

# Metodi Quantitativi per Economia, Finanza e Management

## *Lezione n°8*

Regressione lineare multipla: le ipotesi del modello, la stima del modello

# Il modello di regressione lineare

1. Introduzione ai modelli di regressione
2. Obiettivi
3. Le ipotesi del modello
4. La stima del modello
5. La valutazione del modello
6. Commenti

# L'obiettivo dell'analisi

Prevedere la redditività  
del cliente al tempo  $t+1$

# L'impostazione del problema

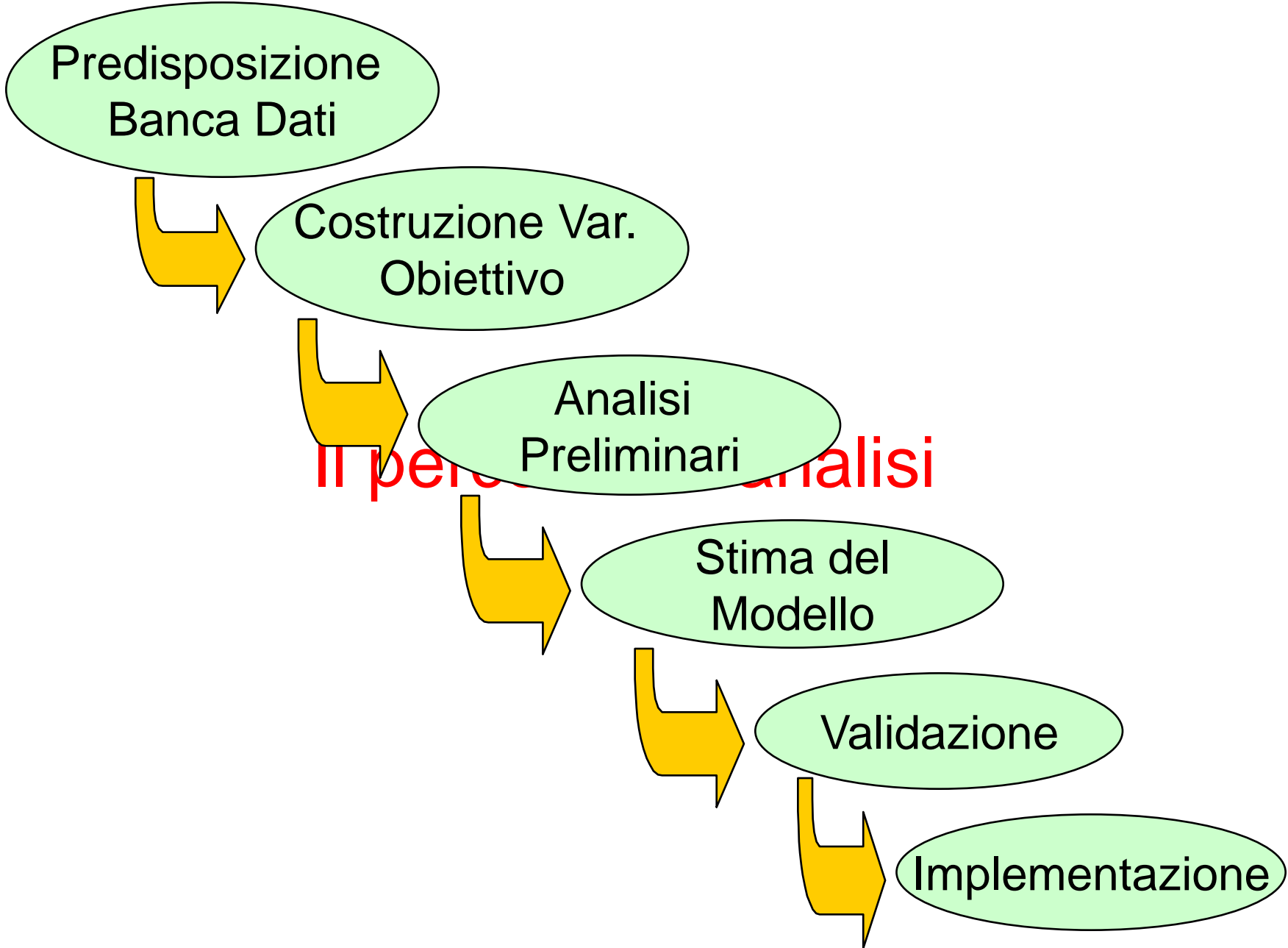
$$\underline{\text{Redditività} = \text{ricavi} - \text{costi}}$$

- ✦ redditività var. continua
- ✦ classi di redditività ( $< 0$  ;  $\geq 0$ )

# I dati di input



- ◆  $Y$  : Redditività consolidata al tempo 't'
- ◆  $\underline{X}$  : # ordini al tempo 't-1,-2,..'  
pagato ordini al tempo 't-1,-2,..'  
pagato rateale mensile ...  
sesso (dicotomica)  
area (dicotomiche)  
.....



Predisposizione  
Banca Dati

Costruzione Var.  
Obiettivo

Analisi  
Preliminari

Stima del  
Modello

Validazione

Implementazione

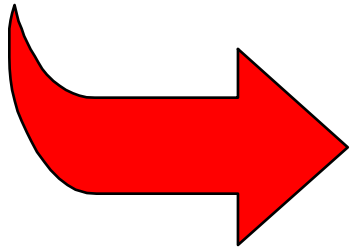
Il percorso dell'analisi

# Analisi preliminari

- ✦ lo studio della distribuzione di  $Y$
- ✦ la struttura di correlazione tra  $Y$  e  $\underline{X}$

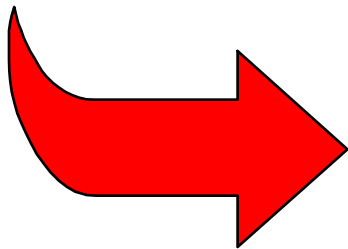
# L'impostazione del problema

✦ Redditività var. continua



Regressione Lineare

✦ Redditività var. dicotomica



Regressione Logistica



# Il modello di regressione lineare

1. Introduzione ai modelli di regressione
2. Obiettivi
3. Le ipotesi del modello
4. La stima del modello
5. La valutazione del modello
6. Commenti

# I modelli di regressione

Modelli di dipendenza per la rappresentazione di relazioni non simmetriche tra le variabili

- $Y$  “variabile dipendente” (variabile target da spiegare)
- $X_1, \dots, X_p$  “variabili indipendenti” (variabili esplicative o regressori)

# Il modello di regressione lineare

Si vuole descrivere la relazione tra  $Y$  e  $X_1, \dots, X_p$  con una funzione lineare

- se  $p=1 \rightarrow$  osservazioni in uno spazio a due dimensioni  
( $i=1, \dots, n$ )

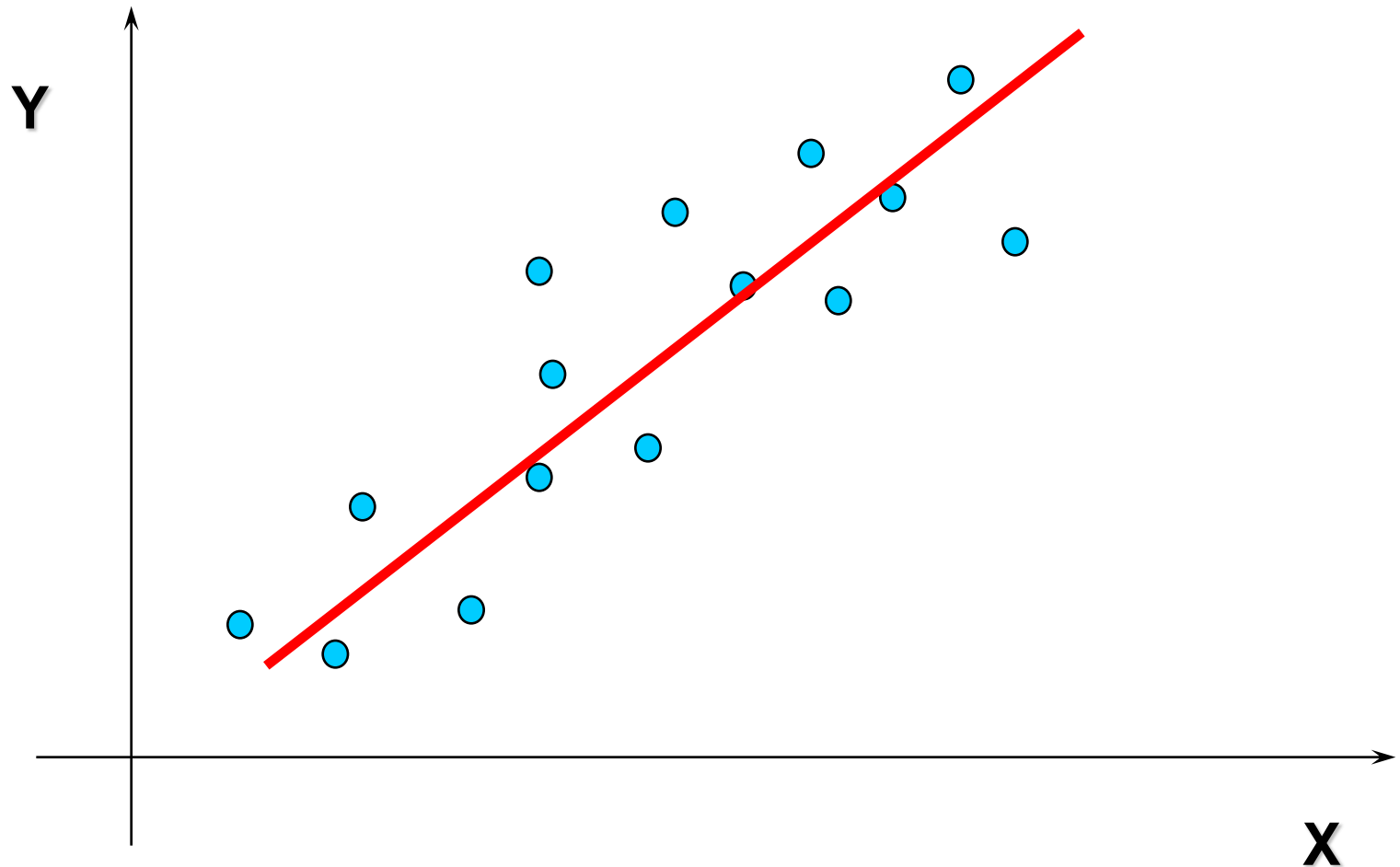
$$Y_i = f(X_{i1})$$

- se  $p>1 \rightarrow$  osservazioni in uno spazio a  $p+1$  dimensioni  
( $i=1, \dots, n$ )

$$Y_i = g(X_{i1}, \dots, X_{ip})$$

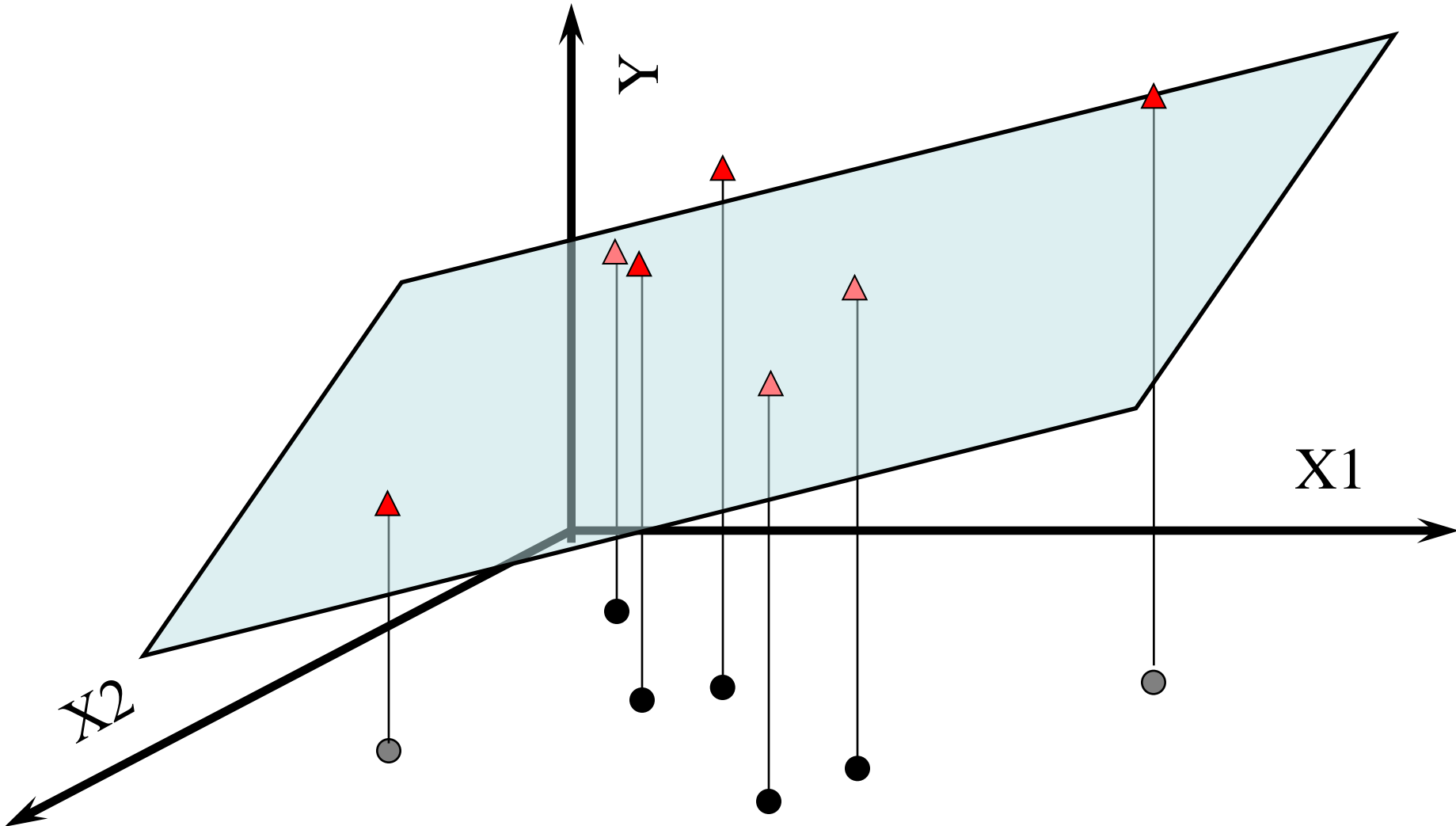
# Il modello di regressione lineare

- se  $p=1 \rightarrow$  spazio a due dimensioni  $\rightarrow$  retta di regressione lineare semplice



# Il modello di regressione lineare

- se  $p > 1 \rightarrow$  spazio a  $p+1$  dimensioni  $\rightarrow$  “retta” di regressione lineare multipla



# Il modello di regressione lineare

## Obiettivi

- **Esplicativo** - Stimare l'influenza dei regressori sulla variabile target.
- **Predittivo** - Stimare il valore non osservato della variabile target in corrispondenza di valori osservati dei regressori.
- **Comparativo** - Confrontare la capacità di più regressori, o di più set di regressori, di influenzare il target (= confronto tra modelli di regressione lineare diversi).

# Il modello di regressione lineare

## Le ipotesi del modello

<u>Y</u>	<u>X<sub>1</sub></u>	<u>X<sub>2</sub></u>	<u>X<sub>3</sub></u>	...	...	...	<u>X<sub>p</sub></u>
$y_1$	$X_{11}$	$X_{12}$	$X_{13}$	...	...	...	$X_{1p}$
$y_2$	$X_{21}$	$X_{22}$	$X_{23}$	...	...	...	$X_{2p}$
$y_3$	$X_{31}$	$X_{32}$	$X_{33}$	...	...	...	$X_{3p}$
...	...	...	...	...	...	...	...
...	...	...	...	...	...	...	...
...	...	...	...	...	...	...	...
$y_n$	$X_{n1}$	$X_{n2}$	$X_{n3}$	...	...	...	$X_{np}$

(nx1) (nxp)

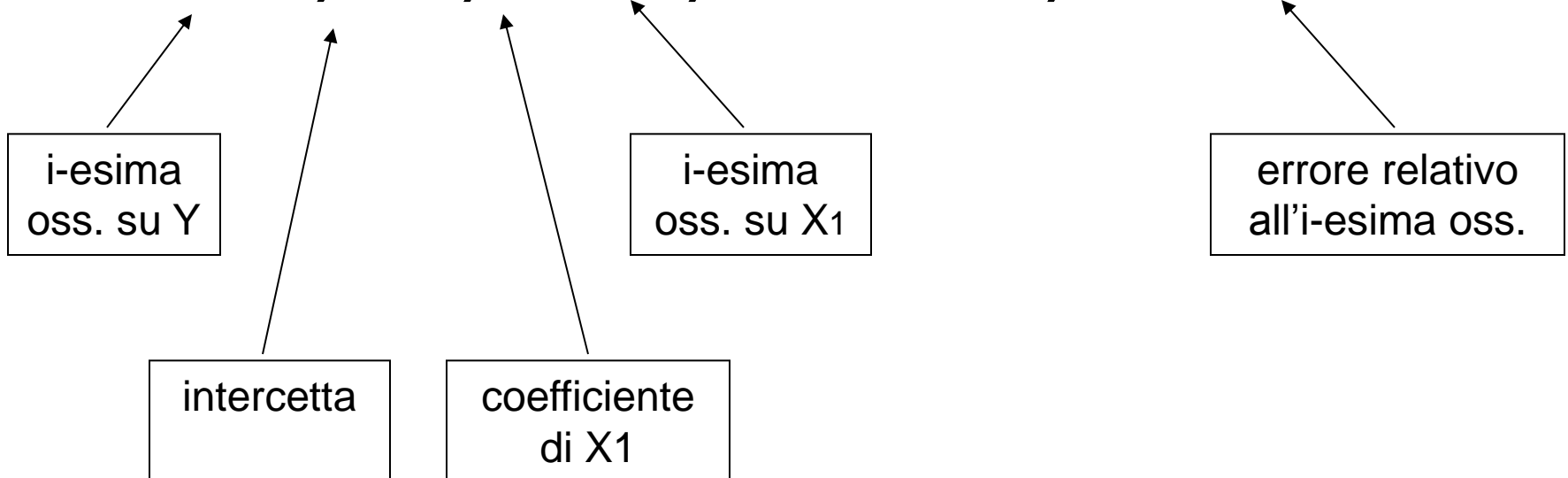
- n unità statistiche
- vettore colonna (nx1) di n misurazioni su una variabile continua (Y)
- matrice (nxp) di n misurazioni su p variabili quantitative ( $X_1, \dots, X_p$ )
- la singola osservazione è il vettore riga ( $y_i, X_{i1}, X_{i2}, X_{i3}, \dots, X_{ip}$ )  
 $i=1, \dots, n$

# Il modello di regressione lineare

## Le ipotesi del modello

Equazione di regressione lineare multipla

$$Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_p X_{ip} + \varepsilon_i$$



La matrice  $X=[1, X_1, \dots, X_p]$  è detta matrice del disegno.



# Il modello di regressione lineare

## Le ipotesi del modello

L'errore presente nel modello si ipotizza essere di natura casuale. Può essere determinato da:

- variabili non considerate
- problemi di misurazione
- modello inadeguato
- effetti puramente casuali

# Il modello di regressione lineare

## Le ipotesi del modello

1. Errori a media nulla

$$E(\varepsilon) = 0$$

2. Errori con varianza costante  
(omoschedasticità)

$$Cov(\varepsilon) = \sigma^2 I_n$$

3. Errori non correlati  
(per ogni  $i \neq j$ )

$$Cov(\varepsilon_i, \varepsilon_j) = 0$$

4. Errori con distribuzione Normale

$$\varepsilon \sim N(0, \sigma^2 I_n)$$

\* 1 – 3  $\rightarrow$  hp deboli

1 – 4  $\rightarrow$  hp forti

# Il modello di regressione lineare

## Le ipotesi del modello

Da un punto di vista statistico

- $Y$  è un vettore aleatorio di cui si osserva una specifica realizzazione campionaria  $\rightarrow$  hp sulla distribuzione
- $X$  è una matrice costante con valore noto  $\rightarrow$  no hp sulla distribuzione
- $\beta$  è un vettore costante non noto
- l'errore è un vettore aleatorio di cui si osserva una specifica realizzazione campionaria  $\rightarrow$  hp sulla distribuzione

# Il modello di regressione lineare

## Le ipotesi del modello

- in media  $Y$  può essere rappresentata come funzione lineare delle sole  $(X_1, \dots, X_p)$

$$\mu = E(Y) = X\beta$$

- ogni osservazione di  $Y$  è uguale ad una combinazione lineare dei regressori con pesi=coefficienti beta + un termine di errore

$$Y = X\beta + \varepsilon$$

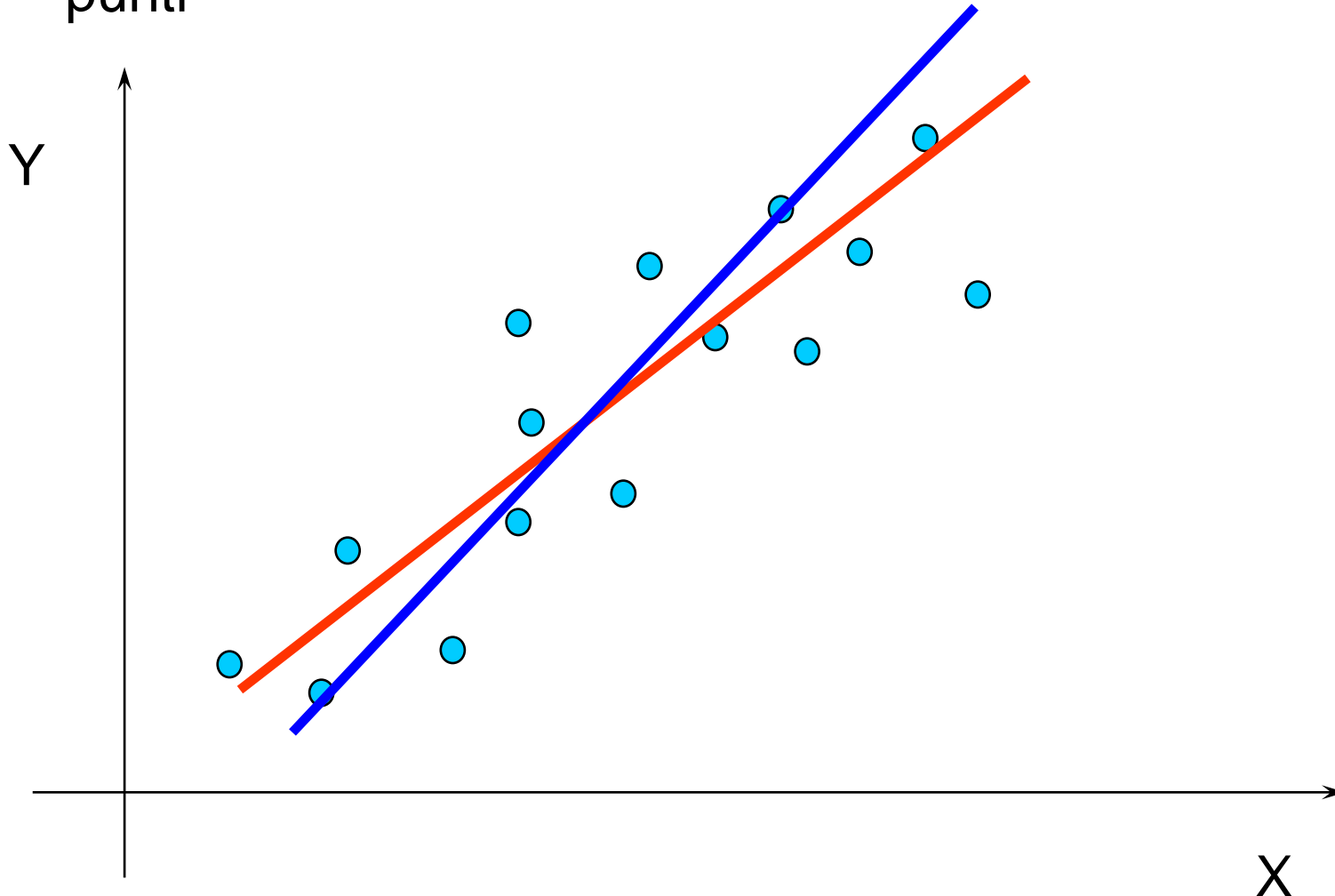
# Il modello di regressione lineare

1. Introduzione ai modelli di regressione
2. Obiettivi
3. Le ipotesi del modello
4. La stima del modello
5. La valutazione del modello
6. Commenti

# Il modello di regressione lineare

## La stima del modello

Si vuole trovare la retta lineare migliore data la nuvola di punti



# Il modello di regressione lineare

## La stima del modello

Equazione teorica → coefficienti non noti

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p + \varepsilon$$

Equazione stimata → coefficienti stimati (una delle infinite rette possibili)

$$Y = b_0 + b_1 X_1 + b_2 X_2 + \dots + b_p X_p + \hat{\varepsilon}$$

stime dei  
coefficienti

previsione

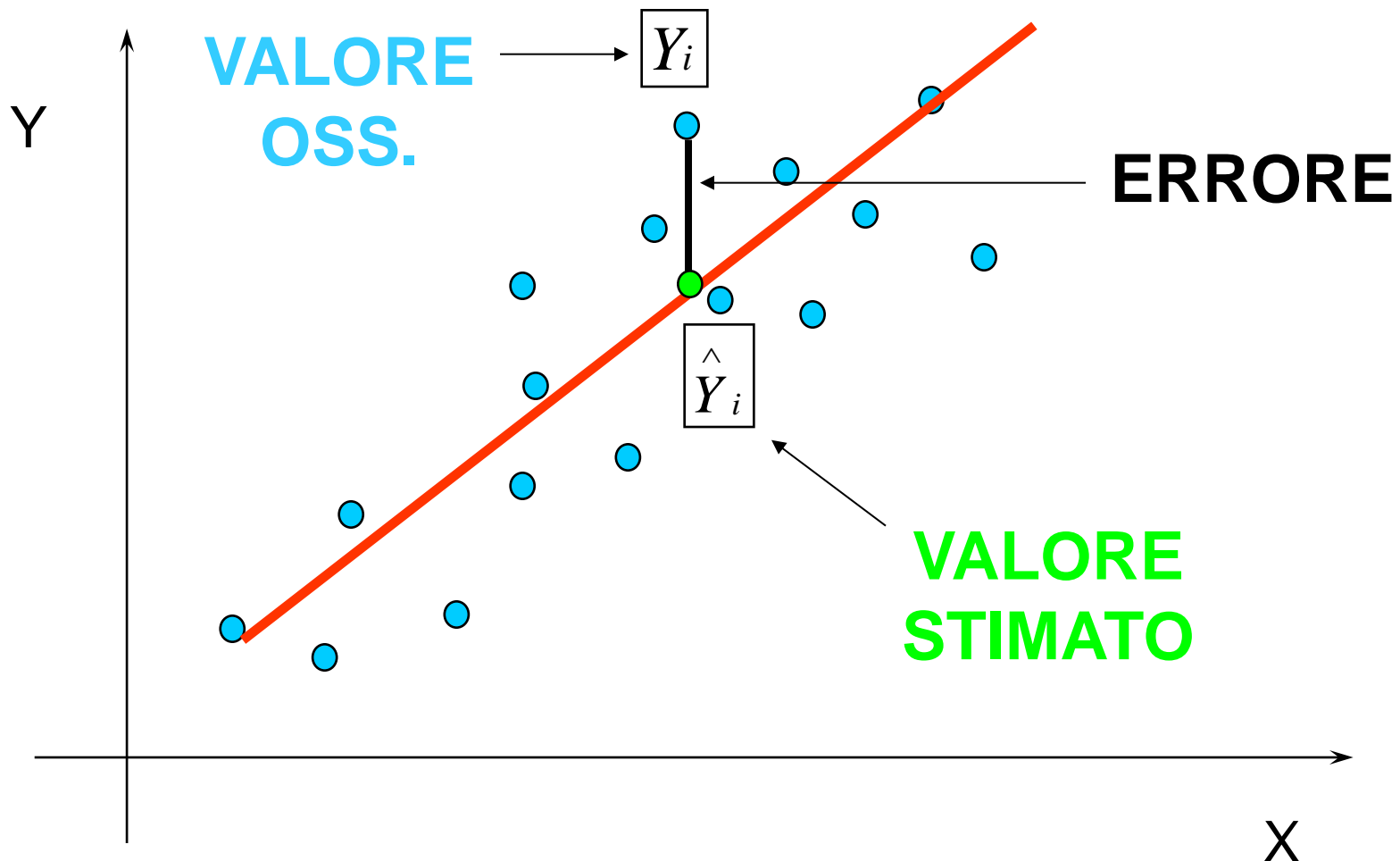
errore di  
previsione

$$Y = \hat{Y} + \hat{\varepsilon}$$

# Il modello di regressione lineare

## La stima del modello

Stimando la retta di regressione si commette un errore di previsione: Metodo dei Minimi Quadrati





# Il modello di regressione lineare

## La stima del modello

Obiettivo → trovare la miglior approssimazione lineare della relazione tra  $Y$  e  $X_1, \dots, X_p$  (trovare le stime dei parametri beta che identificano la “migliore” retta di regressione)

Metodo dei minimi quadrati → lo stimatore LS è la soluzione al problema

$$\min_{\beta} \sum_{i=1}^n (y_i - X_i \beta)^2 \Leftrightarrow \min_{\beta} \varepsilon' \varepsilon$$

# Il modello di regressione lineare

## La stima del modello

Lo stimatore dei Minimi Quadrati: LS

- è funzione di  $Y$  e  $X$

$$\hat{\beta}_{LS} = (X'X)^{-1} X'Y$$

- ha media

$$E(\hat{\beta}_{LS}) = \beta$$

- ha varianza

$$Var(\hat{\beta}_{LS}) = (X'X)^{-1} \sigma^2$$

# Il modello di regressione lineare

## La stima del modello

### Proprietà dello stimatore LS

- non distorto
- consistente (se valgono certe hp su  $X'X$ )
- coincide con lo stimatore di max verosimiglianza sotto hp forti

→ BLUE (Best Linear Unbiased Estimator)

# Il modello di regressione lineare

## La stima del modello

Equazione teorica → coefficienti non noti

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p + \varepsilon$$

Equazione stimata → coefficienti stimati (una delle infinite rette possibili)

$$\hat{Y} = b_0 + b_1 X_1 + b_2 X_2 + \dots + b_p X_p$$



stime dei  
coefficienti

# Il modello di regressione lineare

## La stima del modello

### Interpretazione dei coefficienti

- impatto di  $X_j$  su  $Y$  posto che nel modello sono presenti altre variabili
- tasso di variazione di  $Y$  al variare di  $X_j$
- come varia  $Y$  al variare di una unità di  $X_j$  se gli altri regressori non variano

# Il modello di regressione lineare

## La stima del modello

### Parameter Estimates

Variable	Label	DF	Parameter Estimate	Standard Error	t Value	Pr >  t
Intercept	Intercept	1	-15016	2324.86370	-6.46	<.0001
PAG_ORD	Pagato in contrassegno	1	1.19433	0.05485	21.78	<.0001
PAG_MES	Pagato con rate mensili	1	2.52341	0.10102	24.98	<.0001
TOT_ORD	Totale ordini	1	14881	683.88703	21.76	<.0001
LISTA	Numero di liste di appartenenza	1	603.36550	1110.84778	0.54	0.5871
SESSO	Sesso	1	3453.14705	1994.83468	1.73	0.0835
CEN	Residenza Centro	1	-6431.88493	2597.25872	-2.48	0.0133
SUD	Residenza Sud	1	-18390	2077.96317	-8.85	<.0001

# Il modello di regressione lineare

## La stima del modello

### Segno del coefficiente

- indica la direzione dell'impatto del regressore a cui è associato

### Valore del coefficiente

- indica l'incremento marginale di  $Y$
- dipende dall'unità di misura di  $X_j$

# Il modello di regressione lineare

## La stima del modello

### Valore del coefficiente

- per valutare l'impatto relativo dei singoli regressori è necessario considerare i coefficienti standardizzati

#### Parameter Estimates

Variable	Label	D F	Parameter Estimate	Standard Error	t Value	Pr >  t	Standardized Estimate	Variance Inflation
Intercept	Intercept	1	30935	869.69238	35.57	<.0001	0	0
Factor1		1	61162	869.81092	70.32	<.0001	0.71583	1.00000
Factor3		1	24154	869.81092	27.77	<.0001	0.28269	1.00000
Factor4		1	3446.48124	869.81092	3.96	<.0001	0.04034	1.00000
Factor6		1	-13861	869.81092	-15.94	<.0001	-0.16223	1.00000