



Università degli Studi di Salerno
Dipartimento di Informatica
Corso di Laurea Triennale in Informatica

Progetto Calcolo Probabilità Statistica Matematica
(CPSM)

Indagine Statistica sulle Morti in incidenti stradali

Tozza Gennaro Carmine
Matricola: 0512120382

Anno Accademico 2024-2025

Indice

| | | |
|----------|--|-----------|
| 1 | Introduzione | 2 |
| 1.1 | Problematica | 2 |
| 1.2 | Scopo del progetto | 2 |
| 2 | Tabelle delle frequenze | 4 |
| 3 | Rappresentazione dei dati mediante grafici | 6 |
| 4 | Indici di posizione | 11 |
| 4.1 | Media campionaria | 11 |
| 4.2 | Mediana campionaria | 12 |
| 4.3 | Moda campionaria | 12 |
| 5 | Indici di variabilità | 13 |
| 5.1 | Varianza campionaria | 13 |
| 5.2 | Deviazione standard campionaria | 13 |
| 5.3 | Scarto medio assoluto | 14 |
| 5.4 | Ampiezza del campo di variazione | 14 |
| 5.5 | Coefficiente di variazione | 15 |
| 6 | Indici di forma | 16 |
| 6.1 | Indice di asimmetria | 16 |
| 6.2 | Indice di curtosi | 18 |
| 7 | Percentili campionari | 20 |
| 7.1 | Box Plot | 21 |
| 7.2 | Disuguaglianza di Chebyshev | 23 |
| 8 | Dati Bivariati | 25 |
| 8.1 | Diagramma a dispersione | 25 |
| 8.2 | Coefficiente di correlazione campionario | 26 |

Capitolo 1

Introduzione

1.1 Problematica

Gli **incidenti stradali** costituiscono una delle principali emergenze di sanità pubblica, in quanto responsabili ogni anno di un elevato numero di decessi, in particolare tra i giovani, oltre al drammatico impatto umano e psicologico sulle loro famiglie.

1.2 Scopo del progetto

Il progetto consiste nel realizzare un'indagine statistica ¹sugli incidenti stradali verificatisi sulla rete stradale del territorio nazionale, tra il 2010 e il 2023 verbalizzati da un'autorità di Polizia o dai Carabinieri, avvenuti su una strada aperta alla circolazione pubblica e che hanno causato morti (entro il 30° giorno) con il coinvolgimento di almeno un veicolo.

La rilevazione è condotta correntemente dall'Istat, con la compartecipazione dell'ACI (Automobile Club d'Italia) e di numerosi Enti pubblici istituzionali.

Per l'analisi dei dati è stato scelto l'ambiente di calcolo statistico **R**.

R fornisce un'ampia varietà di tecniche statistiche (modellazione lineare e non lineare, test statistici classici, analisi delle serie temporali, classificazione, ...) e grafiche ed è altamente estensibile.

Uno dei punti di forza di R è la facilità con cui possono essere prodotti grafici ben progettati e di qualità per la pubblicazione, compresi simboli matematici e formule se necessario.

¹<https://siqua.istat.it/SIQual/visualizza.do?id=7777778&refresh=true&language=IT>

Per semplificare l'analisi, si è scelto di lavorare non sull'intero dataset, ma su un sottoinsieme filtrato di dati, relativo alle **morti per incidenti stradali che riguardano solo i conducenti di età compresa tra i 21 e i 24 anni**.

Per approfondire l'analisi con dati dettagliati e specifici, è possibile consultare il dataset direttamente sul sito dell'ISTAT al seguente link: https://esploradati.istat.it/databrowser/#/it/dw/categories/IT1,Z0810HEA,1.0/HEA_ROAD/IT1,41_270_DF_DCIS_MORTIFERITISTR1_1,1.0

Il dataset in formato **CSV(Comma-Separated Values)** già filtrato è stato ottenuto dalla fonte ISTAT tramite il link indicato, assicurando così l'affidabilità dei dati.

Il formato scelto (CSV) permette un'agevole manipolazione dei dati, essendo compatibile con la maggior parte dei software statistici e dei fogli di calcolo, ottimizzando l'analisi e la visualizzazione delle informazioni.

Successivamente, utilizzando il software R, sono stati applicati ulteriori criteri di selezione per estrarre esclusivamente i dati rilevanti ai fini dell'analisi.

```

1 # inclusione librerie
2 library("tidyverse")
3 require("tidyverse")
4 library("dplyr")
5 library(moments)
6
7 # viene caricato il dataset
8 dati <- read.csv("dati_istat.csv")
9
10 # filtraggio dati
11 dati <- dati %>% select(Intersezione, TIME_PERIOD,
    Osservazione)

```

| Intersezione | '10 | '11 | '12 | '13 | '14 | '15 | '16 | '17 | '18 | '19 | '20 | '21 | '22 | '23 |
|---------------|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| Incrocio | 55 | 45 | 42 | 23 | 34 | 24 | 28 | 25 | 23 | 13 | 15 | 16 | 21 | 10 |
| Rotatoria | 7 | 4 | 2 | 5 | 5 | 1 | 5 | 3 | 1 | 2 | 1 | 2 | — | 3 |
| Rettilineo | 92 | 102 | 89 | 82 | 96 | 95 | 66 | 70 | 60 | 72 | 66 | 69 | 69 | 74 |
| Curva | 58 | 51 | 48 | 46 | 45 | 44 | 41 | 36 | 38 | 42 | 23 | 40 | 27 | 38 |
| Dosso/Pend. | 2 | 7 | 2 | 4 | 2 | 3 | 5 | 1 | 2 | 3 | 2 | 4 | 3 | 1 |
| Galleria | 1 | 1 | 1 | 1 | — | 1 | 1 | 1 | — | — | 3 | 1 | — | — |
| Totale | 215 | 210 | 184 | 161 | 182 | 168 | 146 | 136 | 124 | 132 | 110 | 132 | 120 | 126 |

Capitolo 2

Tabelle delle frequenze

```
1 dati <- dati %>%
2   filter(Intersezione != "Totale")
3
4 # 1. Frequenze assolute per intersezione dell'incidente
5 freq_assolute <- aggregate(Osservazione ~ Intersezione, data
6   =dati, sum)
7 colnames(freq_assolute) <- c("Intersezione", "Frequenza_
8   Assoluta")
9 print(freq_assolute)
10
11 # 2. Frequenze relative per intersezione dell'incidente
12 totale <- sum(freq_assolute$Frequenza_Assoluta)
13 freq_assolute$Frequenza_Relativa <- freq_assolute$Frequenza_
14   Assoluta / totale
15 print(freq_assolute)
16
17 # 3. Frequenze cumulate assolute
18 freq_assolute <- freq_assolute[order(-freq_assolute$
19   Frequenza_Assoluta),]
20 freq_assolute$Frequenza_Cumulata_Assoluta <- cumsum(freq_
21   assolute$Frequenza_Assoluta)
22 print(freq_assolute)
23
24 # 4. Frequenze cumulate relative
25 freq_assolute$Frequenza_Cumulata_Relativa <- freq_assolute$
26   Frequenza_Cumulata_Assoluta / totale
27 print(freq_assolute)
28
29 # Per esportare i risultati
30 #write.csv(freq_assolute, "tabella_frequenza.csv", row.names
31   = FALSE)
```

L'output generato dal codice R riportato è sintetizzato nella seguente tabella:

| Tipo intersezione | Freq. Assoluta | Freq. Relativa | Freq. Cum. Assoluta | Freq. Cum. Relativa |
|--------------------------|---------------------------|---------------------------|------------------------------------|------------------------------------|
| Rettilineo | 1102 | 0.5135 | 1102 | 0.5135 |
| Curva | 577 | 0.2689 | 1679 | 0.7824 |
| Incrocio | 374 | 0.1743 | 2053 | 0.9567 |
| Dosso/Pendenza/Strettoia | 41 | 0.0191 | 2094 | 0.9758 |
| Rotatoria | 41 | 0.0191 | 2135 | 0.9949 |
| Galleria | 11 | 0.0051 | 2146 | 1.0000 |

L'analisi di frequenza, fornisce una chiara mappatura della distribuzione degli incidenti stradali analizzati, per un totale di 2.146 casi.

Dai dati emerge in modo netto che la tipologia di strada più a rischio è il rettilineo, che da solo rappresenta oltre la metà degli incidenti totali (1.102 casi, pari al 51,4%). Questo dato è di particolare rilevanza, in quanto concentra la maggior parte del rischio in un contesto spesso percepito come meno pericoloso.

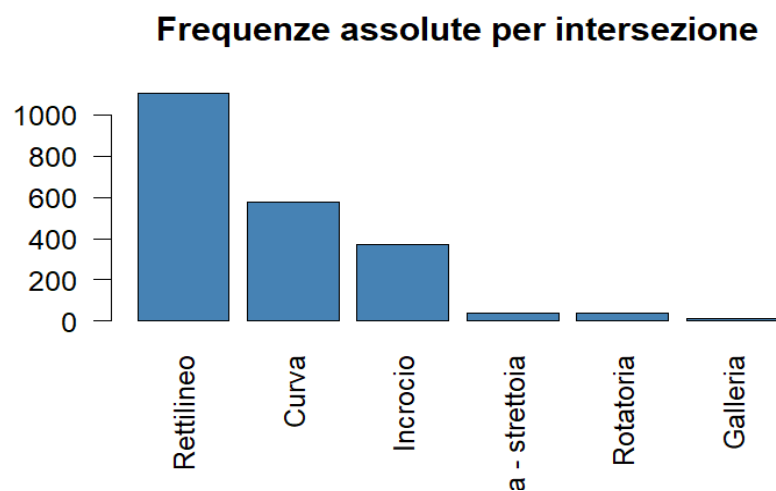
L'analisi cumulativa rafforza questa evidenza. Le prime due categorie, Rettilineo e Curva, coprono complessivamente il 78,2% di tutti gli incidenti. Aggiungendo gli incroci, si arriva a spiegare oltre il 95% del fenomeno (95,7%).

Le restanti tipologie di infrastruttura (dosso/pendenza/strettoia, rotatoria, galleria) hanno un'incidenza marginale, contribuendo complessivamente per meno del 5% al totale degli eventi.

Capitolo 3

Rappresentazione dei dati mediante grafici

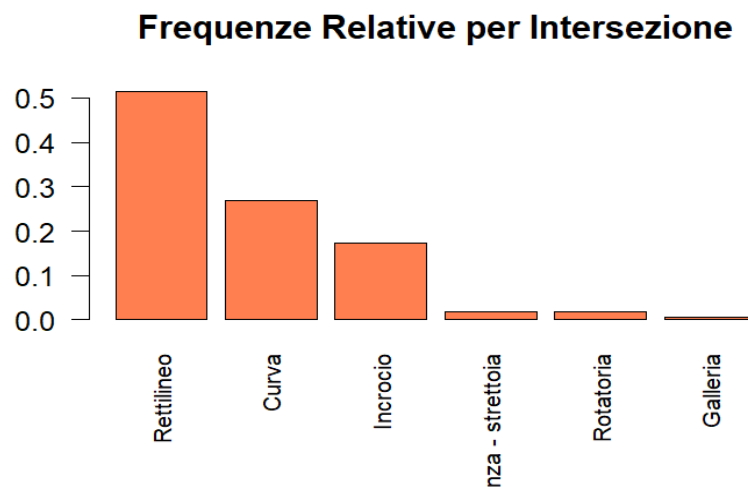
```
1 # 1. Grafico frequenza assolute
2 barplot(freq_assolute$Frequenza_Assoluta,
3         names.arg = freq_assolute$Intersezione,
4         main = "Frequenze assolute per intersezione",
5         xlab = "",
6         ylab = "",
7         col = "steelblue",
8         las = 2) # Etichette verticali
```



```

1 # 2. Grafico frequenze relative
2 barplot(freq_assolute$Frequenza_Relativa,
3         names.arg = freq_assolute$Intersezione,
4         main = "Frequenze_Relative_per_Intersezione",
5         xlab = "",
6         ylab = "",
7         col = "coral",
8         las = 2,
9         cex.names = 0.8)

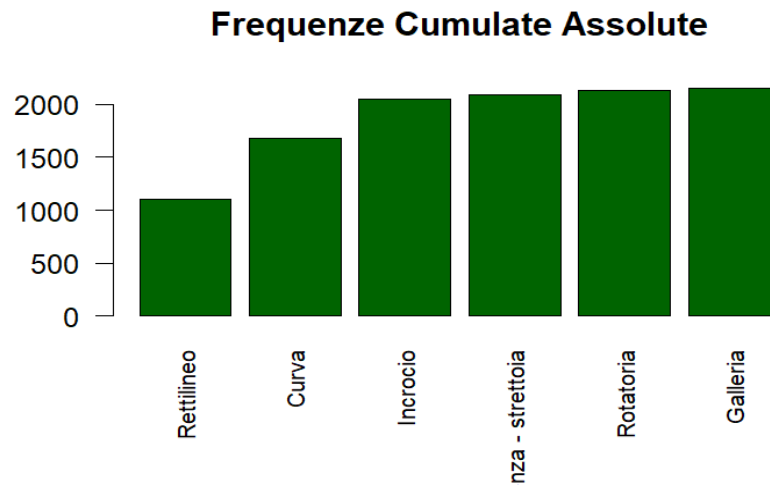
```



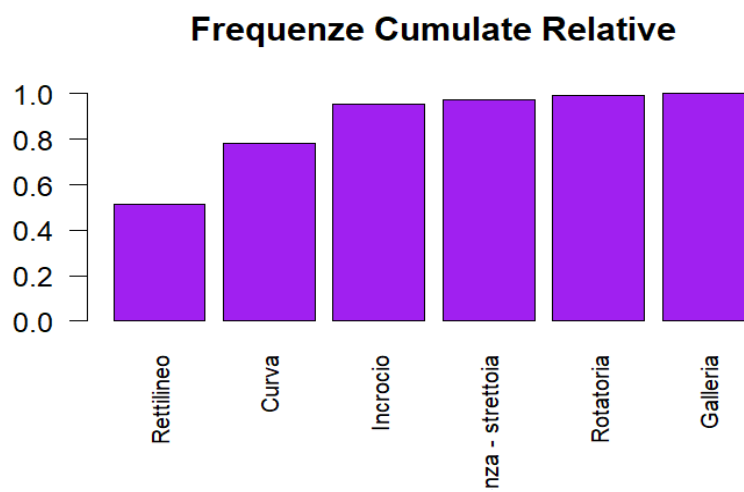
```

1 # 3. Grafico frequenze cumulate assolute
2 barplot(freq_assolute$Frequenza_Cumulata_Assoluta,
3         names.arg = freq_assolute$Intersezione,
4         main = "Frequenze_Cumulate_Assolute",
5         xlab = "",
6         ylab = "",
7         col = "darkgreen",
8         las = 2,
9         cex.names = 0.8)

```

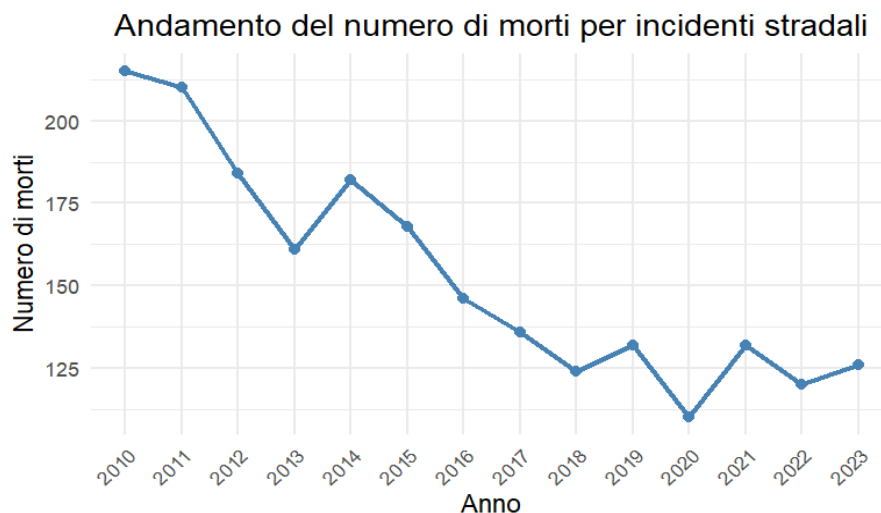


```
1 # 4. Grafico frequenze cumulate relative
2 barplot(freq_assolute$Frequenza_Cumulata_Relativa,
3         names.arg = freq_assolute$Intersezione,
4         main = "Frequenze_Cumulate_Relative",
5         xlab = "",
6         ylab = "",
7         col = "purple",
8         las = 2,
9         cex.names = 0.8)
```



Altri grafici:

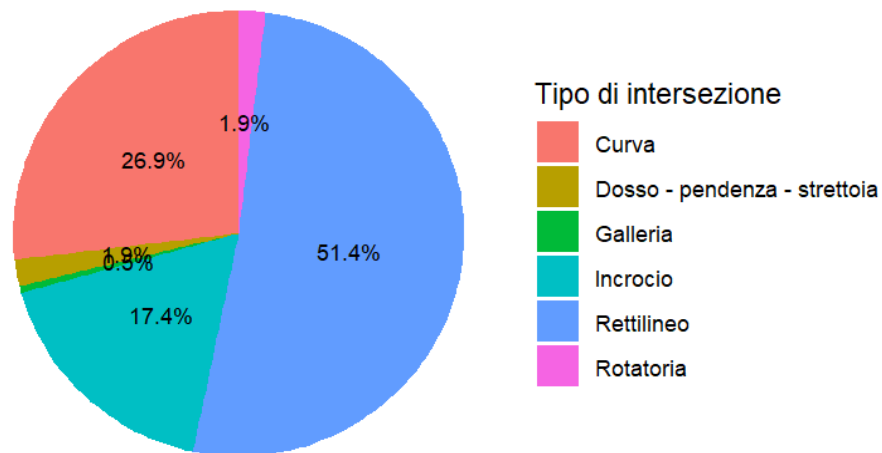
```
1 # Converti TIME_PERIOD in fattore per mantenere l'ordine
  originale
2 dati$TIME_PERIOD <- factor(dati$TIME_PERIOD, levels = unique
  (dati$TIME_PERIOD))
3
4 # Calcola il numero totale di morti per anno
5 morti_per_anno <- dati %>%
6   group_by(TIME_PERIOD) %>%
7   summarise(Totale_Morti = sum(Osservazione, na.rm = TRUE))
8
9 # Crea il grafico a linee con tutti gli anni visibili
10 ggplot(morti_per_anno, aes(x = TIME_PERIOD, y = Totale_Morti
  , group = 1)) +
11   geom_line(color = "steelblue", size = 1) +
12   geom_point(color = "steelblue", size = 2) +
13   labs(title = "Andamento del numero di morti per incidenti
  stradali",
14         x = "Anno",
15         y = "Numero di morti") +
16   theme_minimal() +
17   theme(plot.title = element_text(hjust = 0.5),
18         axis.text.x = element_text(angle = 45, hjust = 1,
19                                     size = 8)) + # Riduci dimensione testo
19   scale_x_discrete(breaks = levels(morti_per_anno$TIME_
  PERIOD)) # Mostra tutti i valori
```



```

1 # Crea il grafico a torta
2 ggplot(freq_assolute, aes(x = "", y = Frequenza_Assoluta,
   fill = Intersezione)) +
3   geom_bar(width = 1, stat = "identity") +
4   coord_polar("y", start = 0) +
5   labs(title = "",
6        fill = "Tipo di intersezione") +
7   theme_void() +
8   theme(plot.title = element_text(hjust = 0.5, size = 14,
9         face = "bold"),
10         legend.position = "right") +
11   geom_text(aes(label = paste0(round(Frequenza_Assoluta/sum(
12     Frequenza_Assoluta)*100, 1), "%")),
13             position = position_stack(vjust = 0.5),
14             size = 3)

```



Capitolo 4

Indici di posizione

In questo capitolo verranno calcolati gli **indici di posizione** sulla variabile numerica `Osservazione`.

Questa variabile rappresenta il numero di morti registrato per una specifica combinazione di tipo di intersezione e anno, limitatamente alla fascia di età dei conducenti tra 21 e 24 anni.

Gli indici di posizione (o di tendenza generale) sono misure che consentono di sintetizzare i dati osservati con un solo valore numerico che sia rappresentativo dei dati stessi. Gli indici di posizione più adoperati sono tre:

- media
- mediana
- moda

4.1 Media campionaria

La **media campionaria** è la somma di tutte le osservazioni divisa per il numero di esse. Fornisce una misura del valore centrale della distribuzione.

```
1 media_generale <- mean(dati$Osservazione, na.rm = TRUE)
2 print(paste("Media_campionaria_generale_delle_osservazioni:",
, round(media_generale, 2)))
```

Output:

```
"Media campionaria generale delle osservazioni: 27.51"
```

Questo valore indica che, in media, per ogni specifica combinazione di tipo di intersezione e anno considerata nel dataset filtrato, si sono registrati circa 27.51 decessi.

4.2 Mediana campionaria

La **mediana** è un numero che precede tanti dati quanti ne segue, ossia il valore centrale dei dati riordinati in senso crescente.

```
1 mediana_generale <- median(dati$Osservazione, na.rm = TRUE)
2 print(paste("Mediana generale delle osservazioni:", mediana_
             generale))
```

Output:

```
"Mediana generale delle osservazioni: 15.5"
```

Il fatto che la mediana (15.5) sia inferiore alla media (27.51) suggerisce una distribuzione asimmetrica a destra.

4.3 Moda campionaria

La **moda** è il valore (o i valori, in caso di distribuzioni multimodali) che appare più frequentemente in un insieme di dati.

```
1 find_mode <- function(x) {
2   u <- unique(x)
3   tab <- tabulate(match(x, u))
4   u[tab == max(tab)]
5 }
6
7 find_mode(dati$Osservazione)
```

Output:

```
1
```

L'output '1', indica che il numero di decessi più frequentemente osservato per una singola combinazione intersezione/anno è 1. Se ci fossero più valori con la stessa frequenza massima, sarebbero elencati tutti.

Capitolo 5

Indici di variabilità

Di seguito calcoliamo gli **indici di variabilità**(o di **dispersione**), che descrivono la variabilità dei dati osservati e consentono di valutare l'informazione fornita dall'indice di posizione utilizzato, dando dei dati più accurati.

5.1 Varianza campionaria

La **varianza campionaria** (s^2) misura la dispersione media quadratica dei dati attorno alla media campionaria. È espressa nell'unità di misura dei dati al quadrato.

```
1 # Calcolo della varianza sulla variabile 'Osservazione'.
2 varianza_campionaria <- var(dati$Osservazione, na.rm = TRUE)
3 print(paste("Varianza_campionaria_delle_osservazioni:",
              round(varianza_campionaria, 2)))
```

Output:

```
"Varianza campionaria delle osservazioni: 859.52"
```

Un valore elevato della varianza indica una notevole dispersione dei dati attorno alla media.

5.2 Deviazione standard campionaria

La **deviazione standard campionaria** (s) è la radice quadrata della varianza campionaria. Fornisce una misura della dispersione media dei dati attorno alla media, espressa nella stessa unità di misura dei dati originali, rendendola più interpretabile della varianza.

```
1 # Calcolo della deviazione standard sulla variabile '
  Osservazione'.
2 dev_std_campionaria <- sd(dati$Osservazione, na.rm = TRUE)
3 print(paste("Deviazione standard campionaria delle
  osservazioni:", round(dev_std_campionaria, 2)))
```

Output (radice quadrata della varianza stimata):

```
"Deviazione standard campionaria delle osservazioni: 29.61"
```

Questo valore indica che, mediamente, i singoli conteggi di decessi si discostano dalla media campionaria (27.51) di circa 29.61 unità.

5.3 Scarto medio assoluto

Lo **scarto medio assoluto** è dato dalla media aritmetica degli scarti assoluti dalla media campionaria. Come la deviazione standard, misura la dispersione media, ma è meno sensibile ai valori anomali perché non eleva al quadrato gli scarti.

```
1 # Calcolo dello scarto medio assoluto dalla media per '
  Osservazione'.
2 media_oss <- mean(dati$Osservazione, na.rm = TRUE)
3 scarto_medio_assoluto <- mean(abs(dati$Osservazione - media_
  oss), na.rm = TRUE)
4 print(paste("Scarto medio assoluto (dalla media) delle
  osservazioni:", round(scarto_medio_assoluto, 2)))
```

Output (stima):

```
"Scarto medio assoluto (dalla media) delle osservazioni: 25.13"
```

In media, le osservazioni si discostano (in valore assoluto) dalla media di circa 25.13 decessi.

5.4 Ampiezza del campo di variazione

L'**ampiezza del campo di variazione** (o semplicemente "range") è la differenza tra il valore massimo e il valore minimo osservato nel dataset. È una misura di variabilità semplice ma molto sensibile ai valori estremi.

```

1 # Calcolo del minimo, massimo e ampiezza del campo di
  variazione per 'Osservazione'.
2 min_oss <- min(dati$Osservazione, na.rm = TRUE)
3 max_oss <- max(dati$Osservazione, na.rm = TRUE)
4 ampiezza_variazione <- max_oss - min_oss
5
6 print(paste("Valore_minimo_delle_osservazioni:", min_oss))
7 print(paste("Valore_massimo_delle_osservazioni:", max_oss))
8 print(paste("Ampiezza_del_campo_di_variazione_delle_
  osservazioni:", ampiezza_variazione))

```

Output:

```

"Valore minimo delle osservazioni: 1"
"Valore massimo delle osservazioni: 102"
"Ampiezza del campo di variazione delle osservazioni: 101"

```

5.5 Coefficiente di variazione

Il **coefficiente di variazione (CV)** è una misura di variabilità relativa, data dal rapporto tra la deviazione standard e la media (in valore assoluto).

```

1 # Calcolo del coefficiente di variazione per 'Osservazione'.
2 media_oss <- mean(dati$Osservazione, na.rm = TRUE)
3 dev_std_oss <- sd(dati$Osservazione, na.rm = TRUE)
4 coeff_variazione <- (dev_std_oss / abs(media_oss)) * 100 #
  abs() per media se potesse essere negativa
5 print(paste("Coefficiente_di_variazione_delle_osservazioni:"
  , round(coeff_variazione, 2), "%"))

```

Output(in percentuale):

```

"Coefficiente di variazione delle osservazioni: 107.6 %"

```

Un CV del 107.6% indica una variabilità molto elevata rispetto alla media. Questo è coerente con il fatto che la deviazione standard (29.61) è addirittura leggermente superiore alla media (27.51), suggerendo una notevole eterogeneità nei conteggi dei decessi.

Capitolo 6

Indici di forma

Gli **indici di forma** misurano caratteristiche relative alla forma della distribuzione dati. I più usati sono l'indice di asimmetria e l'indice di curtosi.

6.1 Indice di asimmetria

L'**indice di asimmetria** (**skewness**) misura il grado di simmetria di una distribuzione di dati attorno alla sua media.

- Un valore di **skewness** > 0 indica un'asimmetria positiva (coda a destra): la distribuzione ha una coda che si estende maggiormente verso i valori più alti. La media è tipicamente maggiore della mediana.
- Un valore di **skewness** < 0 indica un'asimmetria negativa (coda a sinistra): la distribuzione ha una coda che si estende maggiormente verso i valori più bassi. La media è tipicamente minore della mediana.
- Un valore di **skewness** ≈ 0 indica una distribuzione approssimativamente simmetrica (come la distribuzione normale).

```
1 # Calcolo dell'indice di asimmetria per 'Osservazione'.
2 indice_asimmetria <- skewness(dati$Osservazione, na.rm =
  TRUE)
3 print(paste("Indice di asimmetria (Skewness) delle
  osservazioni:", round(indice_asimmetria, 2)))
```

Output:

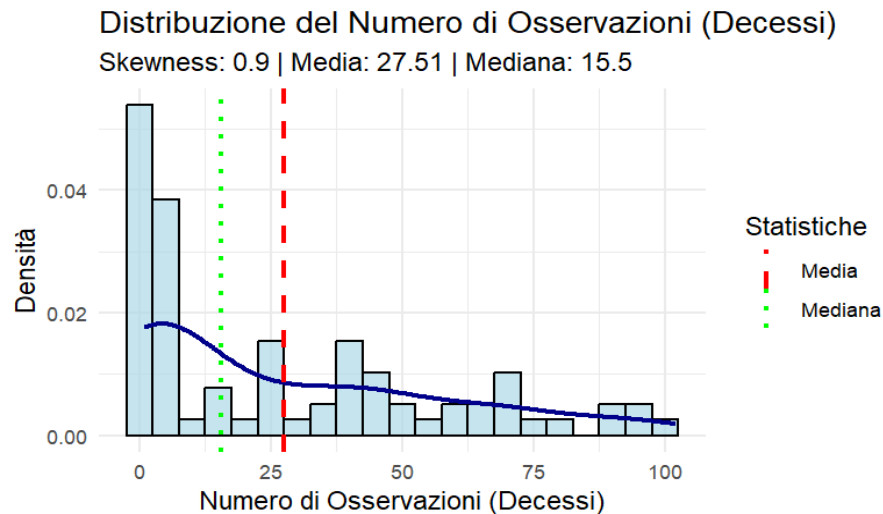
```
"Indice di asimmetria (Skewness) delle osservazioni: 0.9"
```

Un valore di 0.9 indica una moderata asimmetria positiva (coda a destra). Questo è coerente con l'osservazione precedente che la media (27.51) è maggiore della mediana (15.5), suggerendo che ci sono alcune osservazioni con un numero

di decessi relativamente alto che "stirano" la coda destra della distribuzione.

Per visualizzare graficamente questa asimmetria, possiamo generare un istogramma con una curva di densità sovrapposta e linee verticali per la media e la mediana.

```
1 # Creazione del grafico
2 skewness_plot <- ggplot(dati, aes(x = Osservazione)) +
3   geom_histogram(aes(y = ..density..), binwidth = 5, fill =
4     "lightblue", color = "black", alpha = 0.7) +
5   geom_density(color = "darkblue", linewidth = 1) +
6   geom_vline(aes(xintercept = media_oss, color = "Media"),
7     linetype = "dashed", linewidth = 1) +
8   geom_vline(aes(xintercept = mediana_generale, color = "
9     Mediana"), linetype = "dotted", linewidth = 1) +
10  labs(title = "Distribuzione del Numero di Osservazioni (Decessi)",
11    subtitle = paste0("Skewness:", round(indice_
12      asimmetria, 2),
13      "\nMedia:", round(media_oss, 2),
14      "\nMediana:", round(mediana_
15        generale, 2)),
16    x = "Numero di Osservazioni (Decessi)",
17    y = "Densità") +
18  scale_color_manual(name = "Statistiche", values = c("Media
19    " = "red", "Mediana" = "green")) +
20  theme_minimal()
21 print(skewness_plot)
```



6.2 Indice di curtosi

L'**indice di curtosi** misura l'appiattimento di una distribuzione, in particolare la concentrazione dei dati attorno al valore centrale e la "pesantezza" delle code, rispetto a una distribuzione normale.

- **Eccesso di curtosi** > 0 (Curtosi > 3): Distribuzione leptocurtica. È più appuntita al centro e ha code più pesanti (maggior probabilità di valori estremi) rispetto a una normale.
- **Eccesso di curtosi** < 0 (Curtosi < 3): Distribuzione platycurtica. È più piatta al centro e ha code più leggere rispetto a una normale.
- **Eccesso di curtosi** ≈ 0 (Curtosi ≈ 3): Distribuzione normocurtica. Ha una forma simile a quella di una distribuzione normale in termini di appiattimento.

```
1 # Calcolo dell'indice di curtosi (eccesso di curtosi) per '
  Osservazione'.
2 indice_curtosi <- kurtosis(dati$Osservazione, na.rm = TRUE)
3 print(paste("Indice di curtosi (eccesso di curtosi) delle
  osservazioni:", round(indice_curtosi, 2)))
4 # La curtosi "classica" sarebbe approssimativamente: round(
  indice_curtosi, 2) + 3
```

Output:

```
"Indice di curtosi (eccesso di curtosi) delle osservazioni: 2.64"
```

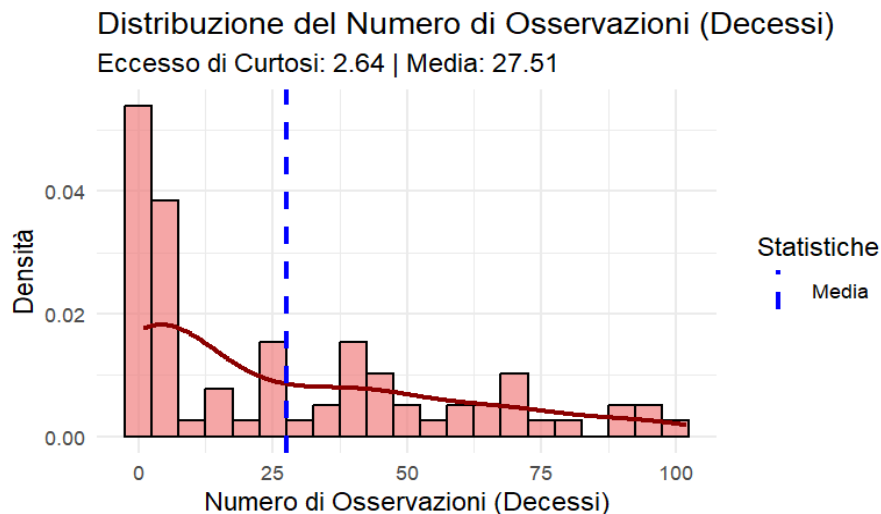
Un eccesso di curtosi di circa 2.64 (corrispondente a una curtosi "classica" di circa $2.64 + 3 = 5.64$) indica una distribuzione marcatamente **leptocurtica**.

Questo significa che la distribuzione del numero di decessi per combinazione intersezione/anno presenta un picco più pronunciato attorno alla moda/media-na e code più "pesanti" rispetto a una distribuzione normale.

Ciò implica una maggiore frequenza di valori vicini al centro, ma anche una maggiore probabilità di osservare valori estremi (molto alti, data l'asimmetria positiva) rispetto a quanto ci si aspetterebbe da una distribuzione normale.

Graficamente si ha

```
1 # Creazione del grafico per la curtosi
2 kurtosis_plot <- ggplot(dati, aes(x = Osservazione)) +
3   geom_histogram(aes(y = ..density..), binwidth = 5, fill =
4     "lightcoral", color = "black", alpha = 0.7) +
5   geom_density(color = "darkred", linewidth = 1) +
6   geom_vline(aes(xintercept = media_oss, color = "Media"),
7     linetype = "dashed", linewidth = 1) +
8   labs(title = "Distribuzione del Numero di Osservazioni (Decessi)",
9     subtitle = paste0("Eccesso di Curtosi: ", round(
10       indice_curtosi, 2),
11       " | Media: ", round(media_oss, 2)),
12     x = "Numero di Osservazioni (Decessi)",
13     y = "Densità") +
14   scale_color_manual(name = "Statistiche", values = c("Media"
15     = "blue")) +
16   theme_minimal()
17 print(kurtosis_plot)
```



Capitolo 7

Percentili campionari

Un **percentile** k-esimo di un campione di dati è un valore che è maggiore di una percentuale k dei dati e minore della restante percentuale.

Il 25-esimo percentile si dice primo **quartile**, il 50-esimo percentile si dice mediana campionaria o secondo quartile, il 75-esimo percentile si dice terzo quartile e si denotano con

$$Q_1, Q_2, Q_3$$

La differenza tra il terzo e il primo quartile è chiamata **scarto interquartile**.

```
1 # Calcolo dei quartili per la variabile 'Osservazione'
2 quartili_osservazione <- quantile(dati$Osservazione, probs =
  c(0.25, 0.50, 0.75), na.rm = TRUE)
3 Q1 <- quartili_osservazione[1]
4 Q2_mediana <- quartili_osservazione[2] # Coincide con la
  mediana già calcolata
5 Q3 <- quartili_osservazione[3]
6 IQR_osservazione <- Q3 - Q1
7
8 print(paste("Primo_Quartile_(Q1)_delle_osservazioni:", round
  (Q1, 2)))
9 print(paste("Secondo_Quartile_(Q2)_delle_osservazioni:",
  round(Q2_mediana, 2)))
10 print(paste("Terzo_Quartile_(Q3)_delle_osservazioni:", round
  (Q3, 2)))
11 print(paste("Scarto_Interquartile_(IQR)_delle_osservazioni:"
  , round(IQR_osservazione, 2)))
12
13 print("Sommario_statistico_della_variabile_Osservazione:")
14 summary_stats <- summary(dati$Osservazione, na.rm = TRUE)
15 print(summary_stats)
```

Output:

```
"Primo Quartile (Q1) delle osservazioni: 2"
"Secondo Quartile (Q2) delle osservazioni: 15.5"
"Terzo Quartile (Q3) delle osservazioni: 45"
"Scarto Interquartile (IQR) delle osservazioni: 43"
```

```
"Sommario statistico della variabile Osservazione:"
Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
1.00    2.00   15.50   27.51   45.00   102.00
```

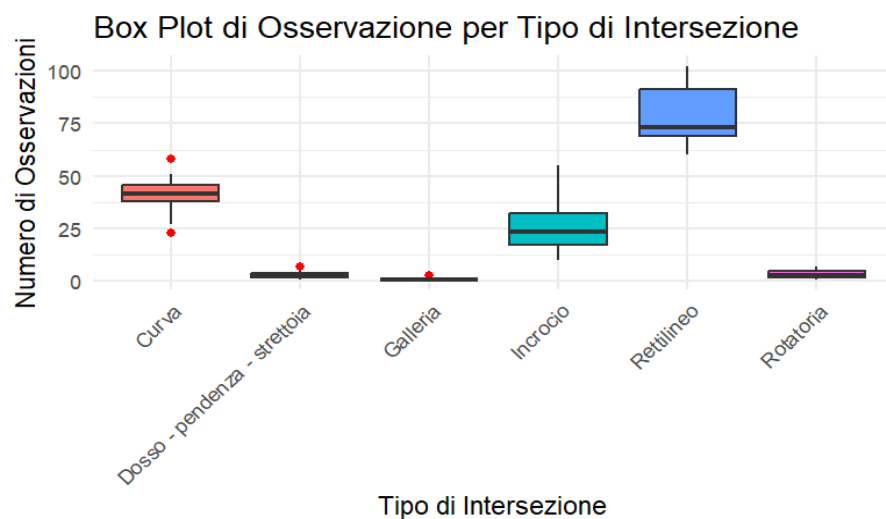
Dall'output del sommario, il 25% delle osservazioni (combinazioni intersezione/anno) ha un numero di decessi inferiore o uguale a 2, il 50% inferiore o uguale a 15.5, e il 75% inferiore o uguale a 45. Lo scarto interquartile indica che il 50% centrale dei dati si distribuisce in un intervallo di ampiezza 43.

7.1 Box Plot

Uno strumento grafico utile per la visualizzare alcuni degli indici rappresentativi dei dati è il **box plot**.

Si ottiene sovrapponendo ad una linea orizzontale che va dal minore al maggiore dei dati, un rettangolo che va dal primo al terzo quartile, con una linea verticale che lo divide al livello del secondo quartile.

```
1 # Box plot di Osservazione per tipo di Intersezione
2 Assicurarsi che la colonna Intersezione sia un fattore per
  un ordinamento corretto (se necessario)
3 dati$Intersezione <- factor(dati$Intersezione, levels = c("
  Rettilineo", "Curva", ...)) # opzionale
4 boxplot_oss_intersezione <- ggplot(dati, aes(x =
  Intersezione, y = Osservazione, fill = Intersezione)) +
5 geom_boxplot(outlier.colour = "red", na.rm = TRUE) +
6 labs(title = "Box Plot di Osservazione per Tipo di
  Intersezione",
7 x = "Tipo di Intersezione",
8 y = "Numero di Osservazioni") +
9 theme_minimal() +
10 theme(axis.text.x = element_text(angle = 45, hjust = 1), #
  Ruota etichette asse x
11 legend.position = "none") # Rimuove la legenda se i colori
  sono distintivi
12
13 print(boxplot_oss_intersezione)
```



Il grafico permette di confrontare la distribuzione del numero di decessi per ciascun tipo di intersezione. Si possono notare differenze nelle mediane, nella variabilità e nella presenza di ¹outlier tra le categorie. Ad esempio, i rettilinei potrebbero mostrare una mediana più alta e una maggiore dispersione.

¹Si dice outlier un dato anomalo, ossia molto distante dagli altri dati. Vengono solitamente rappresentati come punti individuali.

7.2 Disuguaglianza di Chebyshev

La **disuguaglianza di Chebyshev** fornisce un limite inferiore alla proporzione di dati che si trovano entro un certo numero (k) di deviazioni standard dalla media campionaria (s), indipendentemente dalla forma della distribuzione. La disuguaglianza afferma che almeno $1 - 1/k^2$ dei valori si trova nell'intervallo $[\bar{x} - ks, \bar{x} + ks]$, per $k > 1$.

```
1 # Disuguaglianza di Chebyshev
2 # media_generale e dev_std_campionaria sono stati calcolati
  precedentemente
3 # media_generale = 27.51, dev_std_campionaria = 29.61
4 # N_validi = sum(!is.na(dati$Osservazione)) # Numero di
  osservazioni non NA (79)
5 N_validi <- sum(!is.na(dati$Osservazione))
6
7 # Per k = 2 deviazioni standard
8 k2 <- 2
9 lim_inf_k2 <- media_generale - k2 * dev_std_campionaria
10 lim_sup_k2 <- media_generale + k2 * dev_std_campionaria
11 prop_teorica_k2 <- 1 - 1/k2^2
12 prop_osservata_k2 <- sum(dati$Osservazione >= lim_inf_k2 &
  dati$Osservazione <= lim_sup_k2, na.rm = TRUE) / N_validi
13
14 cat(paste0("Disuguaglianza di Chebyshev per k=", k2, ":\n"
  ))
15 cat(paste0("Intervallo (media +/-", k2, "* dev_std): [" ,
  round(lim_inf_k2, 2), ", ", round(lim_sup_k2, 2), "]\n"
  ))
16 cat(paste0("Proporzione minima teorica di dati nell'
  intervallo (almeno): ", round(prop_teorica_k2*100, 2), "
  %\n"))
17 cat(paste0("Proporzione osservata di dati nell'intervallo:
  ", round(prop_osservata_k2*100, 2), "%\n\n"))
18
19 # Per k = 3 deviazioni standard
20 k3 <- 3
21 lim_inf_k3 <- media_generale - k3 * dev_std_campionaria
22 lim_sup_k3 <- media_generale + k3 * dev_std_campionaria
23 prop_teorica_k3 <- 1 - 1/k3^2
24 prop_osservata_k3 <- sum(dati$Osservazione >= lim_inf_k3 &
  dati$Osservazione <= lim_sup_k3, na.rm = TRUE) / N_validi
25
26 cat(paste0("Disuguaglianza di Chebyshev per k=", k3, ":\n"
  ))
27 cat(paste0("Intervallo (media +/-", k3, "* dev_std): [" ,
  round(lim_inf_k3, 2), ", ", round(lim_sup_k3, 2), "]\n"
  ))
```



```

28 cat(paste0("Proporzione minima teorica di dati nell'
           intervallo (almeno):", round(prop_teorica_k3*100, 2), "
           %\n"))
29 cat(paste0("Proporzione osservata di dati nell'intervallo:
           ", round(prop_osservata_k3*100, 2), "%\n"))

```

Output:

Disuguaglianza di Chebyshev per $k = 2$:

Intervallo (media $\pm 2 * dev_std$): [-31.7, 86.72]

Proporzione minima teorica di dati nell'intervallo (almeno): 75%

Proporzione osservata di dati nell'intervallo: 93.59%

Disuguaglianza di Chebyshev per $k = 3$:

Intervallo (media $\pm 3 * dev_std$): [-61.3, 116.33]

Proporzione minima teorica di dati nell'intervallo (almeno): 88.89%

Proporzione osservata di dati nell'intervallo: 100.00%

Per $k = 2$, la disuguaglianza di Chebyshev garantisce che almeno il 75% dei dati si trovi nell'intervallo $[-31.7, 86.72]$. Poiché il numero di osservazioni non può essere negativo, l'intervallo effettivo considerato per i dati è $[1, 86.72]$ (essendo 1 il minimo osservato). La proporzione osservata di dati in questo intervallo è del 93.59%, che soddisfa ampiamente il limite teorico.

Per $k = 3$, almeno l'88.89% dei dati dovrebbe trovarsi nell'intervallo $[-61.3, 116.33]$ (effettivamente $[1, 102]$ per i nostri dati, dato che il massimo è 102). In questo caso, il 100% delle osservazioni valide ricade in questo intervallo, confermando la disuguaglianza.

Capitolo 8

Dati Bivariati

In questa sezione analizzeremo la relazione tra due variabili quantitative: l'Anno di riferimento e il Numero Totale di Morti annuali per incidenti stradali (limitatamente alla fascia di età 21-24 anni, conducenti). Useremo i dati aggregati in `morti_per_anno`.

8.1 Diagramma a dispersione

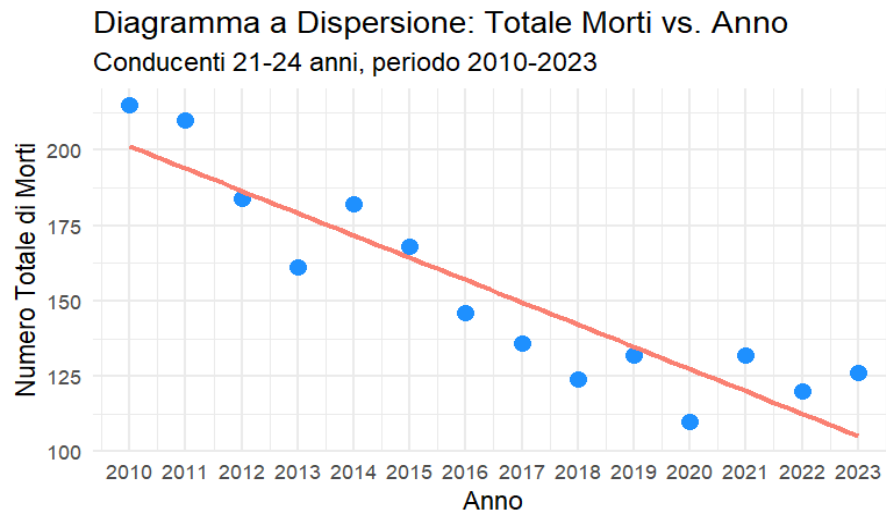
Un diagramma a dispersione (scatter plot) mostra la relazione tra due variabili quantitative.

```
1 # Converti gli anni in numerici per la regressione
2 if (!is.numeric(morti_per_anno$TIME_PERIOD)) {
3   morti_per_anno$TIME_PERIOD_Numeric <- as.numeric(as.
4     character(morti_per_anno$TIME_PERIOD))
5 } else {
6   morti_per_anno$TIME_PERIOD_Numeric <- morti_per_anno$TIME_
7     PERIOD
8 }
9 # Crea lo scatter plot
10 scatter_plot_annuale <- ggplot(morti_per_anno, aes(x = TIME_
11   PERIOD_Numeric, y = Totale_Morti)) +
12   geom_point(color = "dodgerblue", size = 3) +
13   geom_smooth(method = "lm", se = FALSE, color = "salmon") +
14     # Aggiunge retta di regressione (opzionale qui, ma
15     # utile)
16   labs(title = "Diagramma a Dispersione: Totale Morti vs.
17     Anno",
18     subtitle = "Conducenti 21-24 anni, periodo 2010-2023",
19     x = "Anno",
20     y = "Numero Totale di Morti") +
```

```

16     scale_x_continuous(breaks = seq(min(morti_per_anno$TIME_
    PERIOD_Numeric, na.rm=T),
17                                     max(morti_per_anno$TIME_
    PERIOD_Numeric, na.rm=T
    ), by = 1)) + #
    Assicura che gli anni
    siano interi
18     theme_minimal()
19
20     print(scatter_plot_annuale)

```



Il diagramma a dispersione mostra il numero totale di morti per incidenti stradali (nella fascia di età 21-24 anni, conducenti) per ciascun anno dal 2010 al 2023. Ogni punto blu rappresenta un anno. La linea color salmone rappresenta la **retta di regressione** lineare, che suggerisce l'andamento generale dei dati nel tempo. Visivamente, si osserva una tendenza generale alla diminuzione del numero di morti nel corso degli anni considerati, sebbene con alcune fluttuazioni.

8.2 Coefficiente di correlazione campionario

Il **coefficiente di correlazione campionario** (r) è una misura normalizzata della forza e della direzione della relazione lineare tra due variabili quantitative. Varia tra -1 e +1.

- $r = +1$: perfetta correlazione lineare positiva.
- $r = -1$: perfetta correlazione lineare negativa.

- $r \approx 0$: assenza di correlazione lineare (o relazione lineare molto debole).
- Valori vicini a +1 o -1 indicano una forte relazione lineare, mentre valori vicini a 0 indicano una relazione lineare debole o assente.

```

1 # Calcolo del coefficiente di correlazione tra Anno e Totale
  Morti
2 correlazione_annuale <- cor(morti_per_anno$TIME_PERIOD_
  Numeric, morti_per_anno$Totale_Morti, use = "complete.obs",
  method = "pearson")
3 print(paste("Coefficiente di correlazione campionario tra
  Anno e Totale Morti:", round(correlazione_annuale, 4)))

```

Output:

"Coefficiente di correlazione campionario tra Anno e Totale Morti: -0.9132"

Un coefficiente di correlazione campionario di circa -0.9132 indica una **forte relazione lineare negativa** tra l'anno e il numero totale di morti.

Ciò conferma in modo più standardizzato l'osservazione precedente: all'avanzare degli anni, si è osservata una marcata tendenza alla diminuzione del numero di decessi per incidenti stradali tra i conducenti di età compresa tra 21 e 24 anni nel periodo 2010-2023.