

## OBIETTIVI

- ✓ *Introdurre il modello di regressione lineare semplice come mezzo per compiere delle previsioni su una variabile mediante un'altra*
- ✓ *Verificare la capacità di adattamento ai dati del modello di regressione lineare semplice*
- ✓ *Studiare le trappole che si possono incontrare nell'uso del modello di regressione lineare semplice*
- ✓ *Introdurre la correlazione come misura dell'associazione tra due variabili*

## Introduzione

Nei capitoli precedenti abbiamo preso in considerazione una sola variabile quantitativa: ne abbiamo studiato le misure di sintesi (Capitolo 3) ed è stata oggetto di diversi metodi di inferenza statistica volti a determinare delle stime e trarre delle conclusioni su di essa (Capitoli 5-8). Oggetto di questo capitolo è, invece, lo studio della relazione tra due variabili, e a tale scopo introdurremo due tecniche di analisi: la regressione e la correlazione.

La **regressione** ha come scopo principale la previsione: si mira, vale a dire, alla costruzione di un modello attraverso cui prevedere i valori di una **variabile dipendente** o **risposta** a partire dai valori di almeno una **variabile indipendente** o **esplicativa**. In questo capitolo studieremo la regressione lineare semplice in cui si utilizza una *sola* variabile quantitativa indipendente  $X$  per prevedere una variabile quantitativa dipendente,  $Y$ . Nel Capitolo 10, introdurremo il modello di regressione lineare *multipla*, che invece impiega *diverse* variabili esplicative ( $X_1, X_2, \dots, X_p$ ) per prevedere una variabile quantitativa  $Y$ .<sup>1</sup>

La **correlazione** ha, invece, come scopo lo studio dell'associazione tra variabili quantitative. Nel paragrafo 9.11, ad esempio, studieremo la correlazione tra il marco tedesco e lo yen giapponese per un periodo di 10 anni. L'attenzione in questo caso è focalizzata non tanto sulla possibilità di prevedere una variabile mediante un'altra, quanto sullo studio delle relazioni che possono sussistere tra due variabili quantitative.

<sup>1</sup>Tra i modelli di regressione in cui la variabile dipendente è una variabile qualitativa ricordiamo la regressione logistica (riferimento bibliografico 4)

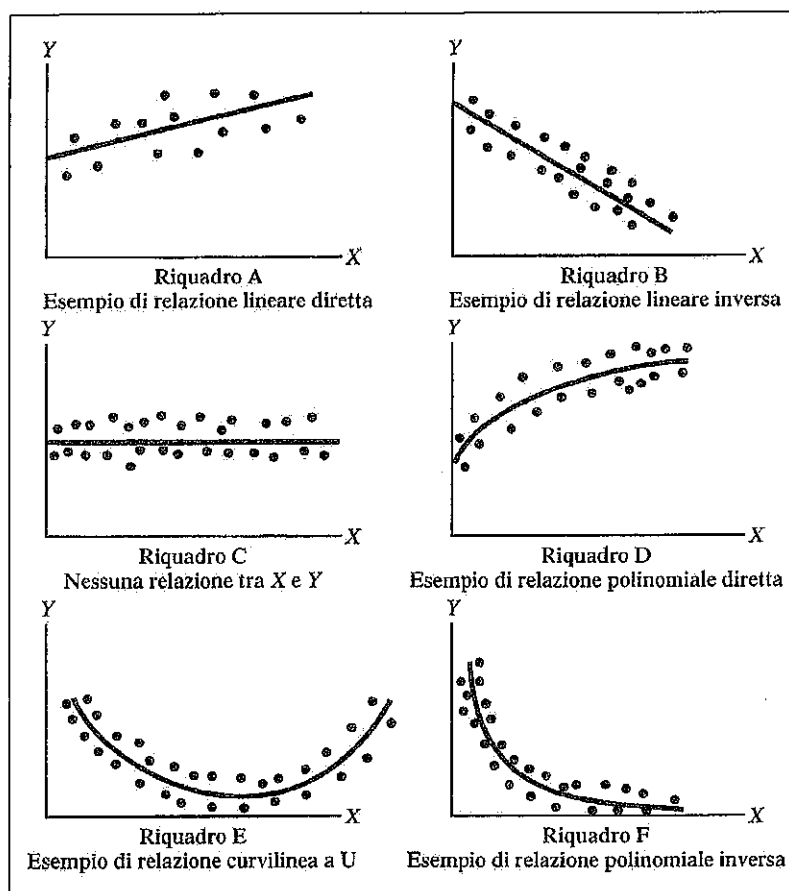
## ◆ APPLICAZIONE *Previsione delle vendite di un negozio di abbigliamento*

Nel corso degli ultimi 25 anni, una catena di negozi di abbigliamento femminile ha accresciuto la sua quota di mercato grazie all'apertura di nuove filiali. La decisione relativa alla dimensione di una nuova filiale non è mai stata oggetto di un approccio sistematico, ma quest'anno il manager responsabile del progetto di apertura di nuove filiali intende basare le sue decisioni su uno studio che gli consenta di prevedere le vendite annuali dei nuovi negozi. ◆

## ◆ 9.1 I MODELLI DI REGRESSIONE

Nel Capitolo 2 abbiamo visto come una variabile possa essere descritta facendo ricorso a diverse rappresentazioni grafiche. In maniera analoga nella regressione, per studiare la relazione tra due variabili, si fa uso di un grafico detto **diagramma di dispersione**, che si

**FIGURA 9.2**  
Esempi di relazioni  
tra variabili  
evidenziate  
dai diagrammi  
di dispersione



Nel Riquadro E si osserva una relazione parabolica o a forma di U tra X e Y: all'aumentare dei valori di X, Y diminuisce sino a un certo punto a partire dal quale comincia a crescere. Una relazione di questo genere sussiste, ad esempio, tra il numero di errori all'ora, che si compiono quando si svolge un certo lavoro, e il numero delle ore di lavoro. Il numero di errori si riduce mano a mano che si acquista dimestichezza con il lavoro, ma poi comincia a crescere come conseguenza della stanchezza o della noia.

Infine, nel Riquadro F si osserva una relazione esponenziale o polinomiale inversa tra X e Y: i valori di Y decrescono al crescere dei valori di X, ma con un tasso di decrescita che si riduce da un certo punto in poi (una volta superati alcuni valori della X). Una relazione di questo tipo sussiste tra il prezzo di vendita di un'automobile e la sua età: il prezzo di un'automobile si riduce drasticamente dopo il primo anno dall'acquisto, ma diminuisce meno rapidamente negli anni successivi.

In questo paragrafo abbiamo introdotto brevemente vari modelli per lo studio della relazione tra due variabili, individuati a mezzo del diagramma di dispersione. Sebbene i diagrammi di dispersione siano dei validi strumenti di analisi, tecniche statistiche più sofisticate consentono di pervenire alla scelta del modello di regressione più appropriato. Nei paragrafi successivi ci concentreremo sul modello di regressione lineare.

## Il metodo dei minimi quadrati

Nel paragrafo precedente abbiamo proposto un modello statistico (il modello di regressione semplice) per studiare la relazione tra due variabili quantitative (la dimensione dei negozi e l'ammontare annuo delle vendite) alla luce dei soli dati campionari. Si dimostra che sotto determinate ipotesi (paragrafo 9.4) l'intercetta campionaria  $b_0$  e l'inclinazione campionaria  $b_1$  si possono usare come stimatori dei parametri della popolazione  $\beta_0$  e  $\beta_1$ : si ottiene in questa maniera la forma campionaria del modello di regressione lineare semplice.

### L'equazione campionaria del modello di regressione lineare

La previsione di  $Y$  in base al modello di regressione lineare è data dalla somma tra l'intercetta campionaria e il prodotto tra il valore di  $X$  e l'inclinazione campionaria

$$\hat{Y}_i = b_0 + b_1 X_i \quad (9.2)$$

dove

$\hat{Y}_i$  = previsione di  $Y$  per l'osservazione  $i$

$X_i$  = valore di  $X$  per l'osservazione  $i$

La previsione di  $Y$  richiede, allora, il calcolo dei due coefficienti di regressione  $b_0$  e  $b_1$ . Una volta determinati, possiamo tracciare la retta di regressione nel diagramma di dispersione e valutare visivamente la capacità esplicativa del modello, osservando se la retta stimata si avvicina o meno ai dati osservati.

La regressione mira a individuare la retta che meglio si adatta ai dati, laddove tale capacità di adattamento può essere valutata in base a criteri diversi. Il criterio più semplice consiste nel valutare le differenze tra i valori osservati ( $Y_i$ ) e i valori previsti in base alla retta stimata ( $\hat{Y}_i$ ) e quindi cercare quella retta che minimizza tali differenze. Tuttavia, siccome le differenze considerate possono essere negative per alcune osservazioni e positive per altre, si tratterà di *minimizzare* la somma dei loro quadrati:

$$\sum_{i=1}^n (Y_i - \hat{Y}_i)^2$$

dove

$Y_i$  = il vero valore di  $Y$  per l'osservazione di  $i$

$\hat{Y}_i$  = il valore previsto di  $Y$  per l'osservazione di  $i$

Dal momento che in base al modello proposto  $\hat{Y}_i = b_0 + b_1 X_i$ , si tratta di minimizzare la seguente espressione:

$$\sum_{i=1}^n [Y_i - (b_0 + b_1 X_i)]^2$$

rispetto alle due incognite  $b_0$  e  $b_1$ .

La tecnica matematica che consiste nel determinare i valori di  $b_0$  e  $b_1$  che rendono minima l'espressione precedente prende il nome di **metodo dei minimi quadrati**.

Nella Figura 9.4 si riporta l'output di Excel relativo al dataset della Tabella 9.1. I valori stimati dei **coefficienti di regressione** sono rispettivamente  $b_0 = 901,247$  e  $b_1 = 1,686$  e quindi la retta stimata ha la seguente espressione:

$$\hat{Y}_i = 901,247 + 1,686 X_i$$

### Esempio 9.1 Interpretazione dell'intercetta $b_0$ e dell'inclinazione $b_1$

Un economista intende usare il tasso di crescita annuo della produttività negli Stati Uniti ( $X$ ) per prevedere la variazione percentuale dell'indice Standard & Poor di 500 azioni ( $Y$ ). Sulla base di dati annui per periodo di 50 anni, l'economista ottiene la seguente stima della retta di regressione di  $Y$  su  $X$ :

$$\hat{Y}_i = -5.0 + 7X_i$$

Come si possono interpretare l'intercetta  $b_0$  e l'inclinazione  $b_1$ ?

#### SOLUZIONE

Alla luce del valore assunto dall'intercetta,  $b_0 = -5$ , ci aspettiamo, in base al modello proposto, che l'indice Standard & Poor diminuisca del 5% se il tasso di produttività è uguale a 0. L'inclinazione  $b_1 = 7$  ci dice che in corrispondenza di una variazione della produttività dell'1% dovremmo osservare una variazione dell'indice del 7%.

### Esempio 9.2 Previsione delle vendite annue di un negozio a partire dalla sua dimensione

Prevedete in base al modello di regressione stimato (tabulato Excel della Figura 9.4), l'ammontare delle vendite per un negozio di 4000 piedi al quadrato.

#### SOLUZIONE

Si tratta di sostituire il valore 4000 al valore di  $X_i$  nella retta di regressione stimata:

$$\hat{Y}_i = 901.247 + 1.686X_i \quad \text{per } 3800$$

$$\hat{Y}_i = 901.247 + 1.686(4000) = 7645.786 \text{ o } \$ 7645.786$$

Pertanto si prevede un ammontare delle vendite pari a \$ 7645.786 per un negozio di 4000 piedi al quadrato.

### La previsione nel modello di regressione: l'interpolazione a confronto con l'estrapolazione

Nel fare previsioni alla luce di un modello di regressione dobbiamo sempre tenere presente quale è l'intervallo dei valori assunti dalla variabile indipendente e questo perché possiamo prevedere  $Y$  solo in corrispondenza dei valori di  $X$  che cadono in questo intervallo. Quando facciamo previsioni in corrispondenza a valori di  $X$  che cadono in questo intervallo, diciamo che stiamo facendo delle *interpolazioni*, mentre quando cerchiamo di prevedere il valore di  $Y$  corrispondente a valori di  $X$  che non cadono in questo intervallo, stiamo cercando di fare delle *estrapolazioni*. Per chiarire meglio la differenza tra le due operazioni, consideriamo i dati nella Tabella 9.1, e supponiamo di utilizzare la superficie del negozio per prevedere l'ammontare delle vendite; i valori osservati con riferimento alla variabile esplicativa vanno da 1102 a 5841 piedi al quadrato. Di conseguenza, potremo fare previsioni sulle vendite di negozi con una superficie compresa in questo intervallo. Qualunque previsione sulle vendite di negozi con superficie al di fuori di questo intervallo può essere fatta soltanto sotto l'ipotesi che la relazione stimata tra le due variabili rimanga la stessa anche al di fuori di questo intervallo (e questa è un'ipotesi che si basa su nostre supposizioni e non su dati osservati).

importanti

MOVIE	BOX OFFICE GROSS (\$ MILLIONI)	HOME VIDEO UNITS SOLD (000)	MOVIE	BOX OFFICE GROSS (\$ MILLIONI)	HOME VIDEO UNITS SOLD (000)
1	1.10	57.18	16	9.36	190.80
2	1.13	26.17	17	9.89	121.57
3	1.18	92.79	18	12.66	183.30
4	1.25	61.60	19	15.35	204.72
5	1.44	46.50	20	17.55	112.47
6	1.53	85.06	21	17.91	162.95
7	1.53	103.52	22	18.25	109.20
8	1.69	30.88	23	23.13	280.79
9	1.74	49.29	24	27.62	229.51
10	1.77	24.14	25	37.09	277.68
11	2.42	115.31	26	40.73	226.73
12	5.34	87.04	27	45.55	365.14
13	5.70	128.45	28	46.62	218.64
14	6.43	126.64	29	54.70	286.31
15	8.59	107.28	30	58.51	254.58



**DATASET  
RENT**

**9.4** Un agente immobiliare intende prevedere gli affitti mensili degli appartamenti sulla base della loro dimensione. Nella tabella seguente si riportano i valori degli affitti e della dimensione di 25 appartamenti di una zona residenziale.

APARTMENT	MONTHLY RENT (\$)	SIZE (SQUARE FEET)	APARTMENT	MONTHLY RENT (\$)	SIZE (SQUARE FEET)
1	950	850	14	1800	1369
2	1600	1450	15	1400	1175
3	1200	1085	16	1450	1225
4	1500	1232	17	1100	1245
5	950	718	18	1700	1259
6	1700	1485	19	1200	1150
7	1650	1136	20	1150	896
8	935	726	21	1600	1361
9	875	700	22	1650	1040
10	1150	956	23	1200	755
11	1400	1100	24	800	1000
12	1650	1285	25	1750	1200
13	2300	1985			

- Create il diagramma di dispersione per i dati della tabella.
- Stimate con il metodo dei minimi quadrati i coefficienti di regressione  $b_0$  e  $b_1$ .
- Fornite un'interpretazione di  $b_0$  e  $b_1$  con riferimento al problema considerato.
- Prevedete l'ammontare dell'affitto mensile per un appartamento di 1000 piedi al quadrato.
- Perché non si può utilizzare la retta stimata per prevedere l'affitto mensile di appartamenti di 500 piedi al quadrato?

**9.3**

**LE MISURE DI VARIABILITÀ**

In questo paragrafo introduciamo alcune misure di variabilità che consentono di valutare le capacità previsive del modello statistico proposto. La somma totale dei quadrati (SQT)

### La somma dei quadrati degli errori (SQE)

La somma dei quadrati degli errori (SQE) è data dalla somma dei quadrati delle differenze tra i valori osservati e i valori previsti di  $Y$

$$SQE = \text{variabilità non spiegata} = \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 \quad (9.6)$$

In base alle Figure 9.4 e 9.5, osserviamo

$$SQR = 106,208,120; \quad SQE = 10,532,255; \quad \text{e} \quad SQT = 116,740,375$$

In base alle Figure 9.4 e 9.5, osserviamo che:

$$SQR = 106\,208\,120; \quad SQE = 10\,532\,255; \quad SQT = 116\,740\,375$$

Inoltre, in base all'equazione (9.3):

$$SQT = SQR + SQE$$

$$116\,740\,375 = 106\,208\,120 + 10\,532\,255$$

La somma dei quadrati delle differenze dalla media, 116 740 375, si suddivide nella somma dei quadrati spiegata dalla regressione 106 208 120 e nella somma dei quadrati residua, 10 532 255.

### COMMENTO: Notazione scientifica

In alcune versioni di Excel, le cifre molto piccole o molto grandi sono espresse non in formato numerico, ma facendo ricorso alla "notazione scientifica". Ad esempio potremmo trovare la  $SQR$  dell'esempio precedente ( $SQR = 106208120$ ) formattata come 1.06E+08. I numeri che seguono la lettera E corrispondono al numero di cifre decimali per le quali si deve spostare la virgola verso sinistra (se negativi) o verso destra (se positivi) per ottenere l'abituale formato numerico. 3.7431E-02 è il numero che si ottiene da 3.7431 spostando la virgola di 2 posti verso sinistra 0.037431 e mentre 3.7431E+02 è il numero che si ottiene da 3.7431 spostando la virgola di 2 posti verso destra, 374.31. Pertanto 1.06E+08 non è altro che 106000000. Osservate come nella notazione scientifica si faccia ricorso a poche cifre significative, cosa che comporta un'approssimazione dei numeri considerati.

### Il coefficiente di determinazione

Le somme dei quadrati precedentemente introdotte ( $SQT$ ,  $SQR$  e  $SQE$ ) forniscono, se considerate da sole, informazioni limitate sulla bontà del modello statistico proposto. Tuttavia il rapporto tra la somma dei quadrati della regressione e la somma totale dei quadrati si configura come una misura utile per valutare il modello di regressione.

Tale misura prende il nome di **coefficiente di determinazione** ed è di seguito definita.

### Il coefficiente di determinazione

Il coefficiente di determinazione è dato dal rapporto tra la somma dei quadrati della regressione e la somma totale dei quadrati.

$$r^2 = \frac{SQR}{SQT} \quad (9.7)$$

Il coefficiente di determinazione misura la parte di variabilità di  $Y$  spiegata dalla variabile indipendente  $X$  nel modello di regressione. Tornando all'esempio riguardante la



DATASET  
MOVIE



DATASET  
RENT

- 9.8** In base all'output di Excel ottenuto per risolvere l'esercizio 9.3:
- calcolate il coefficiente di determinazione  $r^2$  e spiegate il significato;
  - calcolate l'errore standard della stima;
  - ritenete che il modello di regressione sia un utile strumento di previsione delle vendite di videocassette?
- 9.9** In base all'output di Excel ottenuto per risolvere l'esercizio 9.4:
- calcolate il coefficiente di determinazione  $r^2$  e spiegate il significato;
  - calcolate l'errore standard della stima;
  - ritenete che il modello di regressione sia un utile strumento di previsione degli affitti mensili?

## 9.4

### LE ASSUNZIONI DEL MODELLO

Quando abbiamo introdotto la teoria della verifica di ipotesi e dell'analisi della varianza, abbiamo più volte sottolineato come una corretta applicazione delle procedure statistiche dipenda in genere dal soddisfacimento delle ipotesi su cui esse si fondano. Le assunzioni alla base del modello di regressione sono analoghe a quelle su cui si fonda l'analisi della varianza, perché i due modelli di analisi ricadono nello stesso insieme: la classe dei modelli *lineari* (riferimento bibliografico 6).

Nel riquadro 9.1 si riportano le ipotesi del modello di regressione.

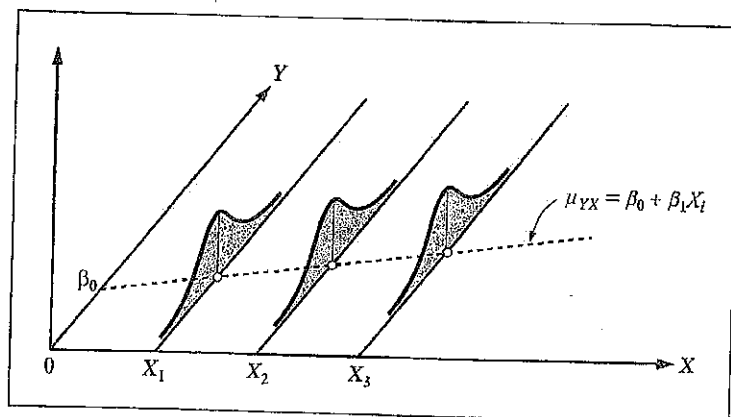
#### Riquadro 9.1 Le ipotesi del modello di regressione

- ✓ 1. Distribuzione normale degli errori.
- ✓ 2. Omoschedasticità.
- ✓ 3. Indipendenza degli errori.

In base alla prima ipotesi, la **normalità**, si richiede che gli errori abbiano, per ogni valore di  $X$ , una distribuzione normale (Figura 9.7). Analogamente al test  $t$  e al test  $F$  dell'A-NOVA, il modello di regressione risulta robusto rispetto a scostamenti dall'ipotesi di normalità: le inferenze sulla retta di regressione e sui coefficienti non risultano seriamente influenzate da una distribuzione degli errori solo approssimativamente normale.

**FIGURA 9.7**

Le ipotesi del modello di regressione



La seconda ipotesi, l'**omoschedasticità**, richiede che la variabilità degli errori sia costante per ciascun valore di  $X$ . Gli errori devono, vale a dire, variare di un medesimo ammontare sia in corrispondenza di valori elevati, che in corrispondenza di valori piccoli di  $X$ . L'omoschedasticità degli errori è cruciale ai fini dell'applicazione del metodo dei

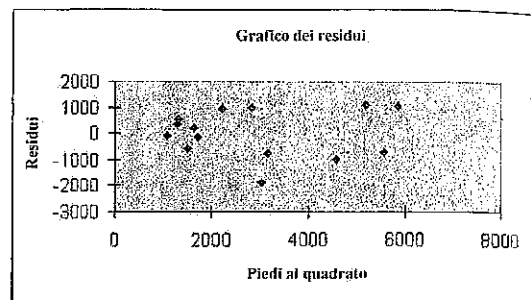
grafico dei residui  $e_i$  rispetto a  $X_i$ . Con il grafico dei residui si rimuove il trend lineare di  $Y$  rispetto a  $X$  e si evidenzia quindi lo scarso adattamento ai dati del modello di regressione (*lack of fit*). Possiamo quindi concludere che un modello polinomiale in questo caso è più appropriato di un modello lineare semplice.

Torniamo all'esempio relativo alla catena di negozi di abbigliamento. Nella Figura 9.9 si riporta l'output di Excel relativo all'analisi dei residui, contenente le previsioni della variabile risposta (l'ammontare delle vendite annue) e i corrispondenti residui.

Nella Figura 9.10, si riporta il grafico dei residui rispetto ai valori della variabile indipendente (la dimensione dei negozi). Il grafico non rivela alcuna relazione particolare tra i residui e  $X_i$ ; i residui sembrano distribuirsi in maniera uguale al di sopra e al di sotto dello 0. Possiamo concludere che, in questo caso, il modello di regressione è adeguato.

Osservazione	Vendite previste	Residui
1	3811,515528	-130,516
2	3669,880191	225,1198
3	5649,402647	1003,597
4	19267,72633	-724,726
5	3079,732952	330,267
6	4624,232585	938,7674
7	3115,141786	544,8582
8	2759,367307	-65,3673
9	8214,257862	-746,258
10	3457,427185	-559,427
11	9603,309152	1070,611
12	8601,82498	-1016,82
13	18749,96093	1010,039
14	5973,140501	-1060,14

**FIGURA 9.9**  
L'output di Excel relativo all'analisi dei residui per il problema di scelta della dimensione delle nuove filiali



**FIGURA 9.10**  
Grafico dei residui rispetto alla dimensione dei negozi ottenuto con Excel

### Valutazione delle ipotesi

◆ **Omoschedasticità** Il grafico dei residui rispetto a  $X_i$  consente di stabilire anche se l'ipotesi di omoschedasticità è soddisfatta. Nel grafico della Figura 9.11, la variabilità dei residui varia a seconda dei valori assunti da  $X$ , segno di una violazione dell'ipotesi di omoschedasticità; i residui sembrano disporsi a ventaglio, a indicare un aumento della variabilità all'aumentare dei valori di  $X$ . Nella Figura 9.10, invece, non si osservano delle differenze significative nella variabilità dei residui in corrispondenza di valori diversi di  $X$ : l'ipotesi di omoschedasticità sembra soddisfatta.

◆ **Normalità** L'analisi dei residui consente di verificare l'ipotesi di normalità degli errori. Si tratta di costruire la distribuzione di frequenza dei residui e di darne una rappresentazione grafica a mezzo dell'istogramma.

Considerando ancora l'esempio relativo alla catena di negozi di abbigliamento, nella Tabella 9.2 si riporta la distribuzione di frequenza dei residui, mentre il corrispondente istogramma è rappresentato nella Figura 9.12.

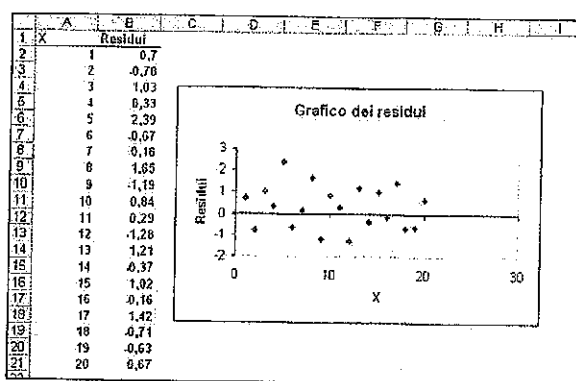
È chiaramente, difficile, valutare l'ipotesi di normalità della distribuzione degli errori alla luce di un campione di sole 14 osservazioni e questo indipendentemente dallo strumento a cui si ricorre (l'istogramma, il diagramma ramo-foglia, il diagramma scatola e baffi o il *normality plot*). Dalla Figura 9.12, possiamo solo osservare che la distribuzione dei dati, sebbene non sembri normale, non è particolarmente asimmetrica. Pertanto la robustezza della regressione rispetto agli



◆ **Indipendenza** L'ipotesi di indipendenza degli errori può essere verificata rappresentando i residui nell'ordine con cui i dati sono stati raccolti: nei dati raccolti nel corso del tempo è spesso presente un'*autocorrelazione* tra osservazioni successive. Il grafico dei residui rispetto al tempo consente allora di evidenziare la presenza di una relazione di questo genere. L'autocorrelazione dei residui viene misurata a mezzo della statistica di Durbin-Watson di cui ci occuperemo nel paragrafo 9.6.

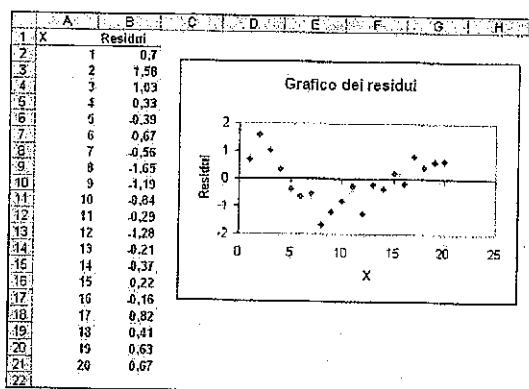
### Esercizi del paragrafo 9.5

- **9.10** Di seguito si riporta la tabella contenente i valori di  $X$  e dei residui ottenuti dalla stima di un modello di regressione, e il relativo grafico dei residui.



È riconoscibile una struttura nei residui? Commentate.

- **9.11** Di seguito si riporta la tabella contenente i valori di  $X$  e dei residui ottenuti dalla stima di un modello di regressione, e il relativo grafico dei residui.



È riconoscibile una struttura nei residui? Commentate.

- **9.12** Tornate all'esercizio 9.2 e conducete un'analisi dei residui. Alla luce dei risultati ottenuti:
  - valutate la capacità di adattamento ai dati del modello;
  - verificate se le ipotesi alla base del modello di regressione sono soddisfatte.



	A	B	C	D	E	F	G	
1	Analisi di regressione per il negozio di consegna a domicilio							
2								
3	<b>Statistica della regressione</b>							
4	R multiplo	0,810629997						
5	R al quadrato	0,657445284						
6	R al quadrato corretto	0,631094922						
7	Errore standard	0,936036681						
8	Osservazioni	15						
9								
10	<b>ANALISI VARIANZA</b>							
11		gdf	SQ	MQ	F	Significatività F		
12	Regressione	1	21,860433	21,86043264	24,95014171	0,000245105		
13	Residuo	13	11,390141	0,876164669				
14	Totale	14	33,250573					
15								
16		Coefficiente		errore standa.	Stat t	Valore di significatività	Inferiore 95%	Superiore 95%
17	Intorcutta	-16,0321936	5,3101671	-3,019150493	0,009868641	27,50410993	-4,560277262	
18	Customers	0,030760228	0,0061582	4,995011683	0,000245105	0,017456271	0,044064185	

FIGURA 9.13 Output di Excel relativo ai dati della Tabella 9.3

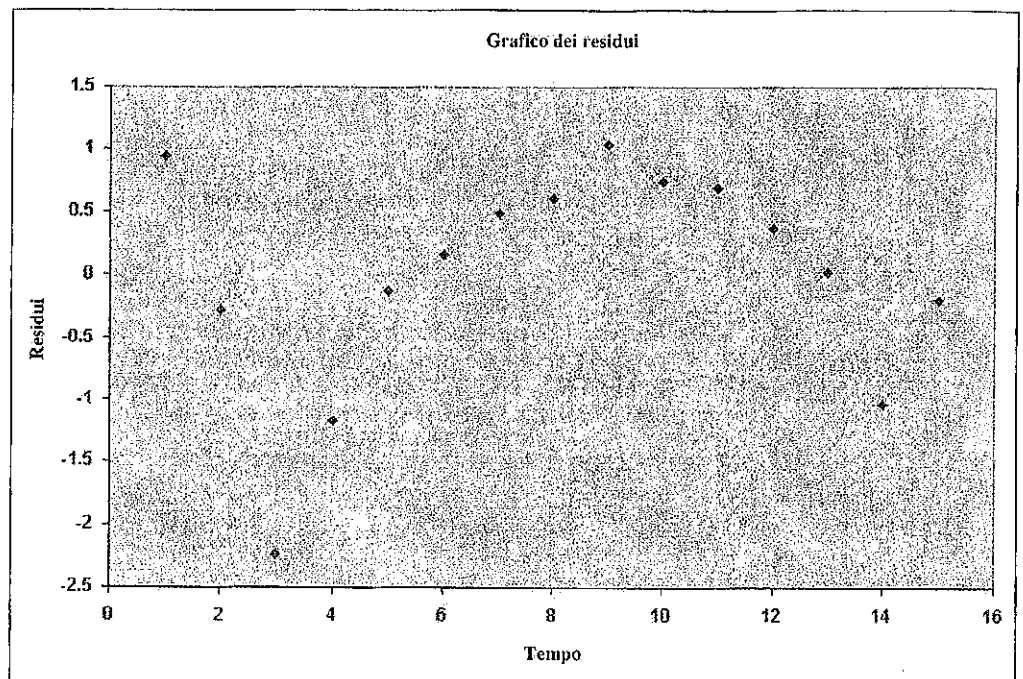


FIGURA 9.14 Grafico dei residui relativo ai dati della Tabella 9.3

il basso, che fa pensare alla presenza di una certa autocorrelazione tra di essi e quindi alla violazione dell'ipotesi di indipendenza degli errori.

### La statistica di Durbin-Watson

L'autocorrelazione dei residui può essere individuata e misurata facendo ricorso a una particolare statistica campionaria, la **statistica di Durbin-Watson**, che misura la correlazione tra ciascun residuo e quello che lo precede.

$D < d_L$   
 presenza di autocorrelazione  
 $D > d_U$  non vi è prova  
 della presenza di  
 autocorrelazione

Nella Tabella 9.4 si riportano, per ciascuna combinazione di  $\alpha$  (il livello di significatività),  $n$  (la dimensione del campione) e  $p$  (il numero delle variabili indipendenti), due valori della statistica  $D$ . Il primo  $d_L$  è il valore critico inferiore di  $D$ , cioè il più piccolo dei valori di  $D$  corrispondenti a una situazione di assenza di autocorrelazione dei residui: se  $D$  è minore di  $d_L$ , concludiamo che vi è prova della presenza di un'autocorrelazione positiva tra i residui. Il secondo valore  $d_U$  è, invece, il valore critico superiore di  $D$ : se  $D$  è maggiore di  $d_U$ , concludiamo che non vi è prova della presenza di un'autocorrelazione positiva tra i residui. Se  $D$  è compreso tra  $d_L$  e  $d_U$  non possiamo giungere a nessuna conclusione.

Pertanto per i dati della Tabella 9.3, siccome vi è una sola variabile esplicativa ( $p = 1$ ) e si considerano 15 osservazioni ( $n = 15$ ), avremo in base alla Tabella 9.4,  $d_L = 1.08$  e  $d_U = 1.36$ . Siccome  $D = 0.883 < 1.08$ , possiamo concludere che vi è autocorrelazione tra i residui. L'analisi di regressione condotta risulta inappropriata e si rende necessario ricorrere a approcci diversi all'analisi della relazione tra variabili (riferimento bibliografico 6).

### Esercizi del paragrafo 9.6

- **9.15** La tabella seguente riporta i valori dei residui per dati raccolti nel corso di 10 periodi di tempo:

PERIODO DI TEMPO	RESIDUI	PERIODO DI TEMPO	RESIDUI
1	-5	6	+1
2	-4	7	+2
3	-3	8	+3
4	-2	9	+4
5	-1	10	+5

- Rappresentate graficamente i residui rispetto al tempo. Commentate.
- Calcolate la statistica di Durbin-Watson.
- Alla luce delle risposte ai punti (a) e (b) a quali conclusioni si può pervenire in merito all'autocorrelazione dei residui?


- **9.16** Riprendete l'esercizio 9.2.

- È necessario in questo caso calcolare la statistica di Durbin-Watson?
- In che caso sarebbe necessario calcolare la statistica di Durbin-Watson, prima di procedere alla stima con il metodo dei minimi quadrati del modello di regressione?

- **9.17** Il proprietario di una casa monofamiliare, in una zona residenziale nel Nord-Est degli Stati Uniti, intende fare ricorso a un modello statistico per prevedere il consumo complessivo di elettricità sulla base della temperatura atmosferica (misurata in gradi Fahrenheit, °F). Nella tabella seguente si riportano i valori dei chilowatt consumati e della temperatura relativi a un periodo di 24 mesi.

- Disegnate il diagramma di dispersione per i dati della tabella.
- Nell'ipotesi che tra le due variabili sussista una relazione lineare, stimate con il metodo dei minimi quadrati i coefficienti di regressione  $b_0$  e  $b_1$ .
- Fornite un'interpretazione di  $b_1$ .
- Prevedete il consumo medio in chilowatt corrispondente a una temperatura media di 50 gradi Fahrenheit.
- Calcolate il coefficiente di determinazione  $r^2$  e interpretatene il significato.
- Calcolate l'errore standard della stima.
- Rappresentate graficamente i residui rispetto alla temperatura atmosferica media.
- Rappresentate graficamente i residui rispetto al tempo.
- Calcolate la statistica di Durbin-Watson. Per  $\alpha = 0.05$ , si può ritenere che vi sia autocorrelazione tra i residui?

 DATASET  
PETFOOD

 DATASET  
ELECUSE

<sup>2</sup>Dettagli ulteriori sul calcolo della statistica  $t$  sono dati nel paragrafo 9.10

annuo delle vendite, per un livello di significatività uguale a 0.05. Dall'output di Excel della Figura 9.4, ricaviamo<sup>2</sup>:

$$b_1 = +1.686 \quad n = 14 \quad S_{b_1} = 0.1533$$

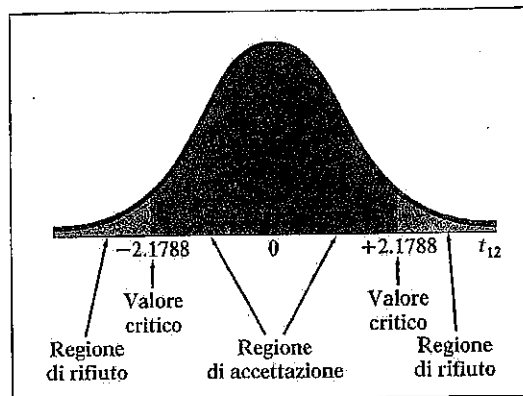
Pertanto il valore della statistica  $t$  è dato da:

$$t = \frac{b_1}{S_{b_1}} = \frac{1.686}{0.1533} = 11.00$$

Il valore della statistica  $t$  è riportato nell'output di Excel nella colonna **Stat  $t$** . Poiché  $t = 11 > t_{12} = 2.1788$ , rifiutiamo  $H_0$ . In base all'approccio del  $p$ -value rifiutiamo  $H_0$  perché il  $p$ -value è approssimativamente uguale a 0. (Nell'output di Excel il valore del  $p$ -value è riportato nella colonna **Significatività** con la notazione scientifica  $1.27E - 07$ , che corrisponde a 0.000000127 e pertanto è inferiore a  $\alpha = 0.05$ ). Possiamo, quindi, concludere che esiste una relazione lineare significativa tra l'ammontare medio annuo delle vendite e la dimensione del negozio. (Figura 9.16).

**FIGURA 9.16**

Verifica di ipotesi sull'inclinazione della retta di regressione, con  $\alpha = 0.05$  e 12 gradi di libertà



### Il test $F$ per l'inclinazione

La significatività dell'inclinazione della retta di regressione può essere sottoposta a verifica anche ricorrendo al test  $F$  (Tabella 9.5). Nel paragrafo 7.2 abbiamo visto che il test  $F$  ha per oggetto il rapporto tra due varianze. Nel verificare la significatività dell'inclinazione, si impiega come misura dell'errore casuale la varianza dei residui (data dalla somma dei quadrati degli errori divisa per il numero dei gradi di libertà). Pertanto il test  $F$  è dato dal rapporto tra la varianza dovuta alla regressione (la somma dei quadrati della regressione divisa per il numero delle variabili indipendenti) e la varianza dei residui (equazione 9.12).

#### Il test $F$ per la verifica di ipotesi sull'inclinazione $\beta_1$

La statistica  $F$  è data dal rapporto tra la media dei quadrati della regressione ( $MQR$ ) e la media dei quadrati dell'errore ( $MQE$ ):

$$F = \frac{MQR}{MQE} \quad (9.12)$$

dove

$$MQR = \frac{SQR}{p}$$

$$MQE = \frac{SQE}{n - p - 1}$$

### L'intervallo di confidenza per l'inclinazione

L'intervallo di confidenza per  $\beta_1$  si ottiene addizionando e sottraendo all'inclinazione campionaria  $b_1$  il prodotto tra il valore critico della statistica  $t$  e l'errore standard dell'inclinazione.

$$b_1 \pm t_{n-2} S_{b_1} \quad (9.13)$$

Dall'output di Excel della Figura 9.4, ricaviamo:

$$b_1 = +1.686 \quad n = 14 \quad S_{b_1} = 0.1533$$

Pertanto

$$\begin{aligned} b_1 \pm t_{n-2} S_{b_1} &= +1.686 \pm (2.1788)(0.1533) \\ &= +1.686 \pm 0.334 \\ &+1.352 \leq \beta_1 \leq +2.02 \end{aligned}$$

L'intervallo di confidenza per l'inclinazione per un livello di confidenza del 95% è compreso tra gli estremi +1.352 e +2.02. Poiché 0 non è compreso nell'intervallo, concludiamo che l'ammontare delle vendite è legato da una relazione lineare significativa alla dimensione del negozio.

### Esercizi del paragrafo 9.7

- **9.19** Intendete verificare la significatività dell'inclinazione della retta di regressione. A tale scopo estraete un campione di ampiezza  $n = 18$  e ottenete i seguenti risultati:

$$b_1 = +4.5 \quad S_{b_1} = 1.5$$

- (a) Calcolate il valore della statistica  $t$ .
- (b) Determinate i valori critici della statistica per  $\alpha = 0.05$ .
- (c) Alla luce dei punti (a) e (b), a quale decisione statistica dovrete pervenire?
- (d) Costruite un intervallo di confidenza di livello 0.95 per l'inclinazione  $\beta_1$ .
- **9.20** Riprendete l'output di Excel ottenuto per risolvere l'esercizio 9.2.
  - (a) Per  $\alpha = 0.05$ , si può ritenere che tra l'ammontare delle vendite di cibo per animali e lo spazio destinato sugli scaffali al prodotto sussista una relazione lineare significativa?
  - (b) Costruite un intervallo di confidenza di livello 95% per l'inclinazione  $\beta_1$ .
- **9.21** Riprendete l'output di Excel ottenuto per risolvere l'esercizio 9.3.
  - (a) Per  $\alpha = 0.05$ , si può ritenere che tra il successo di botteghino e le vendite di videocassette sussista una relazione lineare significativa?
  - (b) Costruite un intervallo di confidenza di livello 95% per l'inclinazione  $\beta_1$ .
- **9.22** Riprendete l'output di Excel ottenuto per risolvere l'esercizio 9.4.
  - (a) Per  $\alpha = 0.05$ , si può ritenere che tra la dimensione degli appartamenti e l'ammontare degli affitti sussista una relazione lineare significativa?
  - (b) Costruite un intervallo di confidenza di livello 95% per l'inclinazione  $\beta_1$ .
- **9.23** La volatilità delle azioni viene spesso misurata facendo ricorso all'indice beta. Tale indice si ottiene stimando un modello di regressione in cui la variabile dipendente è la variazione settimanale dell'azione e la variabile indipendente è la variazione settimanale di un indice di mercato. In genere l'indice di mercato utilizzato è l'indice S&P 500. Pertanto se si volesse calcolare l'indice beta per le azioni dell'IBM, dovremmo stimare il seguente modello di regressione:

$$\begin{aligned} &(\text{variazione \% settimanale di IBM}) = \\ &\beta_0 + \beta_1(\text{variazione \% settimanale di S\&P 500}) + \epsilon \end{aligned}$$



DATASET  
PETFOOD



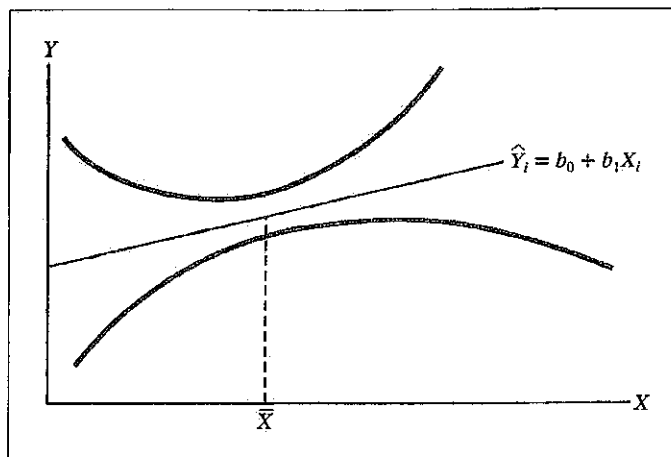
DATASET  
MOVIE



DATASET  
RENT

L'ampiezza dell'intervallo è piccola in prossimità del valor medio  $\bar{X}$  e aumenta mano a mano che ci si allontana da questo. Tale effetto è dovuto alla presenza della radice quadrata del rapporto nell'equazione (9.14) ed è illustrato dalla Figura 9.18.

**FIGURA 9.18**  
Intervallo di confidenza  
per  $\mu_{YX}$



**Esempio 9.3** *Costruzione di intervallo di confidenza del 95% per la risposta media  $\mu_{YX}$*

Per il problema introdotto nell'Applicazione, relativo alla catena di negozi di abbigliamento, abbiamo ottenuto la seguente stima della retta di regressione  $\hat{Y}_i = 901.247 + 1.686 X_i$ . Costruite un intervallo di confidenza del 95% per la media delle vendite di tutti i negozi di 4000 piedi al quadrato.

**SOLUZIONE**

Dalla stima della retta di regressione:

$$\hat{Y}_i = 901.247 + 1.686 X_i$$

e per  $X_i = 4,000$ , otteniamo

$$\hat{Y}_i = 901.247 + 1.686(4000) = 7645.786$$

Inoltre

$$\bar{X} = 2921.2857; \quad S_{YX} = 936.85; \quad \sum_{i=1}^n (X_i - \bar{X})^2 = 37\,357\,090.86$$

e dalla Tavola E.3,  $t_{12} = 2.1788$ . Pertanto,

$$\hat{Y}_i \pm t_{n-2} S_{YX} \sqrt{h_i}$$

dove

$$h_i = \frac{1}{n} + \frac{(X_i - \bar{X})^2}{\sum_{i=1}^n (X_i - \bar{X})^2}$$

di modo che

$$\hat{Y}_i \pm t_{n-2} S_{YX} \sqrt{\frac{1}{n} + \frac{(X_i - \bar{X})^2}{\sum_{i=1}^n (X_i - \bar{X})^2}}$$

dove:

$$h_i = \frac{1}{n} + \frac{(X_i - \bar{X})^2}{\sum_{i=1}^n (X_i - \bar{X})^2}$$

di modo che

$$\hat{Y}_i \pm t_{n-2} S_{YX} \sqrt{1 + \frac{1}{n} + \frac{(X_i - \bar{X})^2}{\sum_{i=1}^n (X_i - \bar{X})^2}}$$

e

$$7645.786 \pm (2.1788)(936.85) \sqrt{1 + \frac{1}{14} + \frac{(4,000 - 2921.2857)^2}{37\,357\,090.86}}$$
$$= 7645.786 \pm 2143.357$$

quindi

$$5502.43 \leq Y_i \leq 9789.143$$

Pertanto possiamo concludere che l'ammontare delle vendite settimanali per un negozio di 4000 piedi al quadrato è compreso tra 5502.43 e 9789.143 migliaia di dollari.

L'esempio 9.4 illustra il calcolo dell'intervallo di confidenza per la previsione di una singola risposta con riferimento ai dati relativi alla catena di negozi di abbigliamento, anche se tale intervallo può essere ottenuto facilmente con l'aggiunta PHStat.

Confrontando i risultati degli esempi 9.3 e 9.4, osserviamo che l'ampiezza dell'intervallo per la previsione di una singola risposta è maggiore dell'ampiezza dell'intervallo della risposta media e questo perché nella previsione di un valore singolo vi è senz'altro una variabilità maggiore di quella che accompagna la previsione di un valore medio.

### Esercizi del paragrafo 9.8

**9.24** Per un campione di 20 osservazioni, si ottiene con il metodo dei minimi quadrati la seguente retta di regressione:  $\hat{Y}_i = 5 + 3X_i$ , inoltre,  $S_{YX} = 1.0$ ,  $\bar{X} = 2$ , e  $\sum_{i=1}^n (X_i - \bar{X})^2 = 20$ .

- Costruite un intervallo di confidenza del 95% per la risposta media corrispondente a  $X = 2$ .
- Costruite un intervallo di confidenza del 95% per il valore della singola risposta corrispondente a  $X = 2$ .

• **9.25** Riprendete l'output di Excel ottenuto per risolvere l'esercizio 9.2.

- Costruite un intervallo di confidenza del 95% per la media delle vendite settimanali di tutti i negozi in cui i cibi per animali occupano uno spazio di 8 piedi.
- Costruite un intervallo di confidenza del 95% per le vendite settimanali di un negozio in cui i cibi per animali occupano uno spazio di 8 piedi.
- Spiegate le differenze tra i risultati dei due punti.

**9.26** Riprendete l'output di Excel ottenuto per risolvere l'esercizio 9.3.

- Costruite un intervallo di confidenza del 95% per le vendite medie delle videocassette di tutti i film per i quali sono stati acquistati biglietti per 10 milioni di dollari.
- Costruite un intervallo di confidenza del 95% per le vendite delle videocassette di un film di cui sono stati acquistati biglietti per 10 milioni di dollari.
- Spiegate le differenze tra i risultati dei due punti.



DATASET  
PETFOOD

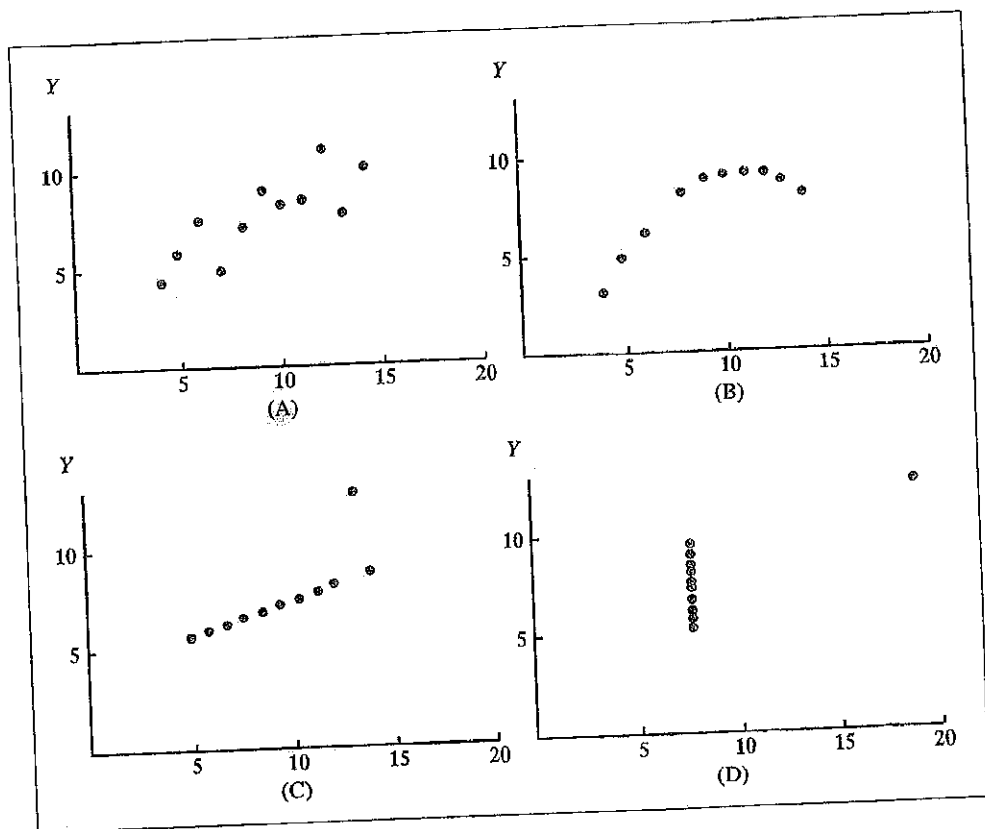


DATASET  
MOVIE

**Tabella 9.6** *Quattro insiemi di dati artificiali*

DATASET A		DATASET B		DATASET C		DATASET D	
$X_i$	$Y_i$	$X_i$	$Y_i$	$X_i$	$Y_i$	$X_i$	$Y_i$
10	8.04	10	9.14	10	7.46	8	6.58
14	9.96	14	8.10	14	8.84	8	5.76
5	5.68	5	4.74	5	5.73	8	7.71
8	6.95	8	8.14	8	6.77	8	8.84
9	8.81	9	8.77	9	7.11	8	8.47
12	10.84	12	9.13	12	8.15	8	7.04
4	4.26	4	3.10	4	5.39	8	5.25
7	4.82	7	7.26	7	6.42	19	12.50
11	8.33	11	9.26	11	7.81	8	5.56
13	7.58	13	8.74	13	12.74	8	7.91
6	7.24	6	6.13	6	6.08	8	6.89

Fonte: "Graphs in Statistical Analysis," F.J. Anscombe,  
*The American Statistician* 27 (1973); 17-21.  
 Copyright © 1973 The American Statistical Association.



**FIGURA 9.19** Diagrammi di dispersione per i quattro dataset



Quest'esempio sottolinea l'importanza dell'analisi dei residui, che si deve sempre accompagnare alla stima del modello. Nel riquadro 9.3 si suggerisce una strategia da seguire per evitare le trappole dell'analisi di regressione evidenziate nel Riquadro 9.2.



### Riquadro 9.3 Una strategia per evitare le trappole della regressione

- ✓ 1. Cominciate l'analisi sempre con un'attenta osservazione del diagramma di dispersione, per cogliere l'eventuale relazione tra  $X$  e  $Y$ .
- ✓ 2. Verificate se le ipotesi alla base del modello di regressione sono soddisfatte dopo la stima del modello e prima di passare a impiegare i risultati.
- ✓ 3. Rappresentate graficamente i residui rispetto alla variabile dipendente per stabilire se il modello si adatta ai dati e se l'ipotesi di omoschedasticità è rispettata.
- ✓ 4. Usate l'istogramma, il diagramma ramo-foglia o il diagramma scatola e baffi dei residui per verificare in quale misura l'ipotesi di normalità degli errori è rispettata.
- ✓ 5. Se i dati sono raccolti in ordine sequenziale, rappresentate graficamente i residui nell'ordine con cui i dati sono stati raccolti e calcolate la statistica di Durbin-Watson.
- ✓ 6. Se alla luce dei punti 3-5 ritenete che le ipotesi alla base del modello di regressione lineare siano violate, ricorrete ad altri metodi di stima del modello o ad altri modelli.
- ✓ 7. Se alla luce dei punti 3-5 ritenete che le ipotesi alla base del modello di regressione lineare non siano violate, potete procedere ad alcune inferenze sul modello. Sottoponete a verifica la significatività dei coefficienti e costruite gli intervalli di confidenza per la risposta media e per la previsione.

## 9.10

### I CALCOLI DELLA REGRESSIONE LINEARE SEMPLICE

In questo paragrafo illustriamo i calcoli che consentono di ottenere le statistiche del modello di regressione prese in considerazione.

#### Il calcolo dell'intercetta $b_0$ e dell'inclinazione $b_1$

Il metodo dei minimi quadrati per la stima dei coefficienti della retta di regressione comporta la risoluzione del sistema dato dalle seguenti equazioni (9.16a) e (9.16b).

**Equazioni da risolvere per applicare il metodo dei minimi quadrati**

$$\sum_{i=1}^n Y_i = nb_0 + b_1 \sum_{i=1}^n X_i \quad (9.16a)$$

$$\sum_{i=1}^n X_i Y_i = b_0 \sum_{i=1}^n X_i + b_1 \sum_{i=1}^n X_i^2 \quad (9.16b)$$

**Tabella 9.7** Calcoli per il problema della dimensione dei negozi

NEGOZIO	PIEDI		$X^2$	$Y^2$	XY
	AL QUADRATO $X$	VENDITE $Y$			
1	1,726	3,681	2,979,076	13,549,761	6,353,406
2	1,642	3,895	2,696,164	15,171,025	6,395,590
3	2,816	6,653	7,929,856	44,262,409	18,734,848
4	5,555	9,543	30,858,025	91,068,849	53,011,365
5	1,292	3,418	1,669,264	11,682,724	4,416,056
6	2,208	5,563	4,875,264	30,946,969	12,283,104
7	1,313	3,660	1,723,969	13,395,600	4,805,580
8	1,102	2,694	1,214,404	7,257,636	2,968,788
9	3,151	5,468	9,928,801	29,899,024	17,229,668
10	1,516	2,898	2,298,256	8,398,404	4,393,368
11	5,161	10,674	26,635,921	113,934,276	55,088,514
12	4,567	7,585	20,857,489	57,532,225	34,640,695
13	5,841	11,760	34,117,281	138,297,600	68,690,160
14	<u>3,008</u>	<u>4,085</u>	<u>9,048,064</u>	<u>16,687,225</u>	<u>12,287,680</u>
Totale	40,898	81,577	156,831,834	592,083,727	301,298,822

$$\begin{aligned}
 SQX &= \sum_{i=1}^n (X_i - \bar{X})^2 = \sum_{i=1}^n X_i^2 - \frac{\left(\sum_{i=1}^n X_i\right)^2}{n} \\
 &= 156\,831\,834 - \frac{(40\,898)^2}{14} \\
 &= 156\,831\,834 - 119\,474\,743.1 \\
 &= 37\,357\,090.86
 \end{aligned}$$

di modo che

$$\begin{aligned}
 b_1 &= \frac{62\,989\,097.29}{37\,357\,090.86} \\
 &= 1.68613
 \end{aligned}$$

e

$$b_0 = \bar{Y} - b_1 \bar{X}$$

e

$$\begin{aligned}
 \bar{Y} &= \frac{\sum_{i=1}^n Y_i}{n} = \frac{81\,577}{14} = 5826.929 \\
 \bar{X} &= \frac{\sum_{i=1}^n X_i}{n} = \frac{40\,898}{14} = 2921.2857
 \end{aligned}$$

di modo che

$$\begin{aligned}
 b_0 &= 5,826.929 - (1.68613)(2,921.2857) \\
 &= 901.2
 \end{aligned}$$

$SQE =$  variabilità residua

$$\begin{aligned} &= \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 \\ &= \sum_{i=1}^n Y_i^2 - b_0 \sum_{i=1}^n Y_i - b_1 \sum_{i=1}^n X_i Y_i \\ &= 592\,083\,727 - (901.2)(81\,577) - (1.68613)(301\,298\,822) \\ &= 10\,532\,255 \end{aligned}$$

### Il calcolo dell'errore standard dell'inclinazione

In questo paragrafo illustriamo i calcoli che consentono di ottenere l'errore standard dell'inclinazione, che abbiamo impiegato nella verifica di ipotesi sull'esistenza di una relazione lineare tra  $X$  e  $Y$ .

$$\begin{aligned} S_{b_1} &= \frac{S_{YX}}{\sqrt{SQX}} \\ SQX &= \sum_{i=1}^n (X_i - \bar{X})^2 \\ &= \sum_{i=1}^n X_i^2 - \frac{\left(\sum_{i=1}^n X_i\right)^2}{n} \\ &= 156\,831\,834 - \frac{(40\,898)^2}{14} \\ &= 37\,357\,090.86 \\ S_{b_1} &= \frac{936.85}{\sqrt{37\,357\,090.86}} \\ &= 0.1533 \end{aligned}$$

L'esempio 9.5 presenta un'ulteriore illustrazione dei calcoli necessari per la stima del modello di regressione.

### Esempio 9.5 Il calcolo di $b_0$ , $b_1$ , $SQT$ , $SQR$ , $SQE$ e $r^2$

Prendete in considerazione il dataset A della Tabella 9.6:

$X_i$	$Y_i$
10	8.04
14	9.96
5	5.68
8	6.95
9	8.81
12	10.84
4	4.26
7	4.82
11	8.33
13	7.58
6	7.24

Calcolate  $b_0$ ,  $b_1$ ,  $SQT$ ,  $SQR$ ,  $SQE$  e  $r^2$ .