano in piattaforma)

Consultare il materiale dei notebook, degli esercizi e laboratorio

Per tenere aggiornato il software didattico del mio package MLANP, potete installarlo dal mio repository github. In Rstudio installate devtools e poi `devtools::install_github("marcellochiodi/packages")`

# 1 Generalizzazioni del modello lineare (Regressione logistica, di Poisson e Generalized linear models)

## 1.1 Un approccio un po' antico: trasformazioni di variabili; riduzione dell'eteroscedasticità

Scopo:

- stabilizzazione della varianza
- o comunque riduzione dell'eteroscedasticità

Dall'analisi dei residui o dalla rappresentazione grafica delle varianze in funzione delle medie per dati raggruppati ci si potrebbe rendere conto che la varianza appare funzione della media:

Si veda per esempio la figura 1, ottenuta dai residui di un modello lineare adattato per descrivere i tempi di guarigione in funzione di alcuni fattori prognostici e diagnostici.

L'andamento dei residui in ordinata sembra variare di più al crescere dei valori stimati. Anche il grafico fornito da R (nella figura 2) sembra confermare tale indicazione.

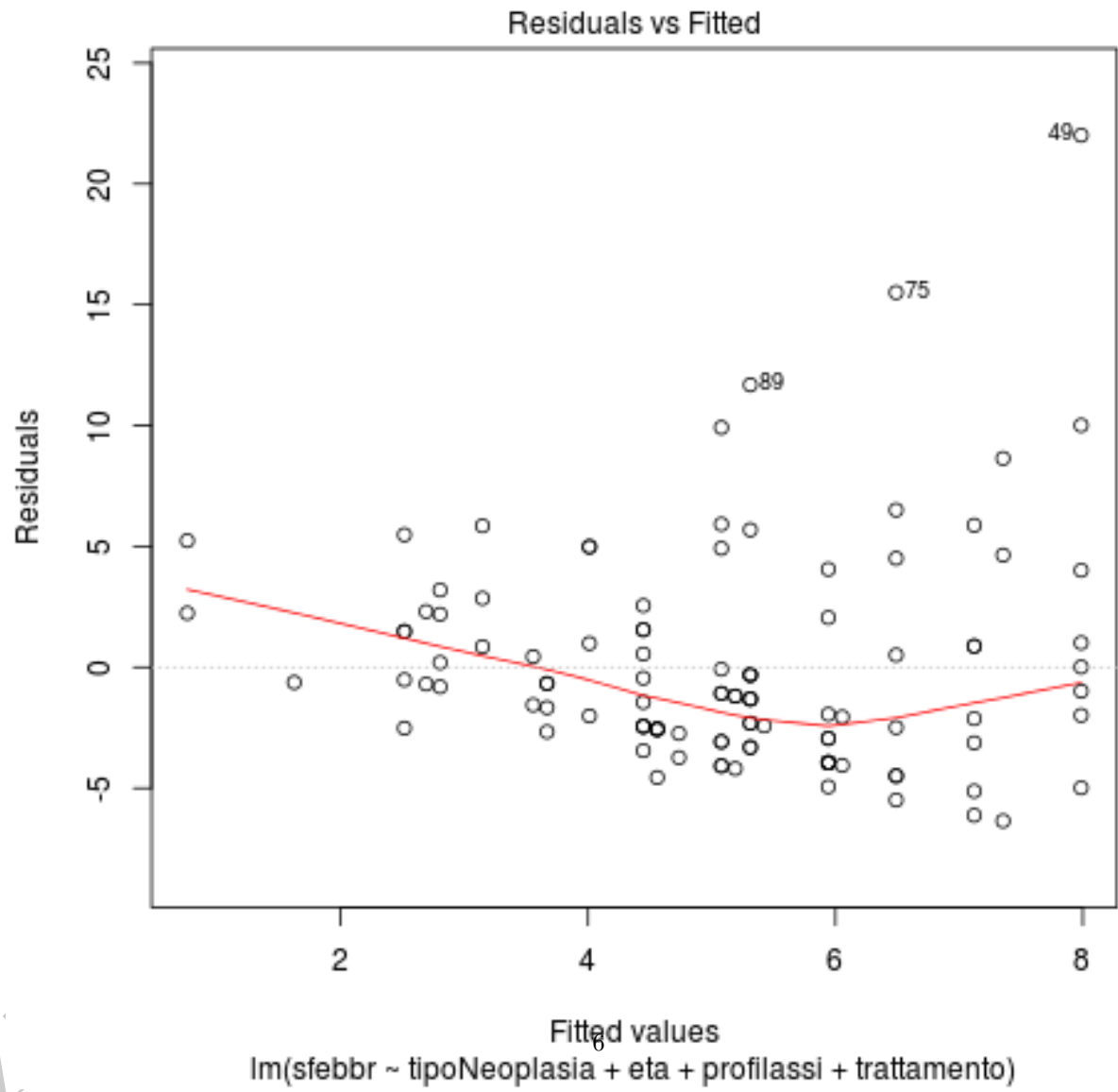


Figure 1: Un grafico di residui in funzione di valori stimati che mostra una possibile eteroscedasticità

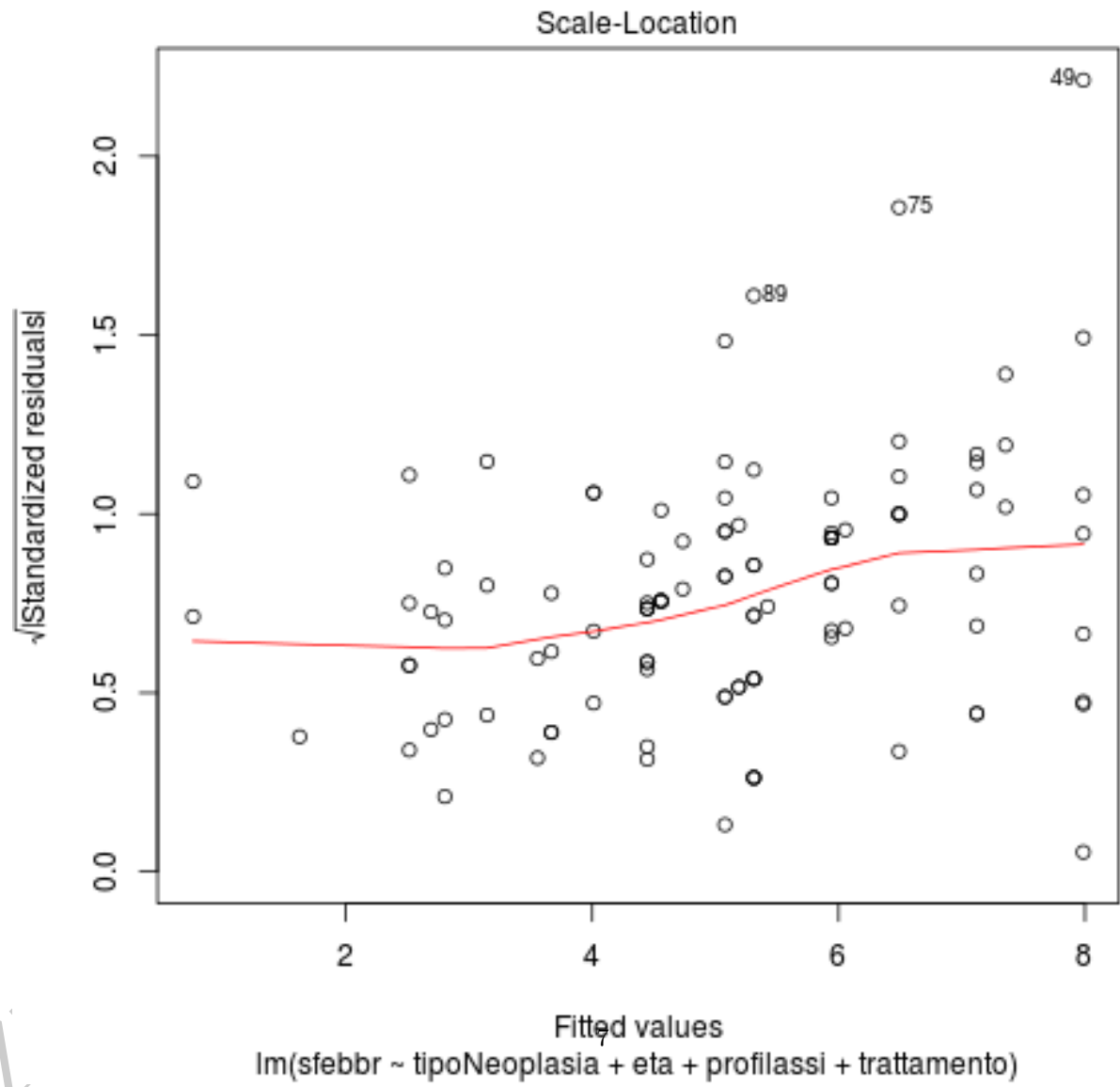


Figure 2: Un grafico location scale che mostra una possibile eteroscedasticità

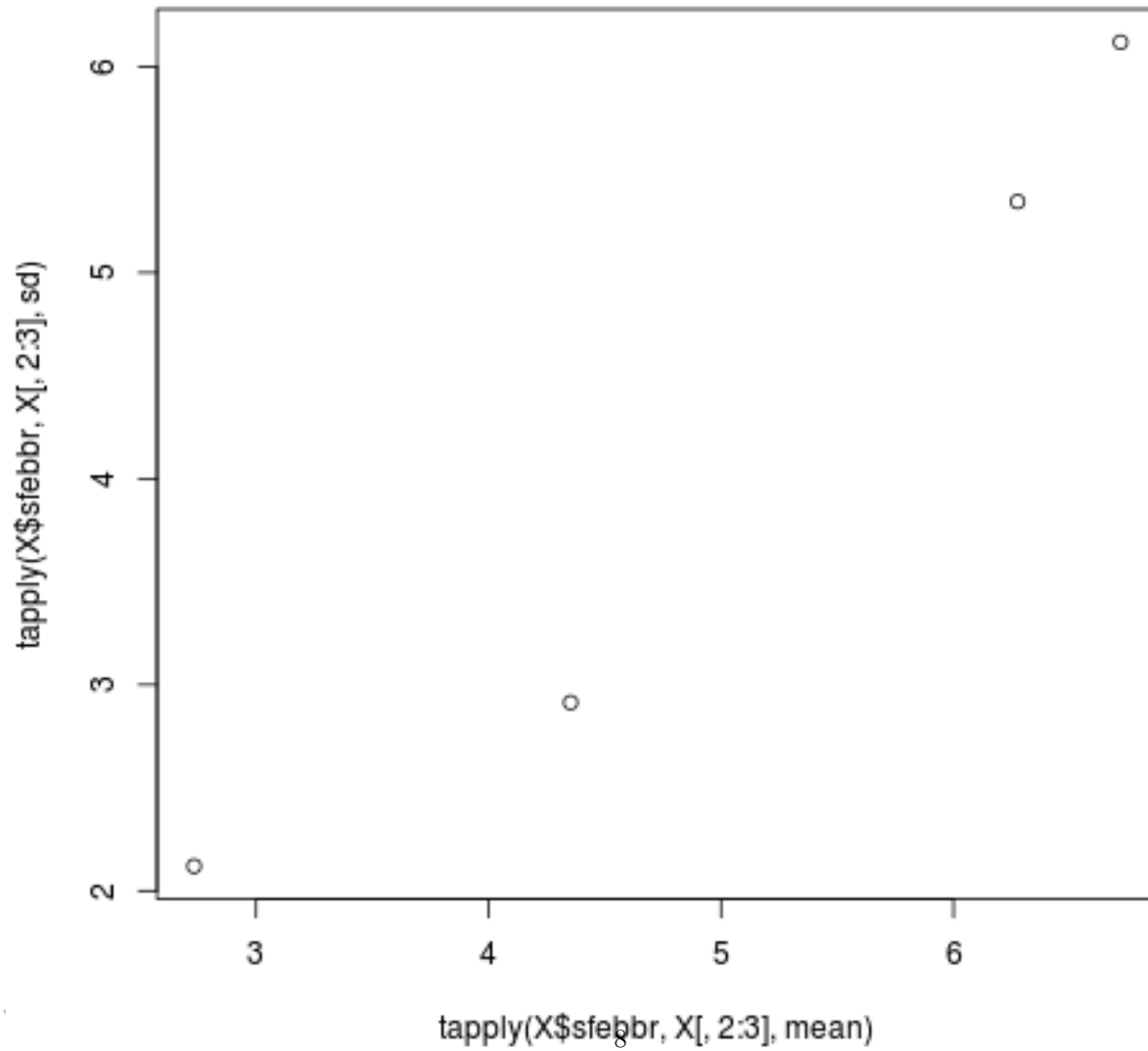


Figure 3: Un grafico deviazione standard vs medie calcolate in 4 gruppi



Nella figura 3 sono stati invece rappresentati le deviazioni standard in funzione delle medie per la stessa variabile in 4 diversi gruppi di pazienti (esempio visto a lezione)

In pratica stiamo vedendo che ciò potrebbe essere un indizio del fatto che da:

$$V [Y_i] = g(E [Y_i]),$$

mentre usualmente assumiamo  $V [Y_i] = \sigma^2$

(p.e.  $V [Y_i] = \sigma^2(E [Y_i])^\theta$ )

si può in prima battuta cercare una trasformazione  $Z_i = f(Y_i)$  che renda le nuove variabili  $Z$  *approssimativamente* con varianza uguale.

Sappiamo che se

$$E[Y_i] = \eta_i \quad \text{e} \quad V[Y_i] = g(\eta_i),$$

possiamo approssimare i primi due momenti di  $Z_i = f(Y_i)$  mediante le note relazioni (ricavate sviluppando in serie in  $\eta_i$  e troncando al primo termine):

$$\begin{aligned} E[f(Y_i)] &\approx f(\eta_i); \\ V[f(Y_i)] &\approx (f'(y_i)_{y_i=\eta_i})^2 V[Y_i] = \\ &= (f'(y_i)_{y_i=\eta_i})^2 g(\eta_i) \end{aligned}$$

per cui occorre (se la funzione  $g()$  è nota!) trovare una  $f(Y_i)$  tale che:

$$V[f(Y_i)] \approx k^2 \quad (\text{costante})$$

e quindi:

$$f'(y_i)_{y_i=\eta_i} = \frac{k}{\sqrt{g(\eta_i)}} \implies f(y_i) = \int \frac{k}{\sqrt{g(y)}} dy$$

## 1.2 Trasformazioni classiche

1. Distribuzione di Poisson:

$$V [f(Y_i)] = E [f(Y_i)] = \eta_i$$

per cui :

$$f(y_i) = \sqrt{y_i}$$

2. Distribuzione di una frequenza relativa di successo:

$$E [f(Y_i)] = p_i$$

$$V [f(Y_i)] = \frac{p_i(1 - p_i)}{n_i}$$

per cui:

$$f(y_i) = \sqrt{n_i} \arcsin(\sqrt{Y_i})$$

3. Dispersione proporzionale alla media:

$$\sqrt{V [f(Y_i)]} = k E [f(Y_i)]$$

per cui:

$$f(y_i) = \log Y_i$$

## 2 Regressione logistica

### 2.1 Risposte singole

Supponiamo di avere una situazione osservata diversa da quelle affrontate nei capitoli precedenti, con una matrice  $\mathbf{X}$  di regressori analoga, ma con una variabile di risposta qualitativa, in particolare dicotoma, diciamo senza perdita di generalità di tipo 0-1. Potrebbe indicare insuccesso-successo, guarigione, ingresso nel mondo del lavoro, successo di un evento sportivo, etc.

Fin qui abbiamo sempre considerato come variabile di risposta una variabile quantitativa: abbiamo cercato di capire se e come il valore atteso di tale variabile dipende da altre variabili osservabili qualitative o quantitative, possibilmente secondo una relazione lineare.

Adesso ci poniamo un problema leggermente diverso e per semplificare le cose supponiamo di avere come variabile esplicativa un solo regressore quantitativo  $X$ , e come variabile di risposta una variabile  $Y$  che può assumere solo i valori 0 e 1. Ad esempio  $Y$  potrebbe indicare la condizione di guarigione da una certa patologia e  $X$  potrebbe essere la dose somministrata. Oppure  $Y$  potrebbe indicare

la condizione lavorativa (0 per i disoccupati e 1 per gli occupati) e  $X$  il numero di anni di studio, o il fatto che una squadra abbia vinto o meno un confronto sportivo in funzione della quotazione fornita dagli allibratori.

Nella figura 4 è mostrato un esempio di questa situazione: *l'evento  $Y = 1$  è più frequente per valori di  $X$  alti, mentre per valori di  $X$  bassi l'evento positivo si presenta raramente.*

Vi è anche una zona intermedia di valori di  $X$  in cui sembra essere maggiore l'incertezza rispetto al verificarsi o meno dell'evento  $Y = 1$ .

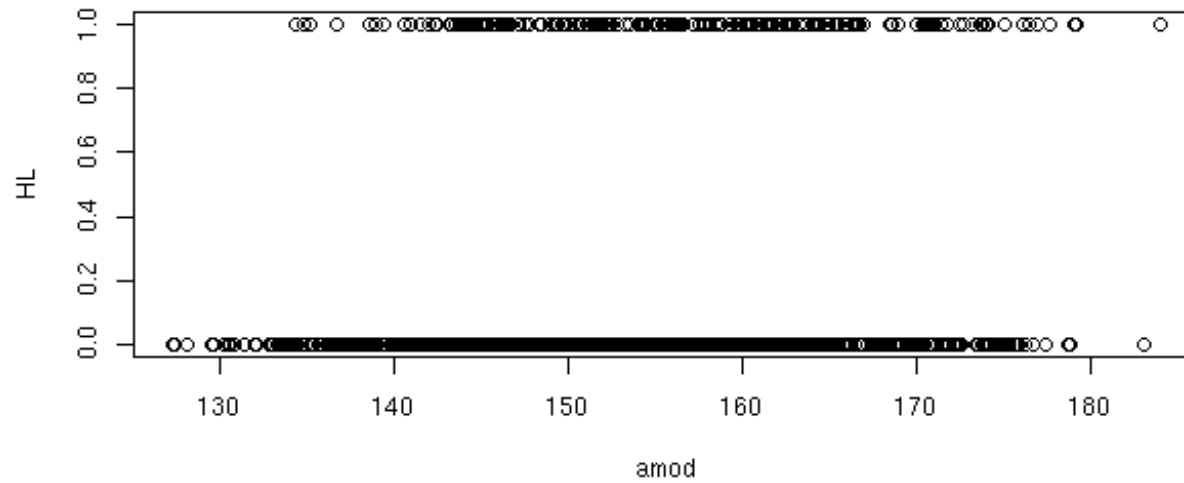


Figure 4: una variabile binaria in funzione di una var. quantitativa

Nella figura 5 è mostrato un esempio in cui la frequenza dell'evento  $Y = 1$  non sembra dipendere in modo marcato dal valore assunto da  $X$

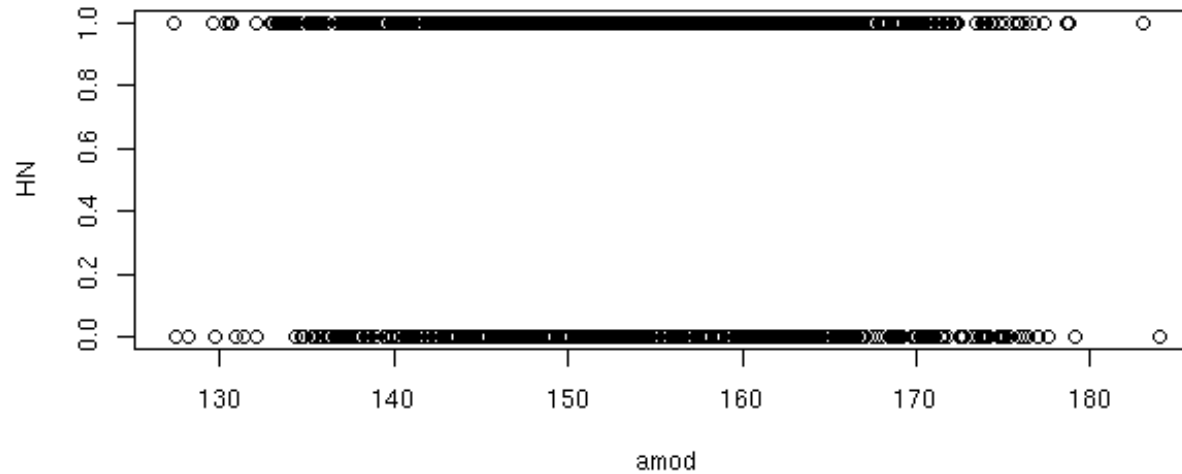


Figure 5: una variabile binaria in funzione di una var. quantitativa

## 2.2 Esempi visti in aula

BOZZE MARCELLO CHIODI 2020



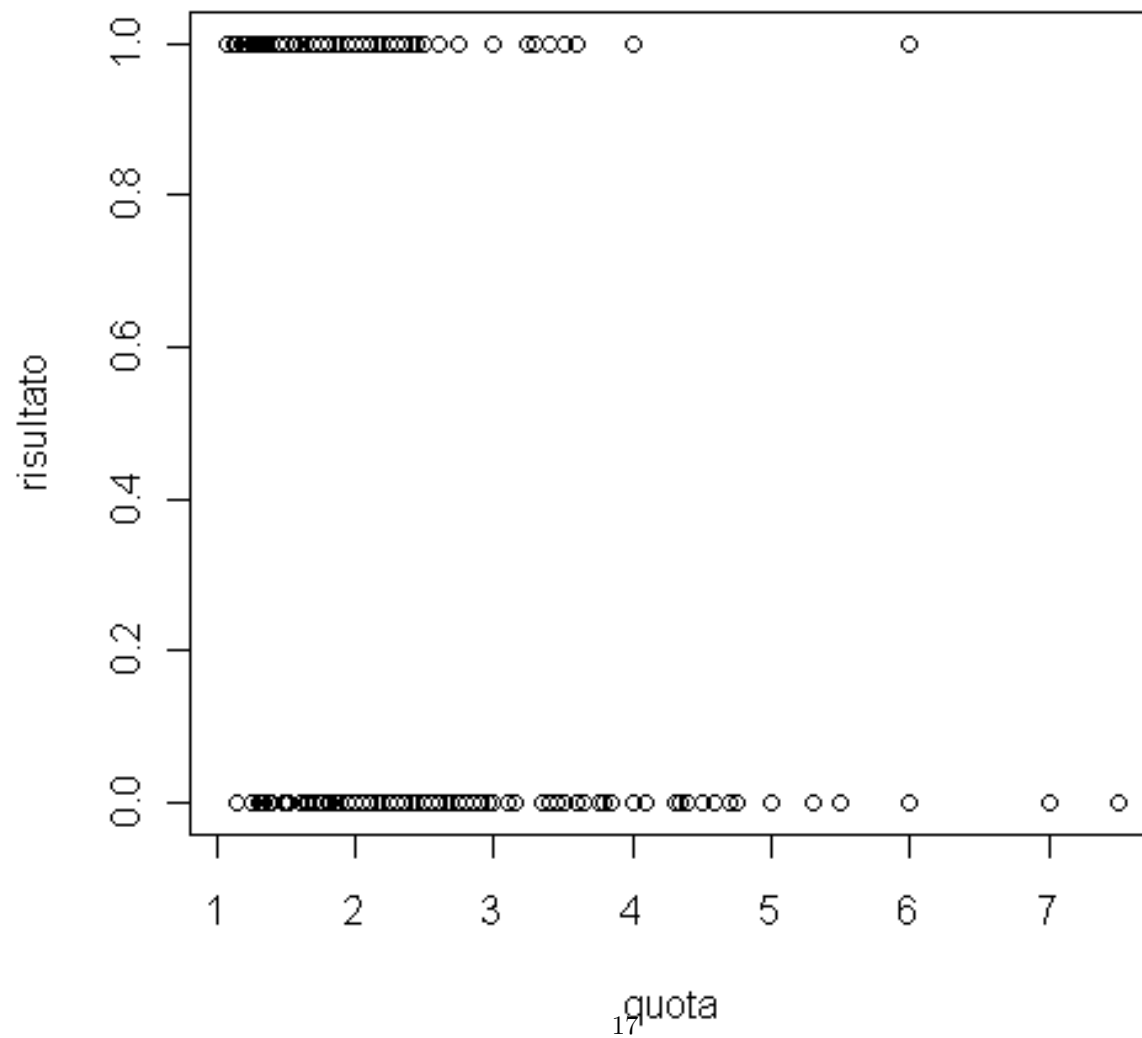


Figure 6: Risultati delle partite di un campionato di calcio in funzione della quota data dagli allibratori per la vittoria della squadra di casa;  $Y = 1$  indica la vittoria della squadra di casa,  $Y = 0$  pareggio o sconfitta

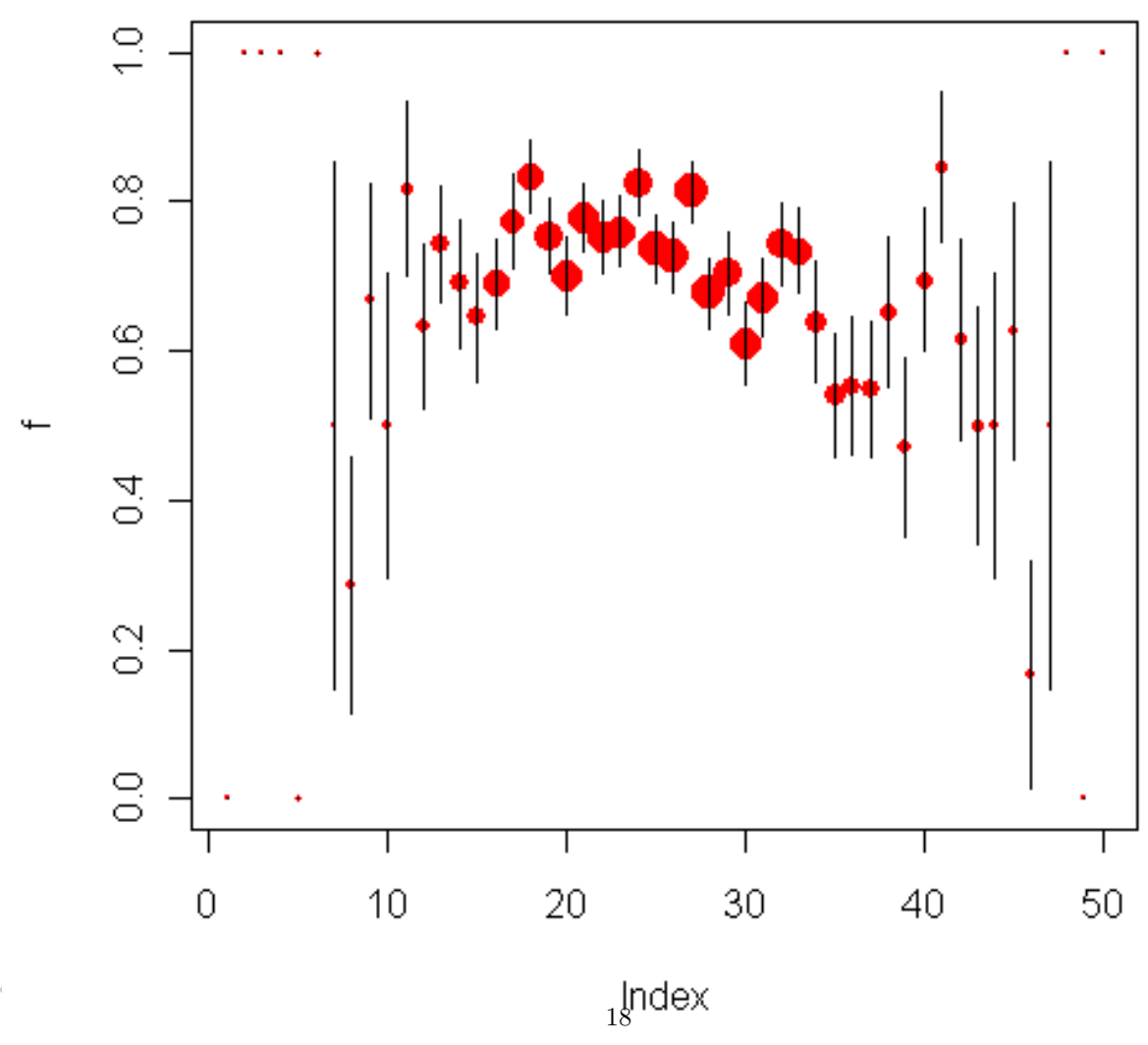


Figure 7: Dati raggruppati: frequenza di immatricolazione alla nostra facoltà nel 2009-2010 in funzione del risultato nel test d'ingresso

Ma potremmo ovviamente essere interessati ad analizzare la dipendenza dell'esito di un esperimento (indicata dalla variabile `esito`) da una **molteplicità di fattori**, quantitativi e non.

```
> str(trial1)

'data.frame': 110 obs. of 11 variables:
 $ id          : int  1 2 3 4 5 6 7 8 9 10 ...
 $ codiceCentro : int 1501 1501 1501 603 1501 1308 1501 501 501 1308 ...
 $ sesso       : Factor w/ 2 levels "f","m": 2 1 2 2 1 2 1 1 1 1 ...
 $ livelloRischio: Factor w/ 2 levels "alto","basso": 2 2 1 2 2 2 2 1 1 2 ...
 $ tipoNeoplasia : int  0 1 1 0 0 1 1 1 1 0 ...
 $ eta         : int  0 1 0 1 1 0 1 1 1 0 ...
 $ profilassi   : int  1 0 1 0 1 1 1 0 1 1 ...
 $ faseMalattia.1: int  0 0 0 2 2 0 0 0 0 0 ...
 $ trattamento  : Factor w/ 3 levels "A","B","C": 2 3 2 1 3 1 2 3 2 1 ...
 $ sfebbr       : int  2 6 7 3 5 2 5 12 5 3 ...
 $ esito        : int  1 1 1 1 1 1 0 1 1 1 ...
```

### 2.3 Modellizzazione della probabilità di un evento

Ipotizzata l'indipendenza fra le osservazioni, possiamo supporre che  $Y$ , essendo un indicatore di evento, abbia una distribuzione bernoulliana.

Se non esistesse una dipendenza da  $X$  avremmo semplicemente:

$$\text{Prob} \{Y = 1\} = p.$$

Adesso vogliamo vedere se tale probabilità (che per la variabile  $Y$  è anche valore atteso) dipende in qualche modo da  $X$ . Possiamo ipotizzare allora che  $Y$ , indicatore di evento, condizionatamente a ciascun valore  $x$  di  $X$ , abbia una distribuzione bernoulliana con parametro funzione di  $x$ <sup>1</sup> :

$$\text{Prob} \{Y = 1 | X = x\} = p(x)$$

---

<sup>1</sup>Va sottolineato che sebbene stia usando il termine *condizionatamente*, non si sta ipotizzando che  $X$  sia una variabile aleatoria, e quindi non occorre neanche fare delle ipotesi sulla distribuzione congiunta di  $X, Y$

Ci proponiamo di studiare la dipendenza da  $X$  di tale probabilità sulla base di un campione di osservazioni (indipendenti). Prima però vediamo di inquadrare il problema in un ambito familiare, magari collegabile poi ai modelli lineari.

Nella distribuzione di Bernoulli sappiamo che il valore atteso coincide con la probabilità del verificarsi dell'evento, per cui per la funzione di regressione si ha:

$$E[Y|X = x] = p(x)$$

**Dunque ci stiamo ponendo sempre la stessa domanda:**

Il valore atteso di  $Y$  dipende da  $x$ ?

Faccio subito notare che con le assunzioni che abbiamo fatto il modello è intrinsecamente eteroscedastico.

Infatti sappiamo che per tale variabile bernoulliana  $V[Y|X = x] = p(x)(1 - p(x))$ , che risulta massima quando  $p(x) = \frac{1}{2}$ ; d'altra parte nel grafico mostrato prima per valori di  $X$  intermedi avevamo una maggiore incertezza, mentre agli estremi della distribuzione di  $X$  sembra esserci meno incertezza, e quindi meno variabilità, sulla determinazione che assume  $Y$ .

Tornando comunque alla dipendenza di  $E[Y|X = x]$  da  $X$ , che forma può assumere questa dipendenza? Teniamo presente che  $0 \leq p(x) \leq 1$ , quindi non potremo assumere una semplice funzione lineare in  $x$  (perchè non potrebbe mai essere limitata, se non nel caso di una funzione costante)!

Come spesso accade, conviene ragionare prima sui casi estremi: *assenza di dipendenza in media e massima dipendenza*.

L'assenza di dipendenza in media è facilmente rappresentata da una relazione costante in  $x$ :

$$p(x) = p_0 \text{ (figura 8).}$$

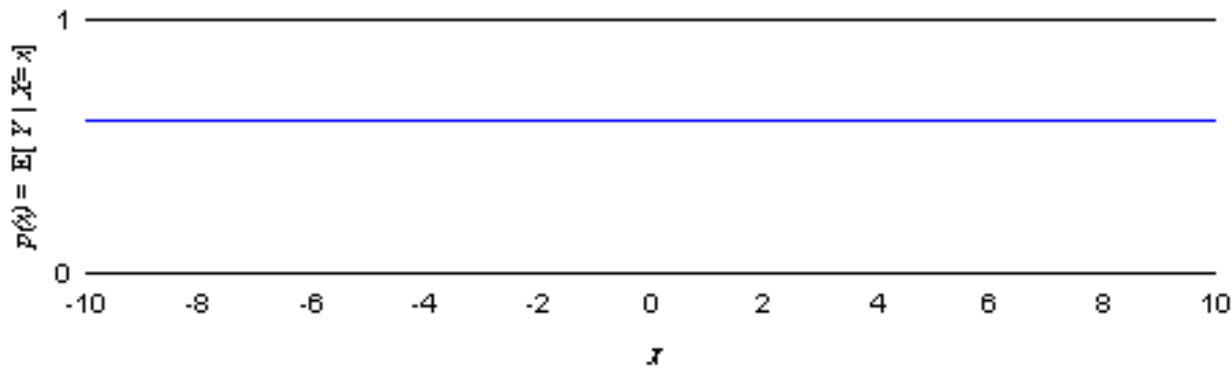


Figure 8: Probabilità costante al variare di  $x$

Per la massima dipendenza possiamo pensare che esista un valore di soglia  $x_0$  tale che:

$$p(x) = \begin{cases} 0 & \text{se } x < x_0 \\ 1 & \text{se } x \geq x_0 \end{cases};$$

Possiamo esprimere graficamente la situazione <sup>2</sup> (figura 9):

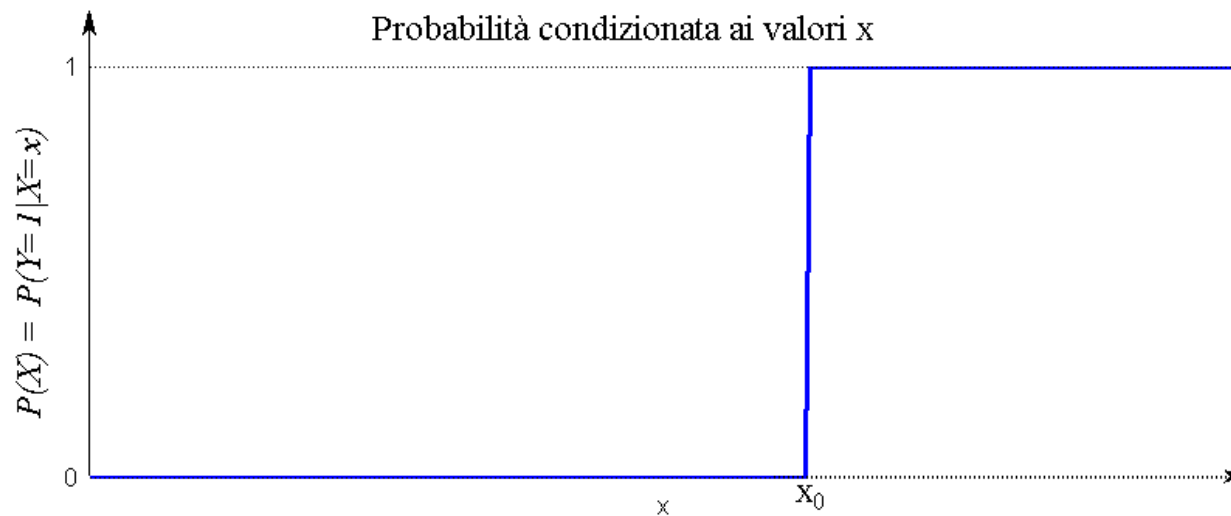


Figure 9: Massima dipendenza della probabilità da  $x$

<sup>2</sup>Ovviamente non è l'unica possibile situazione di dipendenza massima, ma è l'unica in caso di dipendenza monotona. In alternativa potremmo avere probabilità unitaria in un intervallo di valori di  $x$  e 0 fuori dall'intervallo, ma non sarebbe monotona



In pratica il verificarsi dell'evento è determinato con certezza dal valore assunto da  $x$  rispetto a  $x_0$ . Nel caso in cui si abbia

$$p(x) = \begin{cases} 1 & \text{se } x < x_0; \\ 0 & \text{se } x \geq x_0; \end{cases}$$

avremo allora una dipendenza massima di segno negativo.

Ovviamente nasce il problema di come descrivere le situazioni di dipendenza intermedie fra questi due estremi; proviamo a procedere per gradi successivi di approssimazione.

In prima approssimazione magari potremmo pensare una funzione che, anzichè passare da 0 ad 1 con uno scalino in  $x_0$ , (o da 1 a 0), passi da 0 ad 1 con un segmento di retta, con inclinazione tanto maggiore, quanto maggiore è la dipendenza di  $E[Y|X = x]$  da  $X$ , e che sia costituita quindi, complessivamente, da una spezzata come quella della figura qua sotto (figura 10):

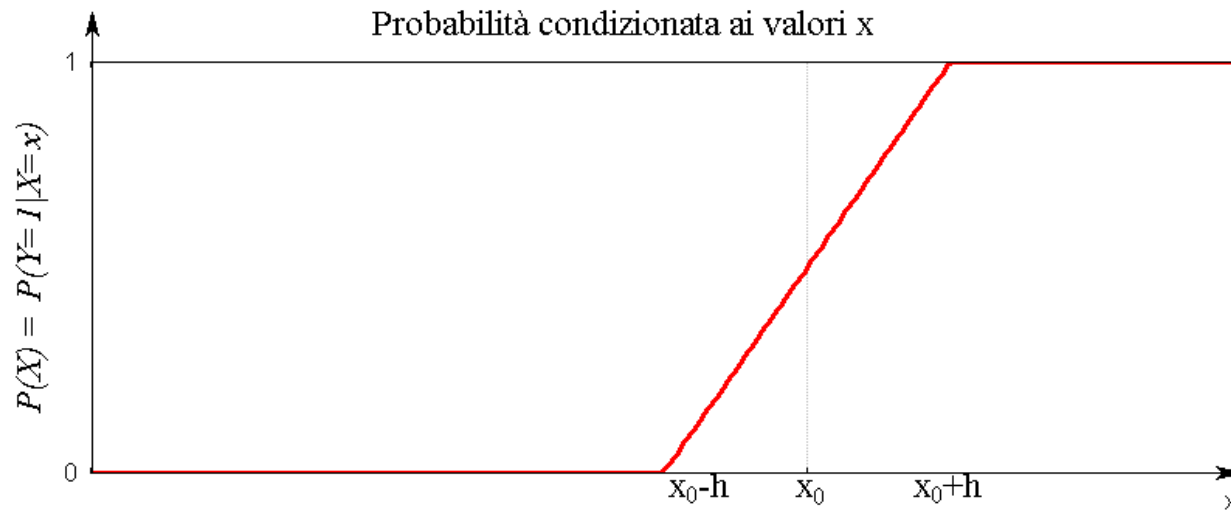


Figure 10: Probabilità funzione lineare a tratti di  $x$

$$p(x) = \begin{cases} 0 & \text{se } x < x_0 - h \\ \frac{x-x_0+h}{2h} & \text{se } x_0 - h \leq x < x_0 + h; \\ 1 & \text{se } x \geq x_0 + h \end{cases}$$

Abbiamo una probabilità nulla per valori di  $X$  inferiori a  $x_0 - h$  e una probabilità unitaria per valori di  $X$  superiori a  $x_0 + h$ , e attorno a  $x_0$  un segmento di retta di pendenza  $\frac{1}{2h}$ ; una forte dipendenza di  $\text{Prob}\{Y = 1|X = x\}$  dal particolare valore di  $x$  corrisponderà ad una forte pendenza.

Una relazione di questo genere può essere sufficiente a modellizzare situazioni semplici, ma si presta poco ad estensioni a casi con più variabili esplicative, e non può essere trattata in modo analiticamente conveniente, dal momento che, per esempio, pur essendo una funzione continua, non possiede derivata continua (l'inclinazione infatti cambia in modo discontinuo nei due punti  $x_0 - h$  e  $x_0 + h$ ). Possiamo però partire da questa prima approssimazione per ricavare delle curve continue dotate di derivate continue.

Intanto stabiliamo alcune caratteristiche che dovrebbe avere  $p(x)$ : sembra ragionevole imporre una condizione di antisimmetria attorno ad un punto sull'asse delle  $x$ , diciamo  $x_0$ :

$$p(x_0 + \delta) = 1 - p(x_0 - \delta) \quad (1)$$

che rispecchia il fatto che ragionare sul verificarsi dell'evento ( $Y = 1$ ) o sul suo non verificarsi ( $Y = 0$ ) deve portare agli stessi risultati, dato che:

$$\text{Prob} \{Y = 1 | X = x\} = p(x) \quad \text{e}$$

$$\text{Prob} \{Y = 0 | X = x\} = 1 - p(x)$$

Si deve dunque avere (dalla 1, per  $\delta = 0$ ):  $p(x_0) = \frac{1}{2}$ . Potremmo anzi definire  $x_0$  come il punto per il quale la probabilità di verificarsi dell'evento è  $\frac{1}{2}$ . In caso di associazione positiva si deve anche avere  $p'(x) > 0$ .

Inoltre appare naturale (ed analiticamente molto comodo), imporre delle condizioni sul comportamento di  $p(x)$  al divergere di  $|x|$ :

$$\lim_{x \rightarrow -\infty} p(x) = 0 \quad \text{e} \quad \lim_{x \rightarrow +\infty} p(x) = 1$$

Attorno ad  $x_0$   $p(x)$  dovrebbe essere approssimabile con una retta e dovrebbe esserci in  $x_0$  un punto di flesso, in particolare per l'ipotesi di antisimmetria, con

$p''(x_0) = 0$  e quindi  $p(x)$  dovrebbe avere il seguente sviluppo in serie di Taylor in  $x_0$  al secondo ordine:

$$\begin{aligned} p(x) &= p(x_0) + p'(x_0)(x - x_0) + \\ &\quad \frac{1}{2}p''(x_0)(x - x_0)^2 + O((x - x_0)^3) = \\ &= \frac{1}{2} + \frac{1}{2h}(x - x_0) + O((x - x_0)^3) \end{aligned}$$

in quanto in  $x_0$  si ha, per i ragionamenti fin qui fatti:  $p(x_0) = \frac{1}{2}$ ,  $p'(x_0) = \frac{1}{2h}$  e  $p''(x_0) = 0$ .

Per farla breve, una funzione<sup>3</sup> che risponde a tutte le condizioni ora viste è la seguente:

$$p(x) = \frac{\exp[\beta(x - x_0)]}{1 + \exp[\beta(x - x_0)]} \quad (2)$$

nota come *funzione logistica*.

---

<sup>3</sup>ma ovviamente non l'unica!

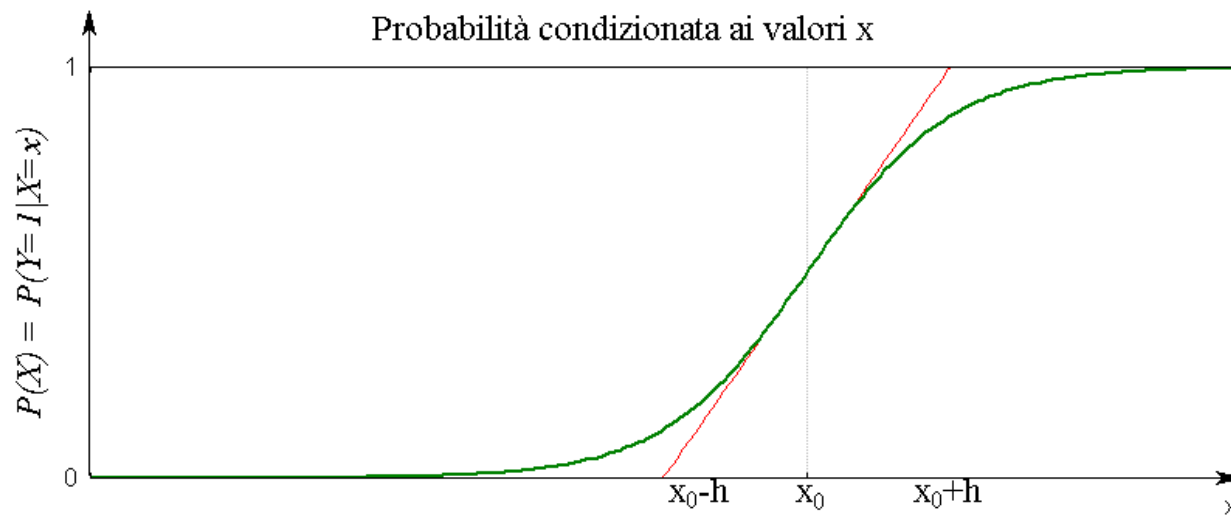


Figure 11: Funzione logistica e approssimazione lineare in  $x_0$

Più avanti farò cenno ad altre possibili scelte di funzioni per  $p(x)$ ; basta per ora segnalare che una qualsiasi funzione di ripartizione di una variabile aleatoria continua dotata di densità, con distribuzione simmetrica, unimodale e di valore atteso  $x_0$  soddisfa le nostre condizioni.

Approfondiamo lo studio della funzione logistica in (2); è facile vedere che:

$$p'(x) = \frac{\beta \exp[\beta(x - x_0)]}{(1 + \exp[\beta(x - x_0)])^2} \quad (3)$$

$$p''(x) = 0 \quad \text{per} \quad x = x_0;$$

$$p'(x_0) = \frac{\beta}{4}$$

Intanto vediamo nella figura 12 l'influenza di un diverso valore di  $\beta$  sulla forza della dipendenza di  $E[Y|X = x]$  da  $x$ .

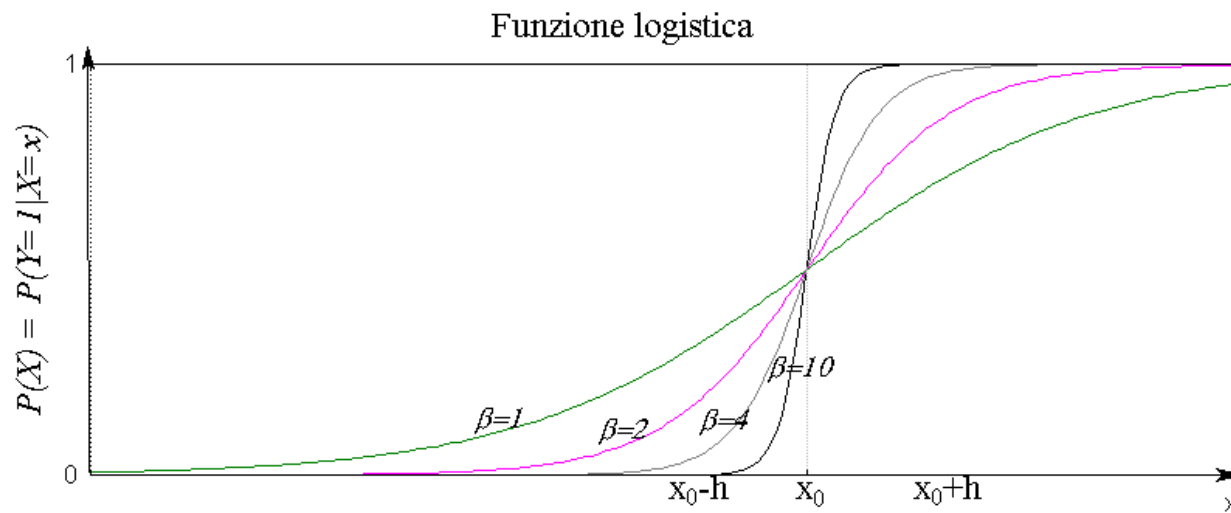


Figure 12: Funzione logistica (2) per quattro diversi valori di  $\beta$ :  $\beta = 1; 2; 4; 10$



Adesso cominciamo ad approfondire il significato di questa relazione e vediamo di dare un significato a  $\beta$  oltre quello geometrico: è facile vedere, con semplici passaggi, che dalla relazione (2) si ricava:

$$1 - p(x) = \frac{1}{1 + \exp[\beta(x - x_0)]}$$

e quindi, dividendo membro a membro la (2) e quest'ultima equazione, si ha l'importante relazione per gli *odds* (in corrispondenza di  $X = x$ ):

$$\text{odds}(x) = \frac{p(x)}{1 - p(x)} = \exp[\beta(x - x_0)]$$

e quindi:

$$\beta(x - x_0) = \log\left(\frac{p(x)}{1 - p(x)}\right) = \log \text{odds}(x)$$

Quindi il logaritmo degli odds è funzione lineare di  $x$ , adottando la funzione (2).

Infine se vogliamo misurare il rapporto fra gli odds per due valori  $x$  e  $x + \delta$ , si

ha facilmente:

$$\frac{\frac{p(x+\delta)}{1-p(x+\delta)}}{\frac{p(x)}{1-p(x)}} = \frac{\exp[\beta(x+\delta-x_0)]}{\exp[\beta(x-x_0)]} = \exp(\beta\delta)$$

che per  $\delta = 1$  diventa:

$$\frac{\frac{p(x+1)}{1-p(x+1)}}{\frac{p(x)}{1-p(x)}} = e^\beta \quad \text{e} \quad \beta = \log \left( \frac{\frac{p(x+1)}{1-p(x+1)}}{\frac{p(x)}{1-p(x)}} \right) \quad (4)$$

Quindi beta è uguale al logaritmo del rapporto degli *odds* (*odds ratio*, indicato spesso con OR) in corrispondenza di un incremento unitario di  $x$ , qualsiasi sia il valore di  $x$ .<sup>4</sup>

$$\beta = \log \text{OR} \quad (5)$$

### 3 Estensioni della regressione logistica semplice

Possiamo anche parametrizzare questa curva in modo leggermente diverso:

$$\beta(x-x_0) = \alpha + \beta x, \quad \text{con} \quad \alpha = -\beta x_0$$

---

<sup>4</sup>Si pensi al caso in cui  $X$  è una variabile che può assumere i valori 0/1, ossia è una variabile dummy associata ad una variabile qualitativa dicotoma.  $\beta$  è il logaritmo dell'odds ratio della tavola 2x2

Più in generale possiamo esprimere il legame fra  $E[Y|X = x]$ , ossia  $p(x)$ , e  $x$  in modo generale come:

$$\log \left( \frac{E[Y|X = x]}{1 - E[Y|X = x]} \right) = \eta(x) \quad (= \alpha + \beta x)$$

Chiamiamo:

- $\eta(x)$  il *predittore lineare*
- e la funzione  $g(E[Y|X = x])$  la *funzione legame* (nel nostro caso  $g(p(x)) = \log \left( \frac{p(x)}{1-p(x)} \right)$ ).

La nostra relazione di dipendenza è allora formalizzabile molto schematicamente con:

$$g(E[Y|X = x]) = \eta(x)$$

Riprenderemo più avanti questa terminologia nel contesto dei modelli lineari generalizzati. (GLM)<sup>5</sup>

Adesso la cosa importante è notare che in effetti abbiamo ancora, come nel modello lineare generale fin qui trattato, un valore atteso condizionato della variabile di

---

<sup>5</sup>si noti che se  $g(\cdot)$  fosse banalmente la funzione identità, avremmo la relazione del modello lineare

risposta, ossia  $E[Y|X = x]$ , espresso come funzione di una trasformazione lineare della variabile esplicativa  $x$ , ossia  $\eta(x)$ ; in modo molto generale:

$$E[Y|X = x] = h(\cdot);$$

ove con  $h(\cdot)$  ho indicato la funzione inversa di  $g(\cdot)$ .

$$(h(\cdot) = g^{-1}(\cdot)).$$

E' una relazione analoga a quella adottata nei modelli lineari generali ma con due differenze sostanziali:

- La distribuzione di  $Y$  non è normale (ma bernoulliana nel caso che stiamo esaminando)
- Il legame con la trasformata lineare non è diretto, ma avviene attraverso una funzione  $h(\cdot)$ , la cui inversa  $g(\cdot)$  abbiamo chiamato *funzione legame*.

Non ci preoccupiamo invece del fatto che nel modello lineare generale abbiamo trattato con  $k$  variabili esplicative: anche nel modello di regressione logistica che stiamo generalizzando possiamo pensare ad una dipendenza da una funzione lineare di più regressori; indichiamo con  $\mathbf{x}^T$  il vettore di variabili esplicative a  $k$

componenti e così il predittore lineare sarà semplicemente:

$$\eta(\mathbf{x}) = \mathbf{x}^\top \boldsymbol{\beta}$$

e avendo indicato al solito con  $\boldsymbol{\beta}$  un vettore di parametri a  $k$  componenti.

Indicando con  $\mathbf{X}$  la variabile a  $k$  componenti, costituita dalle  $k$  variabili esplicative  $X_j$ , possiamo modellizzare nel modo che segue la dipendenza del valore atteso della variabile bernoulliana  $Y$  condizionatamente ai valori assunti dal vettore di variabili esplicative  $\mathbf{X}$  :

$$\log \left( \frac{E[Y|\mathbf{X} = \mathbf{x}]}{1 - E[Y|\mathbf{X} = \mathbf{x}]} \right) = \mathbf{x}^\top \boldsymbol{\beta} \quad (= \eta(\mathbf{x}))$$

A questo punto è semplice pensare di poter modellizzare la dipendenza del valore atteso di  $Y$  da più variabili esplicative sia di tipo quantitativo che di tipo qualitativo, in funzione delle particolari configurazioni di valori osservati  $\mathbf{x}_i^\top$ .

L'interpretazione dei singoli parametri  $\beta_j$  è nuovamente quella di logaritmo dell'odds ratio relativo ad un incremento unitario della variabile  $X_j$ .

Possiamo quindi impostare modelli più o meno complessi, in funzione delle variabili che riteniamo possano influenzare la variabile di risposta (che adesso è una

probabilità del verificarsi di un evento); potremo inserire delle componenti polinomiali, delle interazioni, potremo lavorare con variabili esplicative qualitative, quantitative o miste.

L'interpretazione dei risultati e del contributo delle singole variabili esplicative sarà del tutto simile a quella vista per i modelli lineari.

#### 4 Verosimiglianza ed inferenza nella regressione logistica

Si parte dalla verosimiglianza per una osservazione da una variabile bernoulliana:

$$L(p; y) = p^y (1 - p)^{(1-y)}$$

per cui, con  $n$  osservazioni indipendenti  $y_i$ , con:

$$p_i = E[Y_i] = \frac{\exp\{\mathbf{x}_i^\top \boldsymbol{\beta}\}}{1 + \exp\{\mathbf{x}_i^\top \boldsymbol{\beta}\}}$$

$$\begin{aligned} l(\boldsymbol{\beta}; \mathbf{y}) &= \sum_{i=1}^n [y_i \log p_i + (1 - y_i) \log(1 - p_i)] = \\ &= \sum_{i=1}^n y_i \log \frac{p_i}{1 - p_i} + \sum_{i=1}^n \log(1 - p_i) = \\ &= \sum_{i=1}^n y_i \mathbf{x}_i^\top \boldsymbol{\beta} - \sum_{i=1}^n \log(1 + \exp\{\mathbf{x}_i^\top \boldsymbol{\beta}\}) = \end{aligned}$$

Decremento di verosimiglianza=  
=incremento di devianza

---

Il confronto fra modelli verrà ripreso con i GLM

-----

Bozze MARCELLO CHIODI 2020



## 5 Conteggi: Regressione di Poisson

Analogamente a quanto visto per la regressione logistica, possiamo chiederci se il valore atteso di una variabile  $Y$  espressa da un conteggio dipenda da una variabile esplicativa  $x$  (o in generale, da un predittore lineare  $\eta_i = \mathbf{x}_i^\top \boldsymbol{\beta}$ ). La nostra variabile di risposta è un conteggio, per la cui distribuzione viene spesso utilizzata la distribuzione di Poisson, di parametro  $\mu$  (nella distribuzione di Poisson si ha:  $E[Y] = V[Y] = \mu$ )

Esempi:

- Numero di pazienti affetti da una malattia in una certa località
- Numero di incidenti in funzione dell'orario e del giorno della settimana
- modellizzazione delle frequenze osservate in una tavola a più entrate
- numero di gol segnati da una squadra in funzione della quota proposta

2020

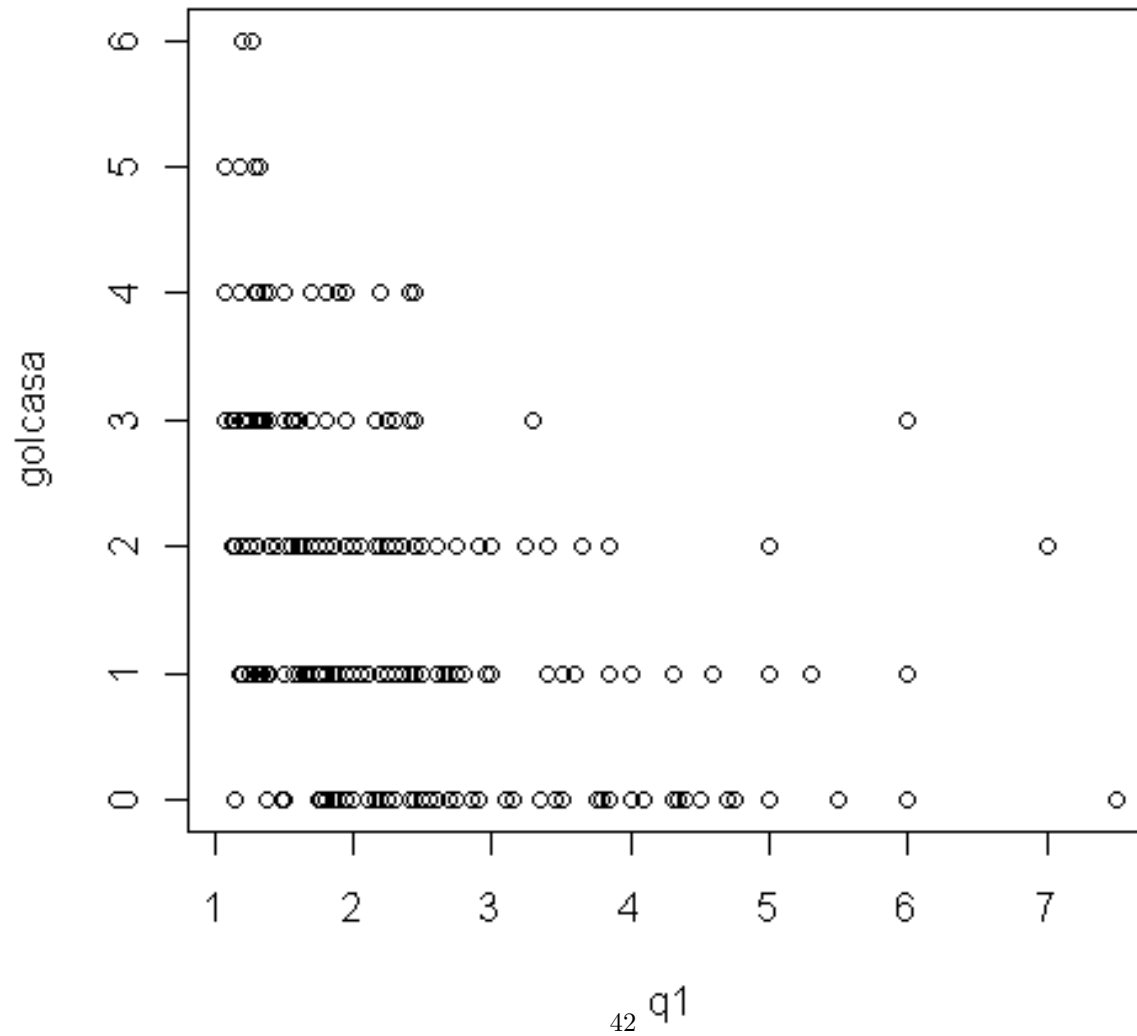


Figure 13: Risultati delle partite di un campionato di calcio: gol segnati dalla squadra di casa in funzione della quota data dagli allibratori per la vittoria della squadra di casa

Possiamo supporre che il valore atteso  $\mu_i$  dipenda dalla variabile esplicativa  $x$  (o in generale, da un predittore lineare  $\eta_i = \mathbf{x}_i^\top \boldsymbol{\beta}$ ) in modo simile a come abbiamo visto per la regressione logistica.

Per la ricerca della funzione legame possiamo partire dal legame della binomiale:

$$p_i = \frac{\mu_i}{n_i}$$

Se  $Y_i$  è distribuito secondo una Poisson, allora  $p_i$  sarà piccolo e possiamo approssimare:

$$\eta_i = \log \frac{p_i}{1 - p_i} \approx \log p_i (= \log \mu_i - \log n_i)$$

per cui possiamo ottenere ancora una relazione fra *valore atteso* e *predittore lineare*:

$$\log \mu_i = \log n_i + \eta_i$$

( $\log n_i$  è chiamato *offset*; se è costante andrà inglobato entro  $\eta_i$  )

In definitiva una soluzione soddisfacente è un legame logaritmico che porta a :

$$E[Y_i] = \mu_i = \exp(\eta_i) = \exp\{\mathbf{x}_i^\top \boldsymbol{\beta}\}$$

Il significato dei singoli parametri  $\beta_j$  parametri è immediato considerando che:

$$\begin{aligned}\frac{\partial E[Y_i]}{\partial x_{ij}} &= \frac{\partial \exp\{\mathbf{x}_i^\top \boldsymbol{\beta}\}}{\partial x_{ij}} = \\ &= \beta_j \exp\{\mathbf{x}_i^\top \boldsymbol{\beta}\}\end{aligned}$$

(Il contributo del singolo regressore alla variazione del valore atteso dipende comunque anche dagli altri regressori)

Un'interpretazione particolare dei parametri si ha per variabili esplicative dicotome:

$$\frac{E[Y_i | x_{ij} = 1]}{E[Y_i | x_{ij} = 0]} = \frac{\exp\{\mathbf{x}_{i(-j)}^\top \boldsymbol{\beta}_{(-j)} + \beta_j\}}{\exp\{\mathbf{x}_{i(-j)}^\top \boldsymbol{\beta}_{(-j)}\}} = \exp\{\beta_j\}$$

## 6 Verosimiglianza nella regressione di Poisson

Dal momento che:

$$\begin{aligned} \text{Prob}(Y = y | \mathbf{x} = \mathbf{x}_i) &= \\ \text{Prob}(Y_i = y_i) &= \frac{\mu_i^{y_i} e^{-\mu_i}}{y_i!} \end{aligned}$$

con  $\mu_i = \exp \{ \mathbf{x}_i^\top \boldsymbol{\beta} \}$ .

La log-verosimiglianza è data da:

$$l(\boldsymbol{\beta}; \mathbf{y}) = \sum_{i=1}^n \{ y_i \log \mu_i - \mu_i \}$$

rispetto al modello saturo (devianza):

$$-2(l(\boldsymbol{\beta}) - l_{max}) = 2 \sum_{i=1}^n \left\{ y_i \left[ \log \frac{y_i}{\mu_i} \right] + (\mu_i - y_i) \right\}$$

Tornando alla verosimiglianza:

$$l(\boldsymbol{\beta}; \mathbf{y}) = \sum_{i=1}^n \{ y_i \mathbf{x}_i^\top \boldsymbol{\beta} - \exp \{ \mathbf{x}_i^\top \boldsymbol{\beta} \} \}$$

$$\frac{\partial l(\boldsymbol{\beta}; \mathbf{y})}{\partial \boldsymbol{\beta}} = \sum_{i=1}^n \mathbf{x}_i \left\{ y_i - \exp \left\{ \mathbf{x}_i^T \boldsymbol{\beta} \right\} \right\}$$

Le condizioni del primo ordine di annullamento delle derivate rispetto a  $\boldsymbol{\beta}$  portano ad un sistema di equazioni non lineari in  $\boldsymbol{\beta}$  che andrà risolto con metodi numerici.

Convorrà affrontare il problema (come pure quello della regressione logistica) direttamente nell'ambito dei modelli lineari generalizzati (GLM)

## 7 Modelli lineari generalizzati

(questi appunti sono in forma molto sommaria e sono strettamente legati alle parti esposte in aula verbalmente, alla lavagna, e agli esempi in R)

### 7.1 La famiglia esponenziale

$$\log(f_Y(y; \theta, \phi)) = \frac{y \theta - b(\theta)}{a(\phi)} + c(y, \phi) \quad (6)$$

$\phi$  è un parametro di disturbo. Se non è noto la (6) potrebbe non essere una famiglia esponenziale

Se  $\phi$  è noto: la (6) è in **forma canonica**

e  $\theta$  è il **parametro naturale**

esempio alla lavagna sulla binomiale e la poisson esempio sulla distr. g

E' facile vedere che la distribuzione di Bernoulli e la distribuzione di Poisson fanno parte della famiglia esponenziale di distribuzioni:

---

Bernoulli:

- $a(\cdot) = 1$
- $b(\theta) = \log(1 + \exp(\theta))$
- $\mu = E[Y] = b'(\theta) = \frac{\exp(\theta)}{1 + \exp(\theta)}$
- $b''(\theta) = \mu(1 - \mu)$  variance function

(Per la binomiale di parametro  $n$  occorre introdurre un termine  $a = \frac{1}{n}$ )

---

Poisson:

- $a(\cdot) = 1$
- $b(\theta) = \exp(\theta)$
- $\mu = E[Y] = b'(\theta) = \exp(\theta)$
- $b''(\theta) = \mu$  variance function



Altre: Normale, Gamma, etc.

---

Per la distribuzione Gamma di parametri  $\lambda, \alpha$ , e di densità:

$$f(y) = \frac{\lambda^\alpha y^{\alpha-1} e^{-\lambda y}}{\Gamma(\alpha)}$$

converrà porre:

$$\lambda = \alpha\theta \quad \Rightarrow \quad \theta = \frac{\lambda}{\alpha} \left( = \frac{1}{E[Y]} \right),$$

in modo da avere (riparametrizzando rispetto a  $\theta, \alpha$ ):

$$f(y) = \frac{(\alpha\theta)^\alpha y^{\alpha-1} e^{-\alpha\theta y}}{\Gamma(\alpha)}$$

e quindi

$$\begin{aligned} \log f(y) &= \alpha \log \alpha\theta - \alpha\theta y + (\alpha - 1) \log y - \log \Gamma(\alpha) = \\ &= \alpha \log \theta - \alpha\theta y + \alpha \log \alpha + (\alpha - 1) \log y - \log \Gamma(\alpha) \end{aligned}$$

e raccogliendo in  $c(y, \alpha)$  tutti i termini che non contengono  $\theta$ , e mettendo poi in

evidenza il termine che contiene  $\alpha$  si ha:

$$\log f(y) = \frac{\theta y - \log \theta}{-\frac{1}{\alpha}} + c(y, \alpha)$$

per cui finalmente vediamo che anche la distribuzione Gamma fa parte della famiglia esponenziale con:

- $\phi = \alpha$  e  $a(\alpha) = -\frac{1}{\alpha}$
- $b(\theta) = \log \theta$
- $E[Y] = b'(\theta) = \frac{1}{\theta}$
- $V[Y] = b''(\theta)a(\alpha) = \frac{1}{\alpha\theta^2}$

Anche la normale rientra nella famiglia, con  $\theta = \mu$ , come è facile vedere, ponendo

:

- $\phi = \sigma^2$  e  $a(\phi) = \sigma^2$
- $b(\theta) = \frac{\theta^2}{2}$
- $E[Y] = b'(\theta) = \theta$
- $V[Y] = b''(\theta)a(\phi) = \sigma^2$

## 7.2 Valore atteso e varianza

Quando la log-densità della nostra distribuzione appartenente alla famiglia esponenziale è espressa nella forma classica vista prima (6), l'elemento essenziale è la funzione  $b(\theta)$ , che, insieme con l'eventuale funzione del parametro di disturbo  $a(\phi)$ , determina i primi due momenti della distribuzione.

Infatti dalle prime due identità di Bartlett si ha:

$$\begin{cases} \mathbb{E} \left[ \frac{\partial \log f}{\partial \theta} \right] = 0 \\ \mathbb{E} \left[ \frac{\partial^2 \log f}{\partial \theta^2} \right] = -\mathbb{E} \left[ \left( \frac{\partial \log f}{\partial \theta} \right)^2 \right] \end{cases}$$

e si ricavano i primi due momenti (derivando la 6 rispetto a  $\theta$  prima una volta e poi due volte e prendendo i valori attesi):

passaggi alla lavagna

$$\begin{cases} E[Y] = b'(\theta) \\ V[Y] = b''(\theta)a(\phi) \text{ (funzione di varianza)} \end{cases} \quad (7)$$

Se indichiamo con  $\mu_i$  il valore atteso dell' $i$ -esima osservazione:

$$\mu_i = E[Y_i] \text{ (ossia } = E[Y|X = \mathbf{x}_i])$$

dalla 7 si ricava facilmente:

$$V[Y_i] = b''(\theta_i)a(\phi) = a(\phi) \frac{\partial b'(\theta_i)}{\partial \theta_i} = a(\phi) \frac{\partial \mu_i}{\partial \theta_i}$$

e quindi anche:

$$\frac{\partial \theta_i}{\partial \mu_i} = \frac{a(\phi)}{V[Y_i]}$$

## 8 Generalized linear models (GLM)

I cosiddetti modelli lineari generalizzati sono dei modelli di dipendenza abbastanza flessibili (di cui segnalo in questi appunti solo gli aspetti essenziali) caratterizzati da alcuni punti comuni:

- Abbiamo una variabile di risposta  $Y$  il cui valore atteso dipende da  $k$  regressori (che possono anche essere delle variabili indicatrici di eventi)
- Ciascuna  $Y_i$  segue una distribuzione appartenente alla **Famiglia esponenziale**
- Definiamo un **Predittore lineare**:  $\eta_i = \mathbf{x}_i^\top \boldsymbol{\beta}$  combinazione lineare dei regressori attraverso dei parametri da stimare
- Il valore atteso di  $Y_i$  e il predittore lineare  $\eta_i$  sono collegati da una **funzione legame** (Link function) definita da:  $\eta_i = g(\mu_i) = g(\mathbb{E}[Y_i])$  ( $h(\cdot)$  è la funzione inversa di  $g(\cdot)$ , tale che:  $\mu_i = h(\eta_i)$ )
- Si ha un **legame canonico** se:  $\eta_i = \theta_i$ , essendo  $\theta$  il parametro naturale della distribuzione della famiglia esponenziale di  $Y$

- p.e.:  $g(\mu_i) = \log \frac{\mu_i}{1-\mu_i}$

legami canonici:

- Normale:  $\eta = \mu$

- Poisson:  $\eta = \log \mu$

- Binomiale:  $\eta = \log \frac{\mu}{1-\mu}$

- Gamma  $\eta = \frac{1}{\mu}$

(per la distribuzione gamma spesso si usa il legame logaritmico  $\eta = \log \mu$ , sebbene non sia canonico, perchè in questo modo si rispetta la condizione di positività del valore atteso)

## 8.1 cenni sull'inferenza nei GLM

Sono soltanto dei cenni: dal punto di vista numerico per quanto riguarda la stima dei parametri utilizziamo in questo corso la funzione `glm` di **R** ed i relativi output

E' comunque importante vedere la base dei procedimenti iterativi di stima almeno per i casi più semplici (legame canonico)

Otteniamo intanto la funzione di verosimiglianza, prendendo le mosse dalla log-densità in 6, e supponendo di avere un campione di  $n$  osservazioni indipendenti  $y_i$  da una famiglia esponenziale di parametro naturale  $\theta_i$  (ma noi poi dovremo stimare gli elementi di  $\boldsymbol{\beta}$ ):

$$l(\cdot; \mathbf{y}) = \sum_{i=1}^n \left[ \frac{y_i \theta_i - b(\theta_i)}{a(\phi)} - c(y_i, \phi) \right]$$

Deriviamo la log-verosimiglianza  $l(\cdot)$  (rispetto a  $\boldsymbol{\beta}$ ), per ottenere la score-function; procediamo per passi successivi:

$$\frac{\partial l}{\partial \theta_i} = \frac{y_i - b'(\theta_i)}{a(\phi)} = \frac{y_i - \mu_i}{a(\phi)}$$

Tratto più estesamente il caso di legame canonico ( $\mathbf{x}_i^T \boldsymbol{\beta} = \eta_i = \theta_i$ ), che porta ad

alcune semplificazioni.

$$\frac{\partial l}{\partial \boldsymbol{\beta}} = \sum_{i=1}^n \frac{\partial l_i}{\partial \theta_i} \frac{\partial \theta_i}{\partial \boldsymbol{\beta}} = \sum_{i=1}^n \frac{y_i - \mu_i}{a(\phi)} \mathbf{x}_i$$

Quindi per trovare  $\hat{\boldsymbol{\beta}}$  occorre risolvere rispetto a  $\boldsymbol{\beta}$  il sistema di  $k$  equazioni:

$$\sum_{i=1}^n \frac{y_i - \mu_i}{a(\phi)} \mathbf{x}_i = \mathbf{0} \quad (8)$$

Usualmente  $a(\phi)$  è costante; inoltre possiamo porre il sistema in forma matriciale (sempre nel caso semplificato di legame canonico)

$$\mathbf{X}^\top (\mathbf{y} - \boldsymbol{\mu}) = \mathbf{0}$$

Per la soluzione di questo sistema di equazioni si può ricorrere al metodo IRLS (Iterated reweighted Least Squares; minimi quadrati iterati riponderati); un alternativa è il metodo di newton-raphson, asintoticamente equivalente.

traccia del metodo iterativo IRLS:



Sviluppiamo prima in serie di Taylor gli scarti della formula di sopra in un intorno della soluzione  $\eta_i^*$ :

$$y_i - \mu_i \approx \frac{\partial \mu_i}{\partial \eta_i} (\eta_i^* - \eta_i)$$

e quindi (dato che nel caso del legame canonico  $\eta_i = \theta_i$ ), esplicitando rispetto a  $(\eta_i^* - \eta_i)$ :

$$\eta_i^* - \eta_i = (y_i - \mu_i) \frac{\partial \theta_i}{\partial \mu_i}$$

e ricordando che

$$\frac{\partial \mu_i}{\partial \theta_i} = \frac{V[Y_i]}{a(\phi)}$$

si ottiene:

$$y_i - \mu_i \approx (\eta_i^* - \eta_i) V[Y_i] \frac{1}{a(\phi)}$$

$$(\mathbf{y} - \boldsymbol{\mu}) = \mathbf{V}(\boldsymbol{\eta}^* - \boldsymbol{\eta}) \frac{1}{a(\phi)}$$

sostituendo questa approssimazione nel sistema di equazioni 8, (e considerando

$a(\phi)$  costante), otteniamo in termini matriciali:

$$\mathbf{X}^T(\mathbf{y} - \boldsymbol{\mu}) = \mathbf{0} \Rightarrow \mathbf{X}^T\mathbf{V}(\boldsymbol{\eta}^* - \boldsymbol{\eta}) = \mathbf{0}$$

e quindi riesprimendo rispetto a  $\boldsymbol{\beta}$

$$\mathbf{X}^T\mathbf{V}(\boldsymbol{\eta}^* - \mathbf{X}\boldsymbol{\beta}) = \mathbf{0}$$

e quindi

$$\mathbf{b} = (\mathbf{X}^T\mathbf{V}\mathbf{X})^{-1}\mathbf{X}^T\mathbf{V}\boldsymbol{\eta}^*$$

In questa espressione  $\boldsymbol{\eta}^*$  non è noto e può essere approssimato con l'approssimazione già usata:

$$\eta_i^* = \eta_i + (y_i - \mu_i) \frac{\partial \eta_i}{\partial \mu_i}$$

e in termini matriciali (sempre ricordando che siamo nel caso di legame canonico e quindi  $\boldsymbol{\eta} = \boldsymbol{\theta}$ ):

$$\boldsymbol{\eta}^* = \boldsymbol{\eta} + (\mathbf{y} - \boldsymbol{\mu})\mathbf{V}^{-1}a(\phi)$$

da cui si imposta uno schema iterativo.

8.2 proprietà asintotiche di  $\mathbf{b}$  (link canonico)

$$\begin{aligned} E[\mathbf{b}] &= \boldsymbol{\beta} \\ V[\mathbf{b}] &= a(\phi)^2 (\mathbf{X}^T \mathbf{V} \mathbf{X})^{-1} \end{aligned}$$

### 8.3 legame non canonico

(se il legame non è canonico possiamo procedere così:)

$$\frac{\partial l_i}{\partial \beta_j} = \frac{\partial l_i}{\partial \eta_i} \frac{\partial \eta_i}{\partial \beta_j} = \frac{\partial l_i}{\partial \mu_i} \frac{\partial \mu_i}{\partial \eta_i} x_{ij} =$$

$$\frac{\partial \mu_i}{\partial \eta_i} x_{ij} \frac{(y - \mu)}{V[Y]}$$

$$\frac{\partial l}{\partial \mu} = \frac{\partial l}{\partial \theta} \frac{\partial \theta}{\partial \mu} = \frac{\frac{\partial l}{\partial \theta}}{\frac{\partial \mu}{\partial \theta}} =$$

$$= \frac{(y - \mu)^{\frac{w}{\phi}}}{b''(\theta)} = \frac{(y - \mu)}{V[Y]}$$

#### 8.4 confronto fra modelli

$$\begin{aligned}
 l(\mathbf{y}; \mu(\boldsymbol{\theta}); \phi) &= \sum_{i=1}^n w_i (y_i \theta_i - b(\theta_i)) / \phi + c(y_i; \phi) \\
 \text{DEVIANZA} &= -2[l(\mathbf{y}; \mu(\boldsymbol{\theta}); \phi) - l(\mathbf{y}; \mathbf{y}; \phi)] = \\
 &= -2 \left[ \underbrace{l(\mathbf{y}; \mu(\boldsymbol{\theta}); \phi)}_{\text{log. ver. con parametro } \boldsymbol{\theta} \text{ e } \mu(\boldsymbol{\theta})} - \underbrace{l(\mathbf{y}; \mathbf{y}; \phi)}_{\text{log. ver. massima } \mu = \mathbf{y}} \right] =
 \end{aligned}$$

Confronto fra la verosimiglianza per il modello in esame e la verosimiglianza massima, relativa al **modello saturo**

si può esprimere anche come:

$$\begin{aligned}
 \text{Dev} &= \sum_{i=1}^n 2w_i (y_i (\tilde{\theta}_i - \hat{\theta}_i) - b(\tilde{\theta}_i) + b(\hat{\theta}_i)) / \phi \\
 &= D(y; \hat{\boldsymbol{\mu}}) / \phi
 \end{aligned}$$

$\tilde{\theta}_i$  stima col modello saturo  $\hat{\theta}_i$  stima col modello in esame

Per esempio:

- Normale:  $D(\mathbf{y}; \hat{\boldsymbol{\mu}}) = \sum_{i=1}^n (y_i - \hat{\mu}_i)^2$
- Poisson:  $D(\mathbf{y}; \hat{\boldsymbol{\mu}}) = 2 \sum_{i=1}^n y_i (\log(y_i / \hat{\mu}_i)) - (y_i - \hat{\mu}_i)$

I residui vengono definiti a partire degli addendi della devianza:

se  $D(\mathbf{y}; \hat{\boldsymbol{\mu}}) = \sum_{i=1}^n d_i$  allora:

Residuo =  $\text{sign}(y_i - \hat{\mu}_i) \sqrt{d_i}$

Queste quantità vengono utilizzate per validare a posteriori il modello adattato

---

Per confrontare modelli, possibilmente nidificati, si usa il rapporto delle verosimiglianze, ovviamente collegato alla differenza fra devianze.

Ha una distribuzione asintotica  $\chi_q^2$