

1. Popolazioni, individui e caratteri. Indicatori sintetici di campioni monovariati

La statistica descrittiva si occupa dell'analisi di dati raccolti da una popolazione, ovvero da un insieme di individui. In sintesi, dato un insieme molto grande di dati, così grande che non è *utile* guardarlo dato per dato, si cerca di estrarne delle informazioni sintetiche e tuttavia significative.

Gli oggetti con cui abbiamo a che fare sono dunque

- gli **individui** oggetto dell'indagine: ciascun individuo è un oggetto singolo dell'indagine.
- la **popolazione**, ovvero l'insieme degli individui oggetto dell'indagine.
- il **carattere** osservato o variabile, che è la quantità misurata o la qualità rilevata su ciascun individuo della popolazione.

Esempio 1.0.1. Rilevo l'altezza di ciascun abitante del Comune di Firenze. Ogni residente del Comune di Firenze è un individuo; la popolazione è l'insieme di tutti i residenti nel Comune di Firenze; il carattere in esame è l'altezza misurata, per esempio, in centimetri.

Esempio 1.0.2. Rilevo il reddito annuo di ciascun nucleo familiare del Comune di Firenze. Ogni nucleo familiare è un individuo; la popolazione è l'insieme dei nuclei familiari registrati all'Anagrafe del Comune di Firenze; il carattere osservato è il reddito annuo familiare misurato in Euro.

Esempio 1.0.3. Rilevo il numero dei componenti di ciascun nucleo familiare del Comune di Firenze. Come nell'esempio precedente ogni nucleo familiare è un individuo; la popolazione è l'insieme dei nuclei familiari registrati all'Anagrafe del Comune di Firenze. Il carattere osservato è il numero dei componenti di ciascun nucleo familiare, cioè un numero intero maggiore-uguale di 1.

Esempio 1.0.4. Per ogni studente presente in aula rilevo il colore degli occhi. Ogni studente presente in aula è un individuo. La popolazione è l'insieme degli studenti presenti ed il carattere osservato è il colore degli occhi.

In questi esempi abbiamo incontrato i due tipi fondamentali di carattere:

- **caratteri numerici o quantitativi** come l'altezza, il reddito familiare, il numero dei componenti del nucleo familiare;
- **caratteri qualitativi** come il colore degli occhi.

I caratteri numerici a loro volta si possono suddividere in

- **caratteri numerici discreti** che possono assumere solo un insieme discreto di valori, come il numero dei componenti dei nuclei familiari;
- **caratteri numerici continui** che variano con continuità ovvero con una estrema accuratezza, eccessiva rispetto ai fini dell'indagine, come l'altezza delle persone o il reddito annuo familiare.

1.1 Campione statistico, modalità e classi modali

Supponiamo di aver osservato un certo carattere su una popolazione di n individui. Abbiamo un *vettore delle osservazioni*

$$x = (x_1, \dots, x_n)$$

che chiamiamo **campione statistico** di cardinalità n .

Se il campione è relativo ad un carattere qualitativo o numerico discreto, chiamo **modalità** i valori che esso assume su un campione.

Se il campione è relativo ad un carattere numerico continuo si procede nel seguente modo: la popolazione in esame è comunque un insieme finito, quindi il carattere, per quanto continuo, nel campione assume solo un numero finito di valori. Sia $[a, b)$ un intervallo che contiene tutti i valori x_i , $i = 1, \dots, n$ assunti dal carattere sugli individui della popolazione. Suddividiamo l'intervallo $[a, b)$ in N parti uguali (N sarà suggerito dall'esperienza). Otteniamo N intervalli

$$I_j := \left[a + (j-1) \frac{b-a}{N}, a + j \frac{b-a}{N} \right), \quad j = 1, \dots, N.$$

Se I_j contiene almeno un'osservazione, dico che I_j è una **classe di modalità** del campione.

1.2 Frequenza assoluta e frequenza relativa

Consideriamo un campione $x = (x_1, \dots, x_n)$ relativo ad un carattere qualitativo o numerico discreto. Nel campione, cioè nella popolazione in esame, il carattere osservato assume un certo numero di valori distinti

$$z_1, \dots, z_k, \quad 1 \leq k \leq n.$$

Per ogni $j = 1, \dots, k$ chiamo **effettivo** o **frequenza assoluta** della modalità z_j il numero

$$n_j := \#\{i \in \{1, \dots, n\} : x_i = z_j\}$$

mentre chiamo **frequenza relativa** della modalità z_j il numero

$$p_j := \frac{n_j}{n}.$$

Se il carattere osservato è numerico continuo, si considera ciascuna classe di modalità individuata

$$I_j := \left[a + (j-1) \frac{b-a}{N}, a + j \frac{b-a}{N} \right), \quad j = 1, \dots, N$$

e si chiama **frequenza assoluta o effettivo** della classe di modalità I_j il numero

$$n_j := \#\{i \in \{1, \dots, n\} : x_i \in I_j\}.$$

Come prima definiamo **frequenza relativa** della classe I_j il numero $p_j := \frac{n_j}{n}$.

1.3 Moda e valori modalì

Sia $x = (x_1, \dots, x_n)$ un campione statistico e siano z_1, z_2, \dots, z_k le modalit  assunte (o I_1, \dots, I_k le classi di modalit  assunte) e siano p_1, \dots, p_k le relative frequenze relative.

Se esiste uno ed un solo indice $\bar{j} \in \{1, 2, \dots, k\}$ tale che la modalit  $z_{\bar{j}}$ (o la classe $I_{\bar{j}}$) ha frequenza massima, ovvero se esiste un unico $\bar{j} \in \{1, 2, \dots, k\}$ tale che $p_{\bar{j}} \geq p_j \forall j = 1, \dots, k$, allora la modalit  $z_{\bar{j}}$ (o la classe $I_{\bar{j}}$) si dice **moda** del campione x .

Se esistono due o pi  indici $\bar{j}_1, \bar{j}_2, \dots, \bar{j}_s$ tali che le modalit  $z_{\bar{j}_1}, z_{\bar{j}_2}, \dots, z_{\bar{j}_s}$ (o le classi $I_{\bar{j}_1}, I_{\bar{j}_2}, \dots, I_{\bar{j}_s}$) hanno frequenza massima, allora queste modalit  (o classi) si dicono **valori (o classi) modalì**.

Possiamo visualizzare con degli istogrammi, vedi Figura 1.3

1.4 Mediana

D’ora innanzi consideriamo solo caratteri numerici.

Sia dunque $x = (x_1, \dots, x_n)$ un campione relativo ad un carattere numerico. Ordiniamo i dati del campione in ordine crescente:

$$x_{(1)} \leq x_{(2)} \leq \dots \leq x_{(n-1)} \leq x_{(n)}$$

e distinguiamo due casi:

- n dispari: $n = 2m + 1$

$$x_{(1)} \leq x_{(2)} \leq \dots \leq x_{(m)} \leq x_{(m+1)} \leq x_{(m+2)} \leq \dots \leq x_{(2m)} \leq x_{(2m+1)}$$

Il dato $x_{(m+1)}$   maggiore-uguale di m dati e minore-uguale di altrettanti dati. Diciamo che il dato $x_{(m+1)}$   la **mediana** del campione.

- n pari: $n = 2m$

$$x_{(1)} \leq x_{(2)} \leq \dots \leq x_{(m-1)} \leq x_{(m)} \leq x_{(m+1)} \leq x_{(m+2)} \leq \dots \leq x_{(2m)}$$

Il dato $x_{(m)}$   maggiore-uguale di $m - 1$ dati e minore-uguale di m dati. Il dato $x_{(m+1)}$   maggiore-uguale di m dati e minore-uguale di $m - 1$ dati.

Chiamiamo **mediana** del campione il numero $\frac{x_{(m)} + x_{(m+1)}}{2}$.

1.5 Media e varianza campionaria. Scarto quadratico medio (o deviazione standard)

Consideriamo un campione relativo ad un carattere numerico

$$x = (x_1, \dots, x_n).$$

Chiamo **media aritmetica** o, pi  semplicemente, **media** il numero

$$\bar{x} := \frac{1}{n} \sum_{i=1}^n x_i.$$

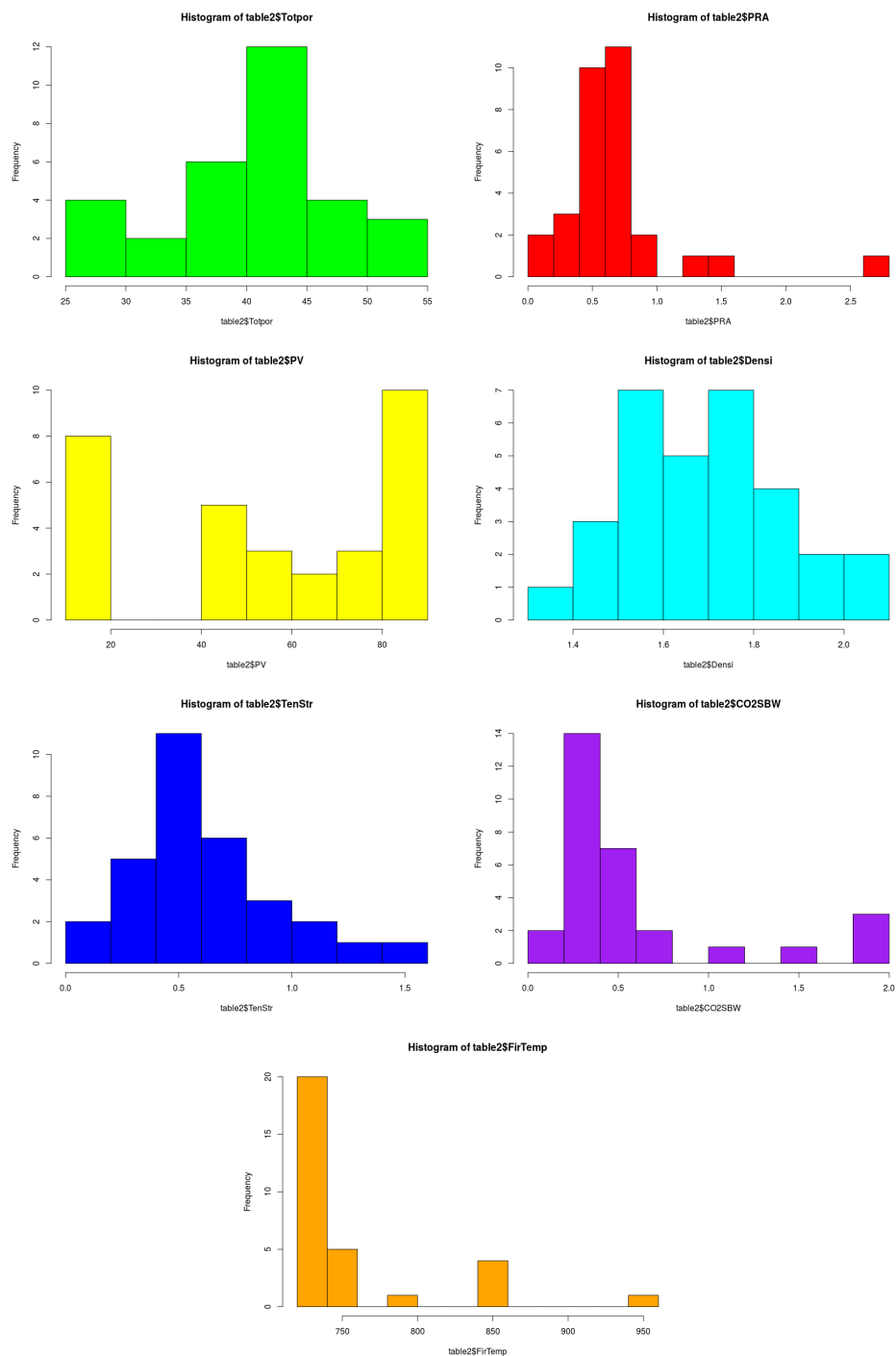


Figura 1.1: Alcuni istogrammi dall'Esempio 1.5.1

Supponiamo che nel campione siano presenti k modalità z_1, z_2, \dots, z_k con rispettive frequenze assolute N_1, N_2, \dots, N_k e frequenze relative p_1, p_2, \dots, p_k . Allora

$$\begin{aligned} \bar{x} &= \frac{1}{n} (x_1 + x_2 + \dots + x_n) = \frac{1}{n} (N_1 z_1 + N_2 z_2 + \dots + N_k z_k) = \\ &= p_1 z_1 + p_2 z_2 + \dots + p_k z_k = \sum_{j=1}^k p_j z_j. \end{aligned}$$

Chiamo **varianza campionaria** di x il numero non-negativo

$$s_x^2 := \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2.$$

Osserviamo che la media è un valore centrale attorno al quale si dispongono i dati x_1, \dots, x_n mentre la varianza è un *indice di dispersione*: la varianza è nulla se e solo se tutti i dati del campione sono uguali (e dunque coincidono con la media). Una varianza bassa indica che comunque i dati sono *vicini* al valore medio \bar{x} mentre una varianza alta indica una maggiore dispersione dei dati.

La radice quadrata della varianza campionaria

$$s_x = \text{Std}[x] := \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2}$$

si chiama **scarto quadratico medio** o **deviazione standard** del campione x .

Anche per la varianza campionaria possiamo scrivere una formula che coinvolga solo le modalità e le rispettive frequenze.

$$\begin{aligned} s_x^2 &= \frac{1}{n-1} ((x_1 - \bar{x})^2 + (x_2 - \bar{x})^2 + \dots + (x_n - \bar{x})^2) = \\ &= \frac{1}{n-1} (N_1(z_1 - \bar{x})^2 + N_2(z_2 - \bar{x})^2 + \dots + N_k(z_k - \bar{x})^2) = \\ &= \frac{n}{n-1} (p_1(z_1 - \bar{x})^2 + p_2(z_2 - \bar{x})^2 + \dots + p_k(z_k - \bar{x})^2) = \\ &= \frac{n}{n-1} \sum_{j=1}^k p_j(z_j - \bar{x})^2. \end{aligned}$$

Esempio 1.5.1. Nella tabella che segue, tratta da [2], riportiamo alcuni dati relativi a campioni di laterizio e che useremo per fare alcuni esempi relativi alle nozioni introdotte mediante il software R <http://cran.r-project.org/>. Per una introduzione si rimanda ai manuali [3] e [1].

SAMPLE CODE	POROSITÀ TOTALE (%)	RAGGIO MEDIO DEL PORO (μm)	VOLUME DEI PORI SU DIMEN- SIONE DEI PORI 0.3–0.8 μm	DENSITÀ (g/cm^3)	RESISTENZA ALLA TRA- ZIONE (MPa)	CO ₂ /SBW	TEMPERATURA DI COTTURA (DTA)
AS1	41.460	0.528	80.0	1.550	0.403	0.38	740
AS2	47.210	0.467	81.2	1.650	0.645	0.70	740
AS3	43.670	0.697	78.5	1.710	0.527	0.46	740
AS4	52.390	0.422	77.3	1.520	0.143	0.48	740
AS5	44.700	0.411	87.4	1.500	0.593	0.29	740
AS6	51.330	0.422	88.6	1.480	0.463	0.33	740
AS7	31.460	0.718	80.6	1.900	0.955	0.23	740
AS8	40.900	0.458	80.4	1.680	0.195	0.41	740
AS9	45.540	0.492	80.8	1.620	1.328	0.50	750
AS10	45.620	0.734	86.2	1.620	1.405	0.34	750
AS11	44.140	0.730	85.7	1.590	0.256	0.42	750
AS12	40.710	0.543	87.8	1.750	0.309	0.20	750
AS13	35.700	0.686	84.3	1.520	0.472	0.05	740
C1	40.290	0.306	43.5	1.760	0.520	0.43	740
C2	36.570	0.625	42.3	1.750	0.738	0.36	740
C3	42.130	0.249	63.2	1.630	0.410	0.25	740
C4	37.830	0.731	47.9	2.020	0.601	0.28	740
C5	42.180	0.407	59.4	1.580	0.376	0.34	740
C6	41.600	0.446	42.8	1.850	0.473	0.26	740
C7	32.660	0.664	64.3	1.850	0.695	0.25	740
C8	36.070	0.673	58.2	1.780	0.624	0.29	740
C9	36.040	1.397	55.6	1.730	0.582	0.38	740
C10	36.640	0.861	45.2	1.750	0.650	0.47	740
R1	42.890	0.785	10.2	1.540	0.453	1.04	850
R2	26.850	0.315	14.7	2.010	1.124	1.86	960
R3	28.550	0.158	18.6	1.920	0.937	1.96	850
R4	29.860	0.158	15.3	1.890	1.020	1.48	850
R5	45.700	0.984	12.8	1.500	0.328	–	800
R6	54.640	1.525	12.5	1.340	0.267	0.67	750
R7	27.550	2.657	14.6	1.920	0.892	0.40	730
R8	40.820	0.622	15.3	1.570	0.502	1.94	860

Inseriamo la tabella in R

```
> library(readr)
> table2 <- read_delim("~/Documents/didattica/2017-18_analisi_reale/alcuni_appunti/table2.
+   "\t", escape_double = FALSE, trim_ws = TRUE)
Parsed with column specification:
cols(
  Code = col_character(),
  Totpor = col_double(),
  PRA = col_double(),
  PV = col_double(),
  Densi = col_double(),
  TenStr = col_double(),
  CO2SBW = col_character(),
  FirTemp = col_integer()
```

```
)
> View(table2)
```

	Code	Totpor	PRA	PV	Densi	TenStr	CO2SBW	FirTemp
1	AS1	41.46	0.528	80.0	1.55	0.403	0.38	740
2	AS2	47.21	0.467	81.2	1.65	0.645	0.70	740
3	AS3	43.67	0.697	78.5	1.71	0.527	0.46	740
4	AS4	52.39	0.422	77.3	1.52	0.143	0.48	740
5	AS5	44.70	0.411	87.4	1.50	0.593	0.29	740
6	AS6	51.33	0.422	88.6	1.48	0.463	0.33	740
7	AS7	31.46	0.718	80.6	1.90	0.955	0.23	740
8	AS8	40.90	0.458	80.4	1.68	0.195	0.41	740
9	AS9	45.54	0.492	80.8	1.62	1.328	0.50	750
10	AS10	45.62	0.734	86.2	1.62	1.405	0.34	750
11	AS11	44.14	0.730	85.7	1.59	0.256	0.42	750
12	AS12	40.71	0.543	87.8	1.75	0.309	0.20	750
13	AS13	35.70	0.686	84.3	1.52	0.472	0.05	740
14	C1	40.29	0.306	43.5	1.76	0.520	0.43	740
15	C2	36.57	0.625	42.3	1.75	0.738	0.36	740
16	C3	42.13	0.249	63.2	1.63	0.410	0.25	740
17	C4	37.83	0.731	47.9	2.02	0.601	0.28	740
18	C5	42.18	0.407	59.4	1.58	0.376	0.34	740
19	C6	41.60	0.446	42.8	1.85	0.473	0.26	740
20	C7	32.66	0.664	64.3	1.85	0.695	0.25	740
21	C8	36.07	0.673	58.2	1.78	0.624	0.29	740
22	C9	36.04	1.397	55.6	1.73	0.582	0.38	740
23	C10	36.64	0.861	45.2	1.75	0.650	0.47	740
24	R1	42.89	0.785	10.2	1.54	0.453	1.04	850
25	R2	26.85	0.315	14.7	2.01	1.124	1.86	960
26	R3	28.55	0.158	18.6	1.92	0.937	1.96	850
27	R4	29.86	0.158	15.3	1.89	1.020	1.48	850
28	R5	45.70	0.984	12.8	1.50	0.328	--	800
29	R6	54.64	1.525	12.5	1.34	0.267	0.67	750
30	R7	27.55	2.657	14.6	1.92	0.892	0.40	730
31	R8	40.82	0.622	15.3	1.57	0.502	1.94	860

Per ciascun carattere definiamo una variabile che contenga la mediana, una per la media, una per la Varianza e una per la deviazione standard e poi stampiamo i valori (tratteremo il carattere di nome CO2SBW con attenzione perché su un individuo non è stato rilevato)

Il comando `summary` indica il numero di dati mancanti, ci dà gli indicatori di centralità ma non quelli di dispersione

```
> summary(table2)
      Code      Totpor      PRA      PV      Densi      TenStr      CO2SBW      FirTemp
Length:31  Min. :26.85  Min. :0.1580  Min. :10.20  Min. :1.340  Min. :0.1430  Min. :0.0500  Min. :730.0
Class :character 1st Qu.:36.05  1st Qu.:0.4220  1st Qu.:30.45  1st Qu.:1.560  1st Qu.:0.4065  1st Qu.:0.2900  1st Qu.:740.0
Mode  :character Median :40.90  Median :0.6220  Median :59.40  Median :1.680  Median :0.5270  Median :0.3900  Median :740.0
      Mean :40.12  Mean :0.6733  Mean :55.33  Mean :1.693  Mean :0.6092  Mean :0.5817  Mean :764.8
      3rd Qu.:44.42  3rd Qu.:0.7305  3rd Qu.:80.70  3rd Qu.:1.815  3rd Qu.:0.7165  3rd Qu.:0.4950  3rd Qu.:750.0
      Max. :54.64  Max. :2.6570  Max. :88.60  Max. :2.020  Max. :1.4050  Max. :1.9600  Max. :960.0
      NA's :1
```

Richiediamo anche varianza campionaria e deviazione standard.

```
> medianaTotPor <- median(table2$Totpor);
> meanTotPor <- mean(table2$Totpor);
> VarTotPor <- var(table2$Totpor);
> StdTotPor <- sd(table2$Totpor)
> medianaTotPor; meanTotPor; VarTotPor; StdTotPor
[1] 40.9
[1] 40.11935
[1] 49.52185
[1] 7.037176
> medianaPRA <- median(table2$PRA);
> meanPRA <- mean(table2$PRA);
VarPRA <- var(table2$PRA);
> StdPRA <- sd(table2$PRA)
> medianaPRA; meanPRA; VarPRA; StdPRA
[1] 0.622
[1] 0.6732581
[1] 0.226613
[1] 0.4760389
> medianaPV <- median(table2$PV);
> meanPV <- mean(table2$PV);
> VarPV <- var(table2$PV);
> StdPV <- sd(table2$PV)
> medianaPV; meanPV; VarPV; StdPV
[1] 59.4
[1] 55.32903
[1] 815.0935
[1] 28.54984
> medianaDensi <- median(table2$Densi);
> meanDensi <- mean(table2$Densi);
> VarDensi <- var(table2$Densi);
> StdDensi <- sd(table2$Densi)
> medianaDensi; meanDensi; VarDensi; StdDensi
[1] 1.68
[1] 1.692903
[1] 0.02894129
[1] 0.1701214
> medianaTenStr <- median(table2$TenStr);
> meanTenStr <- mean(table2$TenStr);
> VarTenStr <- var(table2$TenStr);
> StdTenStr <- sd(table2$TenStr)
> medianaTenStr; meanTenStr; VarTenStr; StdTenStr
[1] 0.527
[1] 0.6092258
[1] 0.09882738
[1] 0.3143682
```



```
> medianaCO2SBW <- median(na.omit(table2$CO2SBW));
> meanCO2SBW <- mean(na.omit(table2$CO2SBW));
> VarCO2SBW <- var(na.omit(table2$CO2SBW));
> StdCO2SBW <- sd(na.omit(table2$CO2SBW))
> medianaCO2SBW; meanCO2SBW; VarCO2SBW; StdCO2SBW
[1] 0.39
[1] 0.5816667
[1] 0.2765868
[1] 0.5259152
> medianaFirTemp <- median(table2$FirTemp);
> meanFirTemp <- mean(table2$FirTemp);
> VarFirTemp <- var(table2$FirTemp);
> StdFirTemp <- sd(table2$FirTemp)
> medianaFirTemp; meanFirTemp; VarFirTemp; StdFirTemp
[1] 740
[1] 764.8387
[1] 2805.806
[1] 52.96986
```


2. Campioni bivariati: covarianza, coefficiente di correlazione e retta di regressione

2.1 Covarianza e coefficiente di correlazione

Supponiamo di avere un **campione bivariato** cioè di rilevare due caratteri sugli individui di una medesima popolazione.

Abbiamo dunque due vettori di dati

$$x = (x_1, \dots, x_n), \quad y = (y_1, \dots, y_n).$$

x_i e y_i sono le rilevazioni dei due caratteri sul medesimo individuo, l'individuo cioè che abbiamo etichettato come *individuo i*.

Chiamiamo **covarianza di x e y** il numero

$$\text{Cov}(x, y) := \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

dove \bar{x} e \bar{y} sono le medie dei campioni x e y , rispettivamente.

Nel caso in cui né x né y siano campioni costanti (ipotesi che sarà sempre sottintesa), definiamo **coefficiente di correlazione di x e y** il numero

$$\rho[x, y] := \frac{\text{Cov}(x, y)}{s_x s_y} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\left(\sum_{i=1}^n (x_i - \bar{x})^2\right) \left(\sum_{i=1}^n (y_i - \bar{y})^2\right)}}.$$

Osservazione 2.1.1. $\text{Cov}(x, x) = s_x^2$; $\rho[x, x] = 1$.

Osservando che $\rho[x, y]$ non è altro che il rapporto tra $\langle x - (\bar{x}, \dots, \bar{x}), y - (\bar{y}, \dots, \bar{y}) \rangle$ (prodotto scalare) e $\|x - (\bar{x}, \dots, \bar{x})\| \|y - (\bar{y}, \dots, \bar{y})\|$ (prodotto delle norme) si ottengono le seguenti proprietà:

1. $-1 \leq \rho[x, y] \leq 1$;
2. $\rho[x, y] = 1$ se e solo se esiste $a > 0, b \in \mathbb{R}$ tale che $y_i = ax_i + b \quad \forall i = 1, \dots, n$. In tal caso i campioni x e y si dicono *positivamente correlati*;
3. $\rho[x, y] = -1$ se e solo se esiste $a < 0, b \in \mathbb{R}$ tale che $y_i = ax_i + b \quad \forall i = 1, \dots, n$. In tal caso i campioni x e y si dicono *negativamente correlati*.

Se $\rho[x, y] = 0$ i campioni x e y si dicono *scorrelati*.

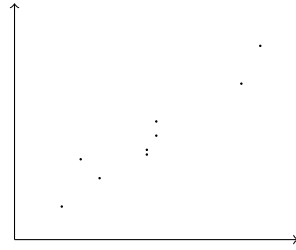


Figura 2.1: Campione bivariato *pressoché lineare*

2.2 Retta di regressione

Supponiamo di avere un campione bivariato

$$x = (x_1, \dots, x_n), \quad y = (y_1, \dots, y_n)$$

dove x_i e y_i sono i dati relativi all' i -esimo individuo. Rappresentiamo i punti (x_i, y_i) sul piano cartesiano Oxy . Capita, molto spesso, di trovarsi a disposizioni *pressoché allineate* come illustrato nella figura 2.1 Si cerca allora una retta che in qualche senso *approssimi* i punti (x_i, y_i) .

Supponiamo che $y = ax + b$ sia l'equazione della retta cercata. Per $x = x_i$ si ottiene il punto sulla retta $(x_i, ax_i + b)$. Cerchiamo la retta (ovvero i parametri a e b) che minimizza la *somma degli errori quadratici nella direzione y*

$$S(a, b) := \sum_{i=1}^n (y_i - (ax_i + b))^2.$$

Si ha

$$\begin{aligned} S(a, b) &= \sum_{i=1}^n (y_i - \bar{y} + \bar{y} - (ax_i - a\bar{x} + a\bar{x} + b))^2 = \\ &= \sum_{i=1}^n ((y_i - \bar{y}) - a(x_i - \bar{x}) + (\bar{y} - a\bar{x} - b))^2 = \\ &= \sum_{i=1}^n ((y_i - \bar{y})^2 + a^2 \sum_{i=1}^n (x_i - \bar{x})^2 + \\ &\quad + n(\bar{y} - a\bar{x} - b)^2 - 2a \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})) = \\ &= (n-1)(s_y^2 + a^2 s_x^2 - 2a \text{Cov}(x, y)) + n(\bar{y} - a\bar{x} - b)^2 \\ &= (n-1) \left(\left(a s_x - \frac{\text{Cov}(x, y)}{s_x} \right)^2 + s_y^2 - \left(\frac{\text{Cov}(x, y)}{s_x} \right)^2 \right) + n(\bar{y} - a\bar{x} - b)^2. \end{aligned}$$

Il minimo dello somma degli errori quadratici $S(a, b)$ si ottiene allora per

$$a = \frac{\text{Cov}(x, y)}{s_x^2}; \quad b = \bar{y} - \frac{\text{Cov}(x, y)}{s_x^2} \bar{x};$$

il minimo dell'errore S vale

$$(n - 1) \left(s_y^2 - \frac{(\text{Cov}(x, y))^2}{s_x^2} \right) = (n - 1) s_y^2 (1 - (\rho[x, y])^2)$$

e la retta ha equazione

$$y = \bar{y} + \frac{\text{Cov}(x, y)}{s_x^2} (x - \bar{x}).$$

Osservazione 2.2.1. La retta così determinata si chiama **retta di regressione del campione y sul campione x** . Osserviamo infine che il punto (\bar{x}, \bar{y}) appartiene alla retta.

Esempio 2.2.1. Riconsideriamo l'esempio 1.5.1. Carichiamo in R la tabella dei dati.

```
> library(readr)
> table2 <- read_delim("~/Documents/didattica/2017-18_analisi_reale/alcuni_appunti/table2.csv",
+   "\t", escape_double = FALSE, trim_ws = TRUE)
Parsed with column specification:
cols(
  Code = col_character(),
  Totpor = col_double(),
  PRA = col_double(),
  PV = col_double(),
  Densi = col_double(),
  TenStr = col_double(),
  CO2SBW = col_character(),
  FirTemp = col_integer()
)
```

Tracciamo sul piano cartesiano i dati relativi ai caratteri porosità totale (in ascissa) e densità (in ordinata) e salviamo la figura in un file.

```
> library(car)
> scatterplot(Densi~Totpor, lm=TRUE, smooth=FALSE, spread=FALSE, boxplots=TRUE, span=0.5, data= table2)
```

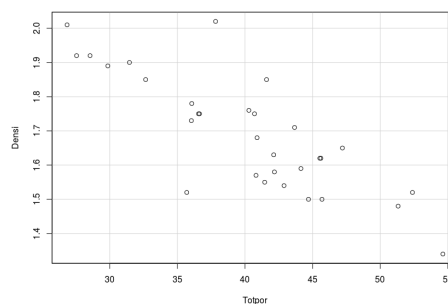


Figura 2.2: Porosità totale versus Densità

Sembrano *ragionevolmente allineati*. Calcoliamo il loro coefficiente di correlazione

```
> CorTotporDensi<- cor(table2$Totpor, table2$Densi)
> CorTotporDensi
[1] -0.8187597
```

Calcoliamo la retta di regressione del carattere Densità sul carattere Porosità Totale

```
> RegModel.Densi.Totpor <- lm(Densi~Totpor, data=table2)
> summary(RegModel.Densi.Totpor)
```

Call:

```
lm(formula = Densi ~ Totpor, data = table2)
```

Residuals:

Min	1Q	Median	3Q	Max
-0.260377	-0.054570	-0.001898	0.045213	0.281783

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	2.486995	0.104930	23.70	< 2e-16 ***
Totpor	-0.019793	0.002577	-7.68	1.81e-08 ***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.09934 on 29 degrees of freedom

Multiple R-squared: 0.6704, Adjusted R-squared: 0.659

F-statistic: 58.98 on 1 and 29 DF, p-value: 1.814e-08

Intercept dice che l'ordinata all'origine (il coefficiente b) della retta di regressione è 2.486995 mentre il coefficiente angolare (cioè a) è -0.019793 . Ridisegniamo i punti sul piano cartesiano, aggiungendo la retta di regressione (e salviamo l'immagine in un file).

```
> abline(lm(Densi ~ Totpor, data=table2), col="red")
```

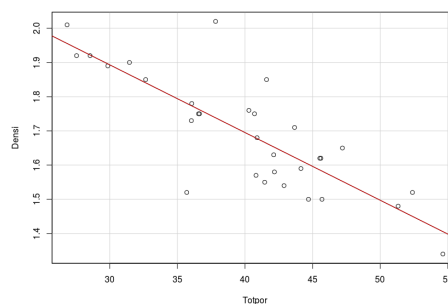


Figura 2.3: Retta di regressione lineare

3. Campioni multivariati. Principal Components Analysis

Lo scopo di questa analisi è il seguente: supponiamo di avere un campione multivariato. Supponiamo cioè di aver raccolto dati relativi a più caratteri, diciamo k caratteri, su una popolazione di n individui.

Riportiamo le informazioni raccolte come nella tabella dell'esempio 1.5.1. Ovvero

- Nella prima riga riportiamo i dati relativi al primo individuo, carattere per carattere

$$x_{11} \quad x_{12} \quad \dots \quad x_{1k}$$

- Nella seconda riga riportiamo i dati relativi al secondo individuo, carattere per carattere

$$x_{21} \quad x_{22} \quad \dots \quad x_{2k}$$

Procedendo di individuo in individuo otteniamo una matrice di n righe e k colonne:

$$X = \begin{pmatrix} x_{11} & x_{12} & \dots & x_{1k} \\ x_{21} & x_{22} & \dots & x_{2k} \\ \vdots & \vdots & & \vdots \\ x_{n1} & x_{n2} & \dots & x_{nk} \end{pmatrix} = (x_{ij})_{\substack{i=1,\dots,n \\ j=1,\dots,k}} \in \mathbb{R}^{n \times k}$$

in cui il numero in posizione (i, j) (i -esima riga e j -esima colonna) è il dato rilevato sull' i -esimo individuo relativamente al j -esimo carattere. Possiamo "leggere" la matrice colonna per colonna e rilevare le informazioni relative ad un singolo carattere. Infatti la prima colonna

$$X_1 := \begin{pmatrix} x_{11} \\ x_{21} \\ \vdots \\ x_{n1} \end{pmatrix}$$

contiene tutti i dati relativi al primo carattere, la seconda colonna

$$X_2 := \begin{pmatrix} x_{12} \\ x_{22} \\ \vdots \\ x_{n2} \end{pmatrix}$$

contiene tutti i dati relativi al secondo carattere e così via.

Per ogni $j = 1, \dots, k$ indichiamo, rispettivamente, con \bar{x}_j e s_{X_j} la media e la deviazione standard del j -esimo carattere. Si ha

$$\bar{x}_j = \frac{1}{n} \sum_{i=1}^n x_{ij}, \quad s_{X_j}^2 = \frac{1}{n-1} \sum_{i=1}^n (x_{ij} - \bar{x}_j)^2.$$

Possiamo anche calcolare la covarianza e il coefficiente di correlazione di due diversi caratteri. Più precisamente la covarianza del carattere j -esimo e del carattere ℓ -esimo è data da

$$\text{Cov}(X_\ell, X_j) = \frac{1}{n-1} \sum_{i=1}^n (x_{i\ell} - \bar{x}_\ell)(x_{ij} - \bar{x}_j).$$

Riportiamo varianze e covarianze in una matrice $k \times k$, detta **matrice di covarianza** del campione X :

$$C = (c_{\ell j})_{\substack{\ell=1, \dots, k \\ j=1, \dots, k}} \in \mathbb{R}^{k \times k} \quad c_{\ell j} := \text{Cov}(X_\ell, X_j), \quad \ell, j = 1, \dots, k.$$

Poiché $\text{Cov}(X_\ell, X_j) = \text{Cov}(X_j, X_\ell)$ la matrice C è simmetrica. Inoltre gli elementi sulla diagonale principale sono le varianze dei caratteri in esame:

$$c_{jj} = \text{Cov}(X_j, X_j) = s_{X_j}^2 \quad \forall j = 1, \dots, k.$$

Supponiamo che i coefficienti di correlazione non siano prossimi a zero, indicando dunque che i caratteri in esame sono legati gli uni agli altri.

Cerchiamo di ridurre il numero di caratteri da osservare sostituendo i caratteri originari con delle loro combinazioni lineari, in modo che i nuovi caratteri siano a due a due scorrelati e la *variabilità* del campione sia concentrata in pochi caratteri. La procedura si compone di due passi. Il primo passo consiste nel rendere le variabili adimensionali (in modo che abbia senso sommarle) e *centrate* (cioè a media nulla).

Primo passo: Standardizzazione del campione

Per ogni $i = 1, \dots, n$ e ogni $j = 1, \dots, k$ pongo

$$y_{ij} := \frac{x_{ij} - \bar{x}_j}{s_{X_j}}.$$

Ovvero il dato relativo a ciascun carattere X_j è stato sostituito da

$$Y_1 = \begin{pmatrix} y_{11} \\ y_{21} \\ \vdots \\ y_{n1} \end{pmatrix}, \quad Y_2 = \begin{pmatrix} y_{12} \\ y_{22} \\ \vdots \\ y_{n2} \end{pmatrix}, \quad \dots, \quad Y_k = \begin{pmatrix} y_{1k} \\ y_{2k} \\ \vdots \\ y_{nk} \end{pmatrix}$$

In che senso i dati Y_j sono standardizzati? Calcoliamone media e varianza campionaria

$$\begin{aligned} \bar{y}_j &= \frac{1}{n} \sum_{i=1}^n y_{ij} = \frac{1}{n} \sum_{i=1}^n \frac{x_{ij} - \bar{x}_j}{s_{X_j}} = \frac{1}{n s_{X_j}} \left(\sum_{i=1}^n x_{ij} - \sum_{i=1}^n \bar{x}_j \right) = \frac{1}{s_{X_j}} (\bar{x}_j - \bar{x}_j) = 0 \\ s_{Y_j}^2 &= \frac{1}{n-1} \sum_{i=1}^n (y_{ij} - \bar{y}_j)^2 = \frac{1}{n-1} \sum_{i=1}^n y_{ij}^2 = \frac{1}{n-1} \sum_{i=1}^n \frac{(x_{ij} - \bar{x}_j)^2}{s_{X_j}^2} = \\ &= \frac{1}{s_{X_j}^2} \frac{1}{n-1} \sum_{i=1}^n (x_{ij} - \bar{x}_j)^2 = \frac{1}{s_{X_j}^2} s_{X_j}^2 = 1. \end{aligned}$$

Calcoliamo la matrice di covarianza $C = (c_{\ell j})$ del campione standardizzato Y . Si ha

$$c_{\ell j} = \text{Cov}(Y_\ell, Y_j) = \frac{1}{n-1} \sum_{i=1}^n (y_{i\ell} - \bar{y}_\ell) (y_{ij} - \bar{y}_j) = \frac{1}{n-1} \sum_{i=1}^n y_{i\ell} y_{ij}$$

ovvero, in termini di matrici

$$C = \frac{1}{n-1} Y^t Y. \tag{3.1}$$

Osservazione 3.0.1. La formula (3.1) è vera tutte le volte che i campioni in esame hanno media nulla.

Se vogliamo calcolare i coefficienti di C in termini del campione X otteniamo anche

$$\begin{aligned} c_{\ell j} &= \text{Cov}(Y_\ell, Y_j) = \frac{1}{n-1} \sum_{i=1}^n (y_{i\ell} - \bar{y}_\ell) (y_{ij} - \bar{y}_j) = \\ &= \frac{1}{n-1} \sum_{i=1}^n y_{i\ell} y_{ij} = \frac{1}{n-1} \sum_{i=1}^n \frac{x_{i\ell} - \bar{x}_\ell}{s_{X_\ell}} \frac{x_{ij} - \bar{x}_j}{s_{X_j}} = \\ &= \frac{1}{s_{X_\ell} s_{X_j}} \frac{1}{n-1} \sum_{i=1}^n (x_{i\ell} - \bar{x}_\ell) (x_{ij} - \bar{x}_j) = \rho[X_\ell, X_j]. \end{aligned}$$

La matrice di covarianza del campione standardizzato Y è dunque simmetrica e gli elementi diagonali $c_{\ell\ell}$ sono tutti uguali ad 1.

Secondo passo: Scorrelazione dei caratteri

Considero una rotazione di \mathbb{R}^k . $A \in \mathbb{R}^{k \times k}$ matrice ortogonale: $A^t = A^{-1}$.

$$Z := YA \in \mathbb{R}^{n \times k} \quad z_{ij} = \sum_{\ell=1}^k y_{i\ell} a_{\ell j}$$

Osserviamo che nel sostituire i caratteri Y_1, \dots, Y_k con i caratteri Z_1, \dots, Z_k , per ciascun individuo $i = 1, \dots, n$ sostituiamo le osservazioni normalizzate $y_{i1}, y_{i2}, \dots, y_{ik}$ con delle loro combinazioni lineari $z_{i1}, z_{i2}, \dots, z_{ik}$ ma che gli individui *non vengono mescolati*. La riga i -esima di Z porta solo informazioni relative all'individuo i . Calcoliamo media e varianza campionaria di ciascun Z_j :

$$\bar{z}_j = \frac{1}{n} \sum_{i=1}^n z_{ij} = \frac{1}{n} \sum_{i=1}^n \sum_{\ell=1}^k y_{i\ell} a_{\ell j} = \frac{1}{n} \sum_{\ell=1}^k \left(\sum_{i=1}^n y_{i\ell} \right) a_{\ell j} = 0$$

Considero dunque la matrice delle covarianze:

$$C_Z = \frac{1}{n-1} Z^t Z = \frac{1}{n-1} (YA)^t YA = A^t \frac{1}{n-1} Y^t YA = A^t C_Y A$$

Ricorriamo ad un famoso risultato di algebra lineare:

Teorema 3.0.1 (Teorema spettrale). *Data $C \in \mathbb{R}^{k \times k}$ matrice simmetrica esiste $A \in \mathbb{R}^{k \times k}$ matrice ortogonale tale che $A^t C A$ è una matrice diagonale*

$$A^t C A = \begin{pmatrix} \lambda_1 & 0 & 0 & \dots & 0 \\ 0 & \lambda_2 & 0 & \dots & 0 \\ \vdots & \ddots & \ddots & \ddots & \vdots \\ \vdots & & \ddots & \ddots & 0 \\ 0 & & \dots & 0 & \lambda_k \end{pmatrix}.$$

Le colonne A_1, A_2, \dots, A_k della matrice A sono gli autovettori di C e gli elementi diagonali $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_k$ sono i rispettivi autovalori cioè $C A_j = \lambda_j A_j \quad \forall j = 1, \dots, k$. Inoltre λ_1 è il massimo della funzione $f(X) := \frac{X^t C X}{X^t X}$ e A_1 è un punto di massimo.

Applico il teorema spettrale: posso scegliere A in modo che la matrice $C_Z = A^t C_Y A$ sia diagonale.

- per ogni $1 \leq \ell < j < k$ il carattere ℓ -esimo Z_ℓ e il carattere j -esimo Z_j sono scorrelati: $\text{cov}(Z_\ell, Z_j) = 0$
- $\lambda_j = s_{Z_j}^2$ è la varianza del carattere j -esimo Z_j e λ_1 è il massimo di tutte le varianze. Questo ci dice che il carattere $Z_1 = Y A_1$ è il carattere che *meglio distingue* un individuo da un altro

$$z_{i1} = \sum_{\ell=1}^k y_{i\ell} a_{\ell 1} = \sum_{\ell=1}^k \frac{x_{i\ell} - \bar{x}_\ell}{s_{X_\ell}} a_{\ell 1}$$

- $\sum_{j=1}^k \lambda_j = \text{traccia}(C_Z) = \text{traccia} C_Y = k$
Posso scegliere di quante colonne Z_1, \dots, Z_j di Z tener conto in base a *quanta variabilità* voglio considerare

Esempio 3.0.1. Ritorniamo all'esempio tratto da [2]. Carichiamo la tabella a cui abbiamo tolto l'individuo R5. e visualizziamo in una *matrice di grafici*

```
> library(readr)
> X <- read_delim("~/Documents/didattica/2017-18_analisi_reale/alcuni_appunti/table2_noR5.csv",
+               "\t", escape_double = FALSE, trim_ws = TRUE)
Parsed with column specification:
cols(
  Totpor = col_double(),
  PRA = col_double(),
  PV = col_double(),
  Densi = col_double(),
  TenStr = col_double(),
  CO2SBW = col_double(),
  FirTemp = col_integer()
)
> View(X)
> plot(X)
```

Calcoliamo la matrice dei coefficienti di correlazione (che abbiamo visto essere la matrice di covarianza del campione standardizzato), con i coefficienti approssimati alla tre cifre decimali.

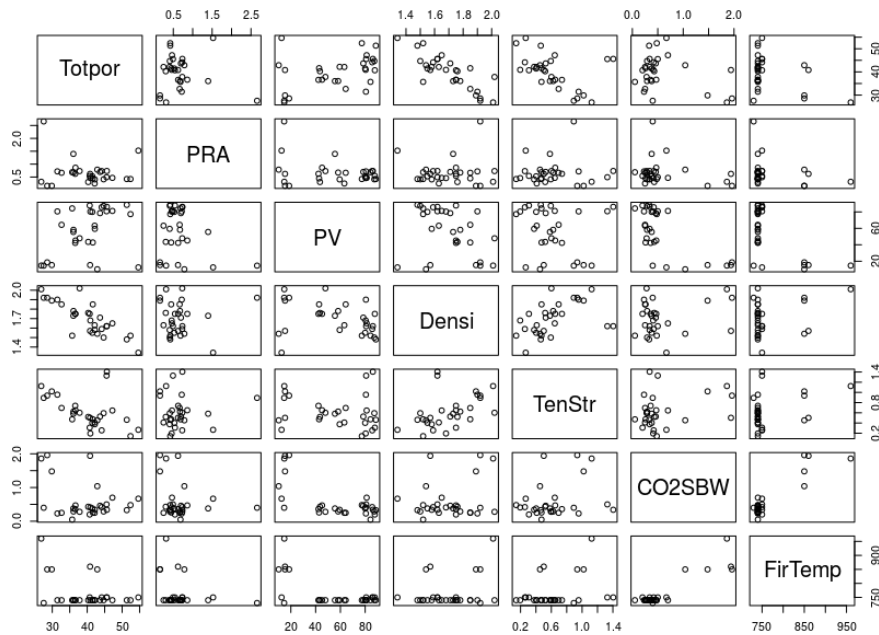


Figura 3.1: Plot dei caratteri, due a due

```
> MatrixCorr <- cor(X)
> round(MatrixCorr , 3)
      Totpor   PRA    PV  Densi TenStr CO2SBW FirTemp
Totpor  1.000 -0.116  0.411 -0.815 -0.461 -0.318 -0.398
PRA    -0.116  1.000 -0.268  0.017  0.024 -0.211 -0.258
PV      0.411 -0.268  1.000 -0.324 -0.162 -0.671 -0.624
Densi  -0.815  0.017 -0.324  1.000  0.467  0.217  0.277
TenStr -0.461  0.024 -0.162  0.467  1.000  0.289  0.328
CO2SBW -0.318 -0.211 -0.671  0.217  0.289  1.000  0.906
FirTemp -0.398 -0.258 -0.624  0.277  0.328  0.906  1.000
```

Visualizziamo i dati normalizzati (arrotondati a tre cifre decimali) e li salviamo in un file

```
> Y <- scale(X, center=TRUE, scale=TRUE)
> round(Y, 3)
      Totpor   PRA    PV  Densi TenStr CO2SBW FirTemp
[1,]  0.216 -0.281  0.833 -0.883 -0.684 -0.383 -0.443
[2,]  1.028 -0.408  0.876 -0.292  0.084  0.225 -0.443
[3,]  0.528  0.071  0.780  0.063 -0.291 -0.231 -0.443
[4,]  1.760 -0.501  0.737 -1.060 -1.508 -0.193 -0.443
[5,]  0.673 -0.524  1.098 -1.178 -0.081 -0.555 -0.443
[6,]  1.610 -0.501  1.141 -1.297 -0.493 -0.479 -0.443
[7,] -1.197  0.115  0.855  1.186  1.067 -0.669 -0.443
[8,]  0.137 -0.426  0.848 -0.114 -1.343 -0.326 -0.443
[9,]  0.792 -0.356  0.862 -0.469  2.250 -0.155 -0.256
[10,] 0.803  0.148  1.055 -0.469  2.494 -0.460 -0.256
[11,] 0.594  0.140  1.038 -0.646 -1.150 -0.307 -0.256
[12,] 0.110 -0.249  1.113  0.300 -0.982 -0.726 -0.256
```

```
[13,] -0.598  0.048  0.987 -1.060 -0.465 -1.011  -0.443
[14,]  0.050 -0.743 -0.475  0.359 -0.313 -0.288  -0.443
[15,] -0.475 -0.079 -0.518  0.300  0.379 -0.421  -0.443
[16,]  0.310 -0.861  0.231 -0.410 -0.662 -0.631  -0.443
[17,] -0.297  0.142 -0.317  1.896 -0.056 -0.574  -0.443
[18,]  0.317 -0.532  0.095 -0.705 -0.769 -0.460  -0.443
[19,]  0.235 -0.451 -0.500  0.891 -0.462 -0.612  -0.443
[20,] -1.027  0.002  0.271  0.891  0.242 -0.631  -0.443
[21,] -0.546  0.021  0.052  0.477  0.017 -0.555  -0.443
[22,] -0.550  1.527 -0.041  0.181 -0.116 -0.383  -0.443
[23,] -0.465  0.412 -0.414  0.300  0.100 -0.212  -0.443
[24,]  0.418  0.254 -1.668 -0.942 -0.525  0.871  1.615
[25,] -1.848 -0.724 -1.507  1.837  1.603  2.431  3.672
[26,] -1.608 -1.051 -1.367  1.305  1.010  2.621  1.615
[27,] -1.423 -1.051 -1.485  1.127  1.273  1.708  1.615
[28,]  2.077  1.794 -1.586 -2.124 -1.115  0.168  -0.256
[29,] -1.749  4.149 -1.510  1.305  0.867 -0.345  -0.630
[30,]  0.125 -0.085 -1.485 -0.765 -0.370  2.583  1.802
attr(,"scaled:center")
      Totpor      PRA      PV      Densi      TenStr      CO2SBW      FirTemp
39.9333333  0.6629000 56.7466667  1.6993333  0.6186000  0.5816667 763.6666667
attr(,"scaled:scale")
      Totpor      PRA      PV      Densi      TenStr      CO2SBW      FirTemp
7.0795326  0.4806106 27.9061201  0.1691548  0.3153048  0.5259152 53.4649955
>
> write.table(Y, "normalizzate.csv",sep="\t", col.names=TRUE, row.names=TRUE, quote=TRUE)
```

Infine facciamo calcolare la matrice A (la matrice Rotation) e stampare un sommario

```
> eigen(MatrixCorr)
$values
[1] 3.2535323 1.5055507 1.1387154 0.6516515 0.2259838 0.1408874 0.0836789

$vectors
      [,1]      [,2]      [,3]      [,4]      [,5]      [,6]      [,7]
[1,] 0.41952730 0.4090418 0.10990336 0.31322074 -0.5122611 -0.5070450 0.164285158
[2,] 0.01780654 -0.4429432 0.74323926 0.24937193 0.2964313 -0.3160646 0.033540667
[3,] 0.41158535 -0.1139636 -0.52677278 0.12259468 0.5921070 -0.4114405 -0.072227697
[4,] -0.37599285 -0.4539186 -0.24367883 -0.36533155 -0.3559843 -0.5724149 0.073192547
[5,] -0.31740811 -0.2746593 -0.30365291 0.82533066 -0.1765009 0.1388565 0.002205038
[6,] -0.44597862 0.4327470 0.07954887 0.08714546 0.1416903 -0.3163847 -0.692629300
[7,] -0.46180279 0.3933811 -0.01070518 0.04944464 0.3480988 -0.1633589 0.693953321

> summary(princomp(Y))
Importance of components:
      Comp.1  Comp.2  Comp.3  Comp.4  Comp.5  Comp.6  Comp.7
Standard deviation  1.7734377 1.2063854 1.0491703 0.79368114 0.4673874 0.36904090 0.28441098
Proportion of Variance 0.4647903 0.2150787 0.1626736 0.09309307 0.0322834 0.02012678 0.01195413
Cumulative Proportion 0.4647903 0.6798690 0.8425426 0.93563570 0.9679191 0.98804587 1.00000000
```

o anche

```
> PCA <- princomp(Y)
> PCA
Call:
princomp(x = Y)
```

```
Standard deviations:
  Comp.1  Comp.2  Comp.3  Comp.4  Comp.5  Comp.6  Comp.7
1.7734377 1.2063854 1.0491703 0.7936811 0.4673874 0.3690409 0.2844110

7 variables and 30 observations.
> A <- unclass(loadings(PCA))
> A
      Comp.1  Comp.2  Comp.3  Comp.4  Comp.5  Comp.6  Comp.7
Totpor  0.41952730 0.4090418 0.10990336 0.31322074 -0.5122611 0.5070450 0.164285158
PRA     0.01780654 -0.4429432 0.74323926 0.24937193 0.2964313 0.3160646 0.033540667
PV      0.41158535 -0.1139636 -0.52677278 0.12259468 0.5921070 0.4114405 -0.072227697
Densi  -0.37599285 -0.4539186 -0.24367883 -0.36533155 -0.3559843 0.5724149 0.073192547
TenStr -0.31740811 -0.2746593 -0.30365291 0.82533066 -0.1765009 -0.1388565 0.002205038
CO2SBW -0.44597862 0.4327470 0.07954887 0.08714546 0.1416903 0.3163847 -0.692629300
FirTemp -0.46180279 0.3933811 -0.01070518 0.04944464 0.3480988 0.1633589 0.693953321
> round(PCA$sd^2, 3) # component variances
Comp.1 Comp.2 Comp.3 Comp.4 Comp.5 Comp.6 Comp.7
3.145  1.455  1.101  0.630  0.218  0.136  0.081

> screeplot(princomp(Y))
```

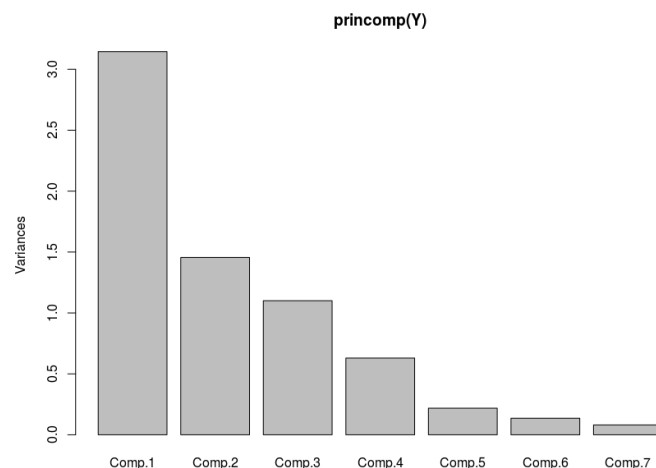


Figura 3.2: Varianza delle componenti principali

Vediamo come leggere questo output. Dalla prima riga del `unclass(loadings(PCA))` vediamo che le componenti principale sono numerate in ordine di deviazione standard decrescente. La prima componente principale ha la deviazione standard massima.

Dalla prima colonna della matrice di rotazione abbiamo che la prima componente principale Z_1 , che qui è indicata con `Comp.1` è pari a

$$\begin{aligned}
 Z_1 = & 0.41952730 \cdot \text{Totpor}_s + 0.01780654 \cdot \text{PRA}_s + 0.41158535 \cdot \text{PV}_s \\
 & - 0.37599285 \cdot \text{Densi}_s - 0.31740811 \cdot \text{TenStr}_s - 0.44597862 \cdot \text{CO2SBW}_s \\
 & - 0.46180279 \cdot \text{FirTemp}_s
 \end{aligned}$$

dove il pedice s indica che dobbiamo prendere il dato standardizzato e non nella sua forma originale. Possiamo ottenere la stessa informazione anche scrivendo

```
> Z1 <- A[,1]
> Z1
      Totpor      PRA      PV      Densi      TenStr      CO2SBW      FirTemp
0.41952730 0.01780654 0.41158535 -0.37599285 -0.31740811 -0.44597862 -0.46180279
```

Possiamo anche visualizzare (approssimiamo a 3 cifre decimali) il valore della prima componente principale su ciascun individuo del campione (numerati da 1 a 30)

```
> round(predict(PCA)[,1],3)
[1] 1.353 0.972 0.920 2.200 1.646 2.198 -0.430 1.218 0.330 0.482
     1.542 1.140 1.358 0.110 -0.254 1.060 -0.487 1.082 0.174
[20] -0.246 0.060 0.124 -0.203 -1.120 -5.388 -3.981 -3.562 1.447 -1.602 -2.139
```

o anche di tutte le sette componenti principali

```
> round(predict(PCA),3)
      Comp.1 Comp.2 Comp.3 Comp.4 Comp.5 Comp.6 Comp.7
[1,] 1.353 0.366 -0.227 -0.197 0.526 -0.241 -0.142
[2,] 0.972 0.534 -0.583 0.501 -0.162 0.573 -0.392
[3,] 0.920 -0.127 -0.241 -0.026 0.054 0.542 -0.110
[4,] 2.200 1.496 0.139 -0.380 -0.152 0.506 -0.035
[5,] 1.646 0.525 -0.622 0.508 0.351 -0.283 0.004
[6,] 2.198 1.093 -0.364 0.522 0.029 0.230 0.093
[7,] -0.430 -1.933 -1.158 0.125 0.294 0.028 -0.009
[8,] 1.218 0.254 -0.334 -1.077 0.383 0.229 -0.146
[9,] 0.330 -0.190 -1.210 2.267 -0.342 -0.028 -0.043
[10,] 0.482 -0.629 -1.035 2.595 -0.170 0.086 0.173
[11,] 1.542 0.439 0.108 -0.404 0.652 0.423 0.013
[12,] 1.140 -0.252 -0.589 -0.887 0.404 0.471 0.274
[13,] 1.358 -0.381 -0.226 -0.161 1.067 -0.816 0.146
[14,] 0.110 0.028 -0.307 -0.664 -0.794 -0.319 -0.064
[15,] -0.254 -0.697 -0.055 -0.087 -0.474 -0.566 -0.036
[16,] 1.060 0.403 -0.472 -0.562 -0.258 -0.434 0.104
[17,] -0.487 -1.416 -0.246 -0.907 -0.894 0.603 0.208
[18,] 1.082 0.513 -0.037 -0.461 -0.096 -0.483 -0.015
[19,] 0.174 -0.363 -0.167 -0.882 -1.027 0.079 0.240
[20,] -0.246 -1.370 -0.590 -0.490 0.084 -0.204 0.007
[21,] 0.060 -0.874 -0.233 -0.390 -0.089 -0.226 0.019
[22,] 0.124 -1.287 1.062 -0.014 0.458 0.113 -0.065
[23,] -0.203 -0.755 0.358 -0.161 -0.193 -0.258 -0.171
[24,] -1.120 1.833 1.554 0.056 -0.013 -0.321 0.644
[25,] -5.388 0.959 -0.728 0.101 0.526 0.412 0.784
[26,] -3.981 0.863 -0.671 -0.268 -0.006 -0.010 -0.798
[27,] -3.562 0.566 -0.697 -0.022 -0.284 -0.392 -0.139
[28,] 1.447 1.479 3.269 0.761 -0.583 -0.082 0.064
[29,] -1.602 -3.609 3.085 0.479 0.346 0.217 -0.139
[30,] -2.139 2.533 1.218 0.124 0.362 0.151 -0.470
```