



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

<b>1</b>	<b>Introduzione</b>	<b>2</b>
1.1	Problematica . . . . .	2
1.2	Scopo del progetto . . . . .	2
<b>2</b>	<b>Tabelle delle frequenze</b>	<b>4</b>
<b>3</b>	<b>Rappresentazione dei dati mediante grafici</b>	<b>6</b>
<b>4</b>	<b>Indici di posizione</b>	<b>10</b>
4.1	Media campionaria . . . . .	10
4.2	Mediana campionaria . . . . .	10
4.3	Moda campionaria . . . . .	11
<b>5</b>	<b>Indici di variabilità</b>	<b>12</b>
5.1	Varianza campionaria . . . . .	12
5.2	Deviazione standard campionaria . . . . .	12
5.3	Scarto medio assoluto . . . . .	13
5.4	Ampiezza del campo di variazione . . . . .	13
5.5	Coefficiente di variazione . . . . .	14
<b>6</b>	<b>Indici di forma</b>	<b>15</b>
6.1	Indice di asimmetria . . . . .	15
6.2	Indice di curtosi . . . . .	15

# 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, e di gravi conseguenze in termini di disabilità temporanee e permanenti, oltre al drammatico impatto umano e psicologico sulle vittime e sulle loro famiglie.

### 1.2 Scopo del progetto

Il progetto consiste nel realizzare un'indagine statistica <sup>1</sup>sugli 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 e di numerosi Enti pubblici istituzionali, è a carattere totale e a cadenza mensile (inserita tra le rilevazioni di interesse pubblico nel Programma Statistico Nazionale - PSN - IST00142).

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

---

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

formule se necessario.

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](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) 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.

<b>Intersezione</b>	<b>'10</b>	<b>'11</b>	<b>'12</b>	<b>'13</b>	<b>'14</b>	<b>'15</b>	<b>'16</b>	<b>'17</b>	<b>'18</b>	<b>'19</b>	<b>'20</b>	<b>'21</b>	<b>'22</b>	<b>'23</b>
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	–	–
<b>Totale</b>	215	210	184	161	182	168	146	136	124	132	110	132	120	126

## Capitolo 2

# Tabelle delle frequenze

---

```
1 # inclusione librerie
2 library("tidyverse")
3 require("tidyverse")
4 library("dplyr")
5
6 # viene caricato il dataset
7 dati <- read.csv("dati_istat.csv")
8
9 # filtraggio dati
10 dati <- dati %>% select(Intersezione, TIME_PERIOD,
11   Osservazione)
12 dati <- dati %>%
13   filter(Intersezione != "Totale")
14
15 # 1. Frequenze assolute per intersezione dell'incidente
16 freq_assolute <- aggregate(Osservazione ~ Intersezione, data
17   =dati, sum)
18
19 colnames(freq_assolute) <- c("Intersezione", "Frequenza_
20   Assoluta")
21 print(freq_assolute)
22
23 # 2. Frequenze relative per intersezione dell'incidente
24 totale <- sum(freq_assolute$Frequenza_Assoluta)
25 freq_assolute$Frequenza_Relativa <- freq_assolute$Frequenza_
26   Assoluta / totale
27 print(freq_assolute)
28
29 # 3. Frequenze cumulate assolute
30 freq_assolute <- freq_assolute[order(-freq_assolute$
31   Frequenza_Assoluta),]
32 freq_assolute$Frequenza_Cumulata_Assoluta <- cumsum(freq_
33   assolute$Frequenza_Assoluta)
34 print(freq_assolute)
```

```

28
29 # 4. Frequenze cumulate relative
30 freq_assolute$Frequenza_Cumulata_Relativa <- freq_assolute$
    Frequenza_Cumulata_Assoluta / totale
31 print(freq_assolute)

```

---

L'output generato dal codice R riportato nel Listing è sintetizzato nella seguente tabella (si noti che i valori esatti dipendono dalla gestione dei dati mancanti "–" nel file CSV originale; i valori qui riportati sono quelli della tabella fornita nel prompt):

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

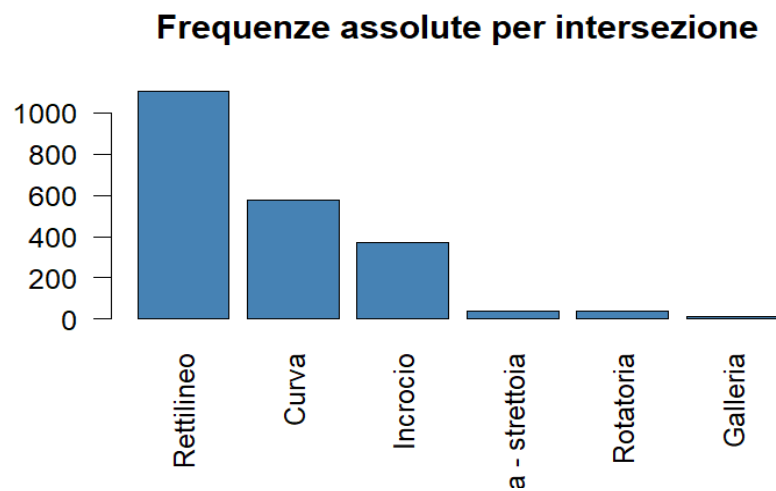
## Capitolo 3

# Rappresentazione dei dati mediante grafici

---

```
1 # 1. Grafico frequenza absolute
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
```

---



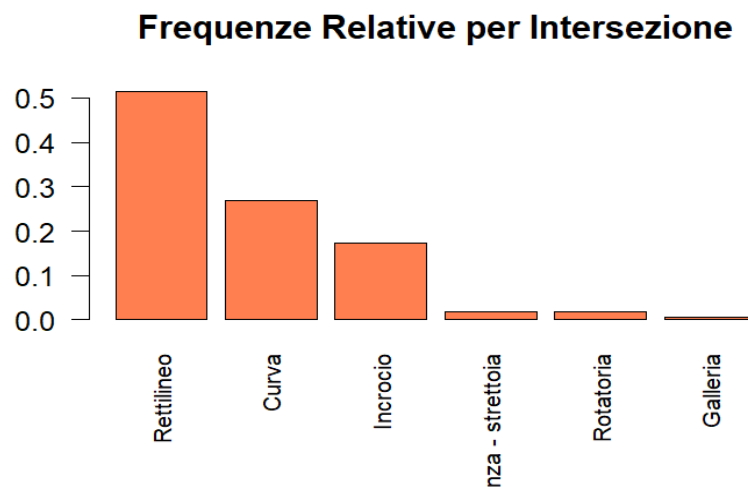
---

```

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)

```

---




---

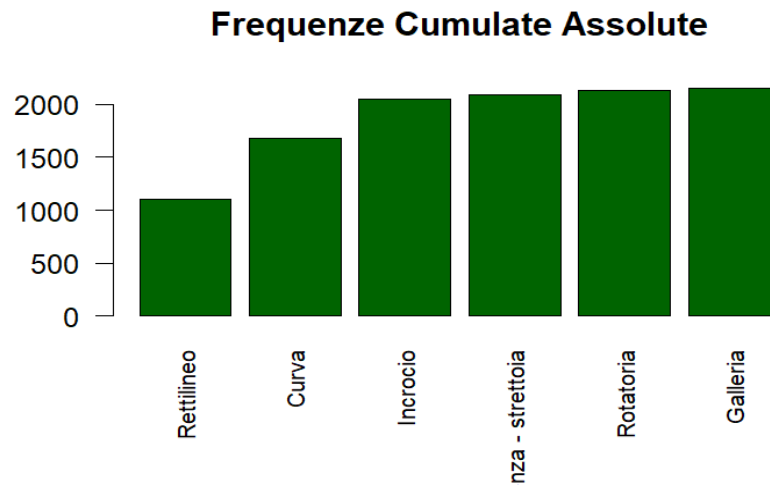
```

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)

```

---





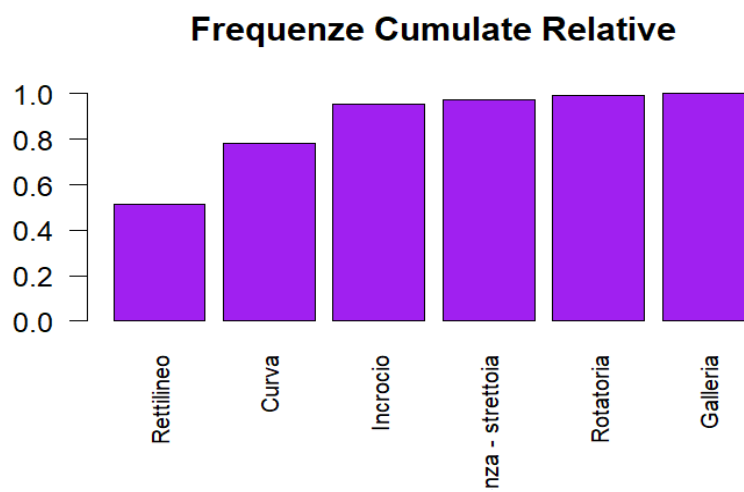

---

```

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)

```

---



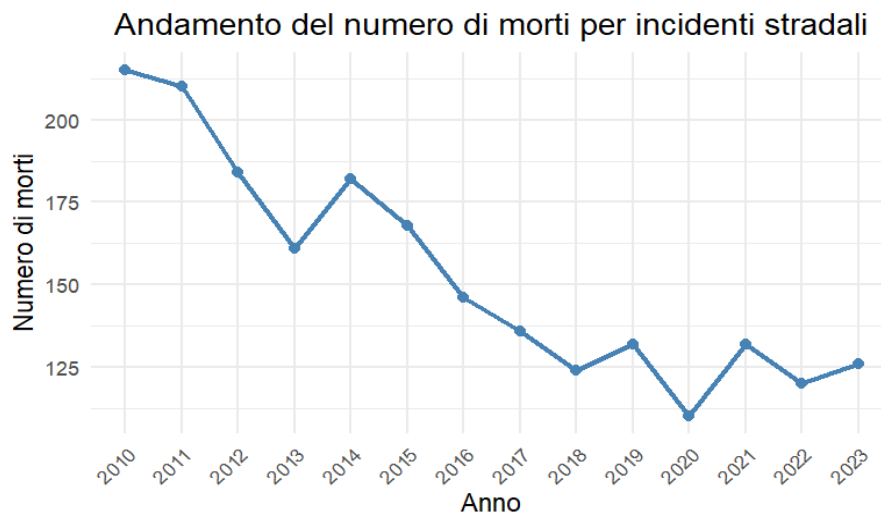
---

```

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_
20                             PERIOD)) # Mostra tutti i valori

```

---



## Capitolo 4

# Indici di posizione

In questo capitolo verranno calcolati gli indici di posizione sulla variabile numerica `Osservazione` del dataset `dati`.

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.

### 4.1 Media campionaria

La media campionaria è la somma di tutte le osservazioni divisa per il numero di osservazioni. 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 il valore centrale di un insieme di dati ordinati. Divide i dati in due metà: il 50% delle osservazioni ha un valore inferiore o uguale alla mediana, e il 50% ha un valore superiore o uguale. È meno sensibile ai valori estremi rispetto alla media.

---

```
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"
```

Metà delle combinazioni intersezione/anno hanno registrato 15.5 morti o meno, e l'altra metà 15.5 morti o più. 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)
```

---

Output:

## Capitolo 5

# Indici di variabilità

Di seguito calcoliamo gli indici di variabilità, 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 (valore basato sulla media e mediana fornite, è una stima):

```
"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 (Mean Absolute Deviation, MAD) dalla media è la media delle deviazioni assolute (cioè, senza segno) dei dati dalla loro media. 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

### 6.1 Indice di asimmetria

### 6.2 Indice di curtosi