

Analisi dei dati

Un goliardico riassunto

Ollari Dmitri

11 maggio 2023

Indice

1	Statistica descrittiva	4
1.1	Grandezze che sintetizzano i dati	4
1.1.1	Media	4
1.1.2	Mediana	4
1.1.3	Moda	4
1.1.4	Varianza	4
1.1.5	Deviazione standard	4
1.1.6	Percentili campionari	5
1.1.7	Box plotting	5
1.2	Disuguaglianza di Chebyshev	5
1.3	Insieme di dati bivariati	5
1.4	Coefficiente di correlazione campionaria	5
2	Elementi di probabilità	7
2.1	Spazio degli esiti ed eventi	7
2.2	Spazio degli esiti equiprobabili	7
2.3	Assiomi della probabilità	7
2.4	Coefficiente binomiale	7
2.5	Probabilità condizionata	7
2.6	Fattorizzazione di un evento e formula di Bayes	8
2.7	Eventi indipendenti	8
3	Variabili aleatorie e valore atteso	9
3.1	Funzione di ripartizione	9
3.2	Variabili aleatorie discrete e continue	10
3.3	Coppie e vettori di variabili aleatorie	10
3.3.1	Distribuzione congiunta - variabili aleatorie discrete	10
3.3.2	Distribuzione congiunta - variabili aleatorie continue	10
3.3.3	Variabili aleatorie indipendenti	10
3.3.4	Distribuzioni condizionali	10
3.4	Valore atteso	11
3.5	Proprietà del valore atteso	11
3.6	Covarianza e varianza della somma di variabili aleatorie	11
4	Funzione generatrice dei momenti	12
5	La legge debole dei grandi numeri	13
6	Modelli di variabili aleatorie	14
6.1	Variabili aleatorie di Bernoulli	14
6.2	Variabili aleatorie binomiale	14
6.2.1	Calcolo esplicito della distribuzione binomiale	14
6.3	Variabile aleatoria di Poisson	15
6.3.1	Calcolo esplicito della distribuzione di Poisson	15
6.4	Variabili aleatorie ipergeometriche	15
6.5	Variabili aleatorie uniformi	16
6.6	Variabili aleatorie normali	16
6.7	Variabili aleatorie esponenziali	17
6.7.1	Processo di Poisson	17
6.8	Variabili aleatorie di tipo gamma	18

6.8.1	Relazione tra chi-quadro e gamma	18
6.9	Le distribuzioni T	18
6.10	Le distribuzioni F	19
6.11	Distribuzione logistica	19

Elenco delle figure

Capitolo 1

Statistica descrittiva

1.1 Grandezze che sintetizzano i dati

1.1.1 Media

Avendo un'insieme di dati x_1, x_2, \dots, x_n , la **media campionaria** è definita da:

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i \quad (1.1)$$

1.1.2 Mediana

Avendo un'insieme di dati **ordinati(crescente)** con ampiezza n , la **mediana** è il valore che **occupa la posizione**:

- $\frac{n+1}{2}$ se n è **dispari**
- Media tra il valore $\frac{n}{2}$ e $\frac{n}{2} + 1$ se n è **pari**

La mediana campionaria (o mediana campionaria ordinata) è un indice di tendenza centrale utilizzato per descrivere il valore centrale di un insieme di dati numerici.

1.1.3 Moda

Unico valore che ha frequenza massima.

1.1.4 Varianza

La varianza è una misura della dispersione o della variabilità dei dati all'interno di un insieme di dati numerici. In altre parole, la varianza indica quanto i dati si distribuiscono intorno alla media.

Dato un'insieme di dati x_1, x_2, \dots, x_n , la **varianza campionaria** (s^2) è:

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 \quad (1.2)$$

1.1.5 Deviazione standard

La deviazione standard campionaria (o deviazione standard campionaria corretta) è una misura della dispersione dei dati all'interno di un campione di dati numerici. Come la varianza, la deviazione standard è un'importante misura di dispersione dei dati, ma a differenza della varianza, ha la stessa unità di misura degli elementi dell'insieme di dati e quindi è più facilmente interpretabile.

$$s = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2} \quad (1.3)$$

1.1.6 Percentili campionari

I percentili campionari sono una misura statistica utilizzata per suddividere un insieme di dati in 100 parti uguali. In altre parole, i percentili campionari suddividono l'insieme di dati in modo tale da poter analizzare la distribuzione dei dati in modo più dettagliato.

Per calcolare i percentili campionari, si segue il seguente procedimento:

1. Si ordina l'insieme di dati in ordine crescente.
2. Si moltiplica il numero totale di elementi nell'insieme di dati per la percentuale desiderata (ad esempio, se si vuole trovare il 25° percentile, si moltiplica il numero totale di elementi per 0,25).
3. Si arrotonda il risultato ottenuto al punto 2 all'intero successivo, se il risultato non è già un intero.
4. Si individua l'elemento nell'insieme di dati che si trova all'indice calcolato al punto 3. Questo valore corrisponde al percentile desiderato.

1.1.7 Box plotting

Il diagramma a scatola si presenta come un rettangolo orientato verticalmente, che ha come estremi inferiori e superiori il primo quartile (Q1) e il terzo quartile (Q3) della distribuzione. All'interno del rettangolo si trova una linea orizzontale che rappresenta la mediana (o secondo quartile, Q2), mentre le linee che si estendono dalla parte superiore e inferiore del rettangolo (i cosiddetti baffi) rappresentano la variazione massima e minima dei dati.

1.2 Disuguaglianza di Chebyshev

La Disuguaglianza di Chebyshev è un importante teorema della teoria della probabilità che fornisce una stima superiore sulla percentuale di dati che si discostano dalla media, per ogni distribuzione di probabilità.

La disuguaglianza afferma che, per qualsiasi distribuzione di probabilità con media μ e varianza finita σ^2 , la percentuale di dati che si discostano dalla media per più di k volte la deviazione standard (cioè la distanza dalla media in termini di deviazioni standard) è inferiore o uguale a $1/k^2$, ovvero:

$$P(|X - \mu| \geq k\sigma) \leq \frac{1}{k^2}$$

dove X è una variabile aleatoria con media μ e varianza σ^2 , e k è un qualsiasi valore positivo.

Questa disuguaglianza implica che, per qualsiasi distribuzione di probabilità, almeno il 75% dei dati si trova entro due deviazioni standard dalla media, e almeno il 89% dei dati si trova entro tre deviazioni standard dalla media. In altre parole, la disuguaglianza di Chebyshev fornisce un limite superiore alla dispersione dei dati, indipendentemente dalla loro distribuzione.

La disuguaglianza di Chebyshev è utile per stimare la probabilità di osservare un valore estremo in un insieme di dati, senza conoscere la loro distribuzione. Ad esempio, se si vuole stimare la percentuale di studenti di una scuola che hanno un QI inferiore a 130, sapendo solo la media e la deviazione standard del QI degli studenti, la disuguaglianza di Chebyshev può essere utilizzata per fornire una stima superiore sulla percentuale di studenti che hanno un QI inferiore a 130.

1.3 Insieme di dati bivariati

In statistica, un insieme di dati bivariati è un insieme di dati che contiene informazioni su due variabili per ogni osservazione nel campione. In altre parole, ogni osservazione nel campione è rappresentata da una coppia ordinata di numeri, uno per ogni variabile. Ad esempio, un insieme di dati bivariati potrebbe contenere informazioni sul reddito e sull'età di un gruppo di persone.

Un insieme di dati bivariati può essere rappresentato graficamente utilizzando un grafico di dispersione, dove i valori delle due variabili sono rappresentati sui due assi cartesiani. Ciò consente di visualizzare la relazione tra le due variabili e identificare eventuali pattern o tendenze.

1.4 Coefficiente di correlazione campionaria

Il coefficiente di correlazione campionaria è una misura di quanto due variabili aleatorie siano linearmente correlate tra loro in un campione di dati. Il coefficiente di correlazione può assumere valori compresi tra -1 e 1, dove un valore di 1 indica una correlazione positiva perfetta, un valore di -1 indica una correlazione negativa perfetta, e un valore di 0 indica assenza di correlazione.

Il coefficiente di correlazione campionaria si calcola come il rapporto tra la covarianza campionaria delle due variabili e il prodotto delle deviazioni standard campionarie delle due variabili. In formule:

$$r = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}}$$

dove n è la dimensione del campione, x_i e y_i sono i valori osservati delle due variabili nel campione, \bar{x} e \bar{y} sono le medie campionarie delle due variabili, e $\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2}$ e $\sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}$ sono le deviazioni standard campionarie delle due variabili.

Un valore di r vicino a 1 indica una forte correlazione positiva tra le due variabili, mentre un valore di r vicino a -1 indica una forte correlazione negativa tra le due variabili. Un valore di r vicino a 0 indica che le due variabili non sono correlate linearmente.

Capitolo 2

Elementi di probabilità

2.1 Spazio degli esiti ed eventi

Lo spazio degli esiti è l'insieme di tutti gli esiti possibili di un esperimento.

2.2 Spazio degli esiti equiprobabili

Lo spazio degli esiti equiprobabili è un concetto fondamentale della probabilità e si riferisce all'insieme di tutti gli esiti possibili di un esperimento casuale in cui ogni esito ha la stessa probabilità di occorrere. In altre parole, è uno spazio campionario in cui ogni esito ha la stessa probabilità di accadere.

La probabilità che accada 1 evento dello spazio degli esiti equiprobabili è:

$$P(E) = \frac{1}{N} \quad (2.1)$$

Dove N è il numero di esiti possibili.

2.3 Assiomi della probabilità

- la probabilità di un evento è compresa tra 0 e 1
- la probabilità dello spazio degli esiti è 1
- unendo due eventi la loro probabilità si somma

2.4 Coefficiente binomiale

Il coefficiente binomiale è il numero di sottoinsiemi di k elementi che si possono formare a partire da un insieme di n elementi.

$$\binom{n}{k} = \frac{n!}{k!(n-k)!} \quad (2.2)$$

2.5 Probabilità condizionata

La probabilità condizionata è una misura di probabilità che tiene conto di una certa informazione o condizione nota. In altre parole, la probabilità di un evento A , dato che un evento B si è verificato, viene calcolata in base alla conoscenza del verificarsi di B .

$$P(A|B) = \frac{P(A \cap B)}{P(B)} \quad (2.3)$$

2.6 Fattorizzazione di un evento e formula di Bayes

La fattorizzazione di un evento è un metodo che consente di scrivere la probabilità di un evento composto come prodotto delle probabilità di eventi più semplici. In altre parole, si tratta di scomporre un evento complesso in eventi più elementari, al fine di semplificare il calcolo della probabilità complessiva.

La formula di Bayes, invece, è un teorema della teoria della probabilità che permette di calcolare la probabilità condizionata di un evento A , data la conoscenza di un evento B . In particolare, la formula di Bayes afferma che:

$$P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)} \quad (2.4)$$

dove $P(A|B)$ rappresenta la probabilità di A dato B , $P(B|A)$ rappresenta la probabilità di B dato A , $P(A)$ rappresenta la probabilità di A , e $P(B)$ rappresenta la probabilità di B .

2.7 Eventi indipendenti

Due eventi A e B si dicono **indipendenti** se la probabilità di A non viene influenzata dalla conoscenza dell'avvenimento B e viceversa, ovvero se:

$$P(A|B) = P(A) \quad (2.5)$$

$$P(B|A) = P(B) \quad (2.6)$$

In altre parole, la conoscenza di uno degli eventi non fornisce alcuna informazione utile per determinare la probabilità dell'altro evento. Ad esempio, se si lancia una moneta equilibrata due volte, il risultato del primo lancio non influisce sulla probabilità di ottenere testa o croce al secondo lancio.

Se due eventi sono indipendenti, allora la probabilità dell'evento composto A e B è data dal prodotto delle probabilità di A e B :

$$P(A \cap B) = P(A) \cdot P(B) \quad (2.7)$$

Capitolo 3

Variabili aleatorie e valore atteso

In teoria delle probabilità, una variabile aleatoria è una funzione che assegna un valore numerico a ciascun possibile esito di un esperimento casuale. In altre parole, una variabile aleatoria è una quantità che può assumere diversi valori a seconda dell'esito dell'esperimento.

Ad esempio, se si lancia un dado equilibrato, il numero che esce è una variabile aleatoria, che può assumere i valori da 1 a 6 con uguale probabilità.

Il valore atteso di una variabile aleatoria è il valore medio che ci si aspetta di ottenere se si ripete l'esperimento un gran numero di volte. In altre parole, è una sorta di media ponderata dei possibili valori che la variabile aleatoria può assumere, in cui i pesi sono le rispettive probabilità.

Il valore atteso di una variabile aleatoria X si indica con $E(X)$ e si calcola come:

$$E(X) = \sum x \times P(X = x) \quad (3.1)$$

dove la somma è estesa a tutti i possibili valori x che X può assumere, e $P(X = x)$ è la probabilità che X assuma il valore x .

Ad esempio, se si considera la variabile aleatoria che indica il numero uscito lanciando un dado equilibrato, il valore atteso è:

$$E(X) = \frac{1}{6} \cdot 1 + \frac{1}{6} \cdot 2 + \frac{1}{6} \cdot 3 + \frac{1}{6} \cdot 4 + \frac{1}{6} \cdot 5 + \frac{1}{6} \cdot 6 = 3.5 \quad (3.2)$$

Questo significa che se si ripete il lancio del dado un gran numero di volte, ci si può aspettare che la media dei risultati si avvicini a 3.5.

Il valore atteso è una misura importante della "centralità" di una variabile aleatoria, e consente di fare previsioni sul comportamento medio della variabile stessa. Ad esempio, se si conosce il valore atteso di una variabile aleatoria X , è possibile stimare la probabilità che X assuma valori superiori o inferiori a una determinata soglia.

3.1 Funzione di ripartizione

La funzione di ripartizione di una variabile aleatoria X si indica con $F(X)$ e si definisce come:

$$F(X) = P(X \leq x)$$

dove x è un valore qualsiasi, e $P(X \leq x)$ rappresenta la probabilità che la variabile aleatoria X assuma un valore minore o uguale a x .

La funzione di ripartizione ha le seguenti proprietà:

1. $F(X)$ è una funzione non-decrescente, ovvero per ogni valore di x_1 e x_2 tali che $x_1 \leq x_2$, si ha $F(x_1) \leq F(x_2)$.
2. $F(X)$ è limitata superiormente da 1 e inferiormente da 0, ovvero $F(-\infty) = 0$ e $F(+\infty) = 1$.
3. La probabilità che la variabile aleatoria X assuma un valore compreso tra due valori x_1 e x_2 si può calcolare come la differenza tra le rispettive probabilità cumulative, ovvero:

$$P(x_1 < X \leq x_2) = F(x_2) - F(x_1)$$

La funzione di ripartizione è una delle proprietà fondamentali di una variabile aleatoria, in quanto consente di calcolare molte altre quantità importanti, come la media e la varianza della variabile stessa. Inoltre, permette di effettuare test di ipotesi e di costruire intervalli di confidenza sulla base dei dati osservati.

3.2 Variabili aleatorie discrete e continue

Una variabile aleatoria si dice discreta se può assumere solo un numero finito o numerabile di valori distinti. Ad esempio, il numero di facce che esce dal lancio di un dado è una variabile aleatoria discreta, in quanto può assumere solo i valori 1, 2, 3, 4, 5, o 6.

Una variabile aleatoria si dice continua se può assumere un numero infinito di valori in un intervallo. Ad esempio, la lunghezza di un filo può essere una variabile aleatoria continua, in quanto può assumere qualunque valore in un intervallo continuo di valori.

- Per le variabili aleatorie discrete, la funzione di probabilità assegna una probabilità a ciascun valore possibile che può assumere la variabile. Questa funzione è spesso rappresentata da un grafico a barre, in cui l'altezza di ogni barra corrisponde alla probabilità associata a un valore particolare.
- Per le variabili aleatorie continue, la funzione di probabilità assume una forma di densità di probabilità, che descrive la probabilità di trovare la variabile in un intervallo di valori. Questa funzione può essere rappresentata da un grafico a linea, in cui l'area sottostante alla curva corrisponde alla probabilità di trovare la variabile in un intervallo specifico.

3.3 Coppie e vettori di variabili aleatorie

La funzione di ripartizione congiunta di X e Y è:

$$F(x, y) := P(X \leq x, Y \leq y) \quad (3.3)$$

Dove la **virgola** denota l'**intersezione** degli eventi.

3.3.1 Distribuzione congiunta - variabili aleatorie discrete

$$p(x_i, y_j) := P(X = x_i, Y = y_j) \quad (3.4)$$

È la **funzione di massa di probabilità congiunta**.

Le funzioni di massa individuali si possono ottenere:

$$p_X(x_i) := P(X = x_i) = \sum_j p(x_i, y_j) \quad (3.5)$$

3.3.2 Distribuzione congiunta - variabili aleatorie continue

$$P(X \in A, Y \in B) = \int_B \int_A f(x, y) dx dy \quad (3.6)$$

È la **densità congiunta**.

Per ricavare le individuali:

$$f_X(x) = \int_{-\infty}^{+\infty} f(x, y) dy \quad (3.7)$$

$$f_Y(y) = \int_{-\infty}^{+\infty} f(x, y) dx \quad (3.8)$$

3.3.3 Variabili aleatorie indipendenti

Due variabili aleatorie sono indipendenti se tutti gli eventi relativi alla prima sono indipendenti dalla seconda e viceversa.

$$P(X \in A, Y \in B) = P(X \in A)P(Y \in B) \quad (3.9)$$

3.3.4 Distribuzioni condizionali

Le distribuzioni condizionali di variabili aleatorie si riferiscono alla distribuzione di una variabile aleatoria data un'informazione o un vincolo su un'altra variabile aleatoria. In altre parole, quando abbiamo informazioni su una variabile aleatoria, possiamo utilizzare tali informazioni per stimare la distribuzione di un'altra variabile aleatoria.

La distribuzione condizionale di X data $Y = y$ è definita come:

$$P(X = x | Y = y) = \frac{P(X = x, Y = y)}{P(Y = y)}$$

3.4 Valore atteso

Il valore atteso (o media) di una variabile aleatoria è una misura della tendenza centrale dei suoi possibili valori. In altre parole, il valore atteso rappresenta il valore medio che ci si aspetta di ottenere da una variabile aleatoria se viene ripetutamente campionata.

Il valore atteso di una variabile aleatoria discreta X è definito come:

$$E[X] = \sum_i x_i P(X = x_i) \quad (3.10)$$

per una variabile aleatoria continua, il valore atteso è definito come:

$$E[X] = \int_{-\infty}^{+\infty} x f(x) dx \quad (3.11)$$

3.5 Proprietà del valore atteso

- Linearità: il valore atteso di una somma di variabili aleatorie è la somma dei loro valori attesi. In altre parole, se X e Y sono variabili aleatorie e a e b sono costanti, allora:

$$E[aX + bY] = aE[X] + bE[Y]$$

- Additività: il valore atteso di una funzione di una variabile aleatoria è la somma dei valori attesi della funzione per ciascun valore della variabile. In altre parole, se $g(X)$ è una funzione di X , allora:

$$E[g(X)] = \sum_x g(x) P(X = x)$$

- Monotonia: se X e Y sono variabili aleatorie tali che $X \leq Y$, allora:

$$E[X] \leq E[Y]$$

- Indipendenza: se X e Y sono variabili aleatorie indipendenti, allora:

$$E[XY] = E[X]E[Y]$$

- Varianza: la varianza di una variabile aleatoria X è definita come:

$$Var[X] = E[(X - E[X])^2]$$

e può essere espressa come:

$$Var[X] = E[X^2] - (E[X])^2$$

3.6 Covarianza e varianza della somma di variabili aleatorie

La covarianza tra due variabili aleatorie X e Y è definita come:

$$Cov[X, Y] = E[(X - E[X])(Y - E[Y])]$$

La varianza della somma di due variabili aleatorie è data da:

$$Var[X + Y] = Var[X] + Var[Y] + 2Cov[X, Y]$$

Se le variabili sono indipendenti, la covarianza è zero e l'equazione si riduce a:

$$Var[X + Y] = Var[X] + Var[Y]$$

Capitolo 4

Funzione generatrice dei momenti

La funzione generatrice dei momenti è utile perché fornisce una rappresentazione compatta delle informazioni sui momenti di una variabile aleatoria. Inoltre, se due variabili aleatorie hanno la stessa funzione generatrice dei momenti, allora hanno gli stessi momenti e quindi la stessa distribuzione di probabilità.

Infine, la funzione generatrice dei momenti è particolarmente utile per variabili aleatorie discrete, poiché consente di calcolare facilmente le probabilità e le statistiche associate a queste variabili.

Nello specifico quando si tratta di **variabili aleatorie discrete**:

$$\Phi(t) = E[e^{tX}] = \sum_x e^{tx} p(x) \quad (4.1)$$

Nello specifico quando si tratta di **variabili aleatorie continue**:

$$\Phi(t) = E[e^{tX}] = \int_{-\infty}^{+\infty} e^{tx} f(x) dx \quad (4.2)$$

Derivando una volta il moto nell'origine si ottiene:

$$\Phi'(0) = E[X] \quad (4.3)$$

In generale si ha che:

$$\Phi^n(0) = E[X^n] \quad (4.4)$$

Se X e Y sono variabili aleatorie **indipendenti** con funzioni generatrici Φ_X e Φ_Y , e se Φ_{X+Y} è la funzione generatrice dei momenti di $X + Y$, allora:

$$\Phi_{X+Y}(t) = \Phi_X(t)\Phi_Y(t) \quad (4.5)$$

Capitolo 5

La legge debole dei grandi numeri

La legge debole dei grandi numeri afferma che, se si considera una sequenza di variabili aleatorie indipendenti ed identicamente distribuite X_1, X_2, \dots, X_n , la media campionaria \bar{X}_n converge in probabilità al valore atteso della variabile aleatoria, ovvero:

$$\lim_{n \rightarrow \infty} P(|\bar{X}_n - \mu| < \epsilon) = 1$$

dove μ è il valore atteso di X_i e ϵ è un valore positivo arbitrario. In altre parole, la probabilità che la media campionaria si discosti dal valore atteso di più di un valore prefissato ϵ diventa sempre più piccola all'aumentare del numero di osservazioni n .

Questa legge fornisce una garanzia formale dell'andamento della media campionaria al crescere del numero di osservazioni e rappresenta un importante risultato in teoria della probabilità e nelle applicazioni pratiche dell'analisi dei dati.

Capitolo 6

Modelli di variabili aleatorie

6.1 Variabili aleatorie di Bernoulli

Una variabile aleatoria di Bernoulli è una variabile aleatoria discreta che assume valore 1 con probabilità p e valore 0 con probabilità $1 - p$, dove $0 \leq p \leq 1$. Ad esempio, la variabile aleatoria che rappresenta il lancio di una moneta equilibrata è di tipo Bernoulli, in quanto assume valore 1 se il risultato del lancio è testa e 0 se è croce.

$$P(X = 0) = 1 - p \quad (6.1)$$

$$P(X = 1) = p \quad (6.2)$$

6.2 Variabili aleatorie binomiale

Una variabile aleatoria binomiale è una variabile aleatoria discreta che rappresenta il numero di successi in una sequenza di n prove indipendenti, ciascuna con probabilità di successo p . La variabile aleatoria binomiale viene indicata con $X \sim B(n, p)$ e assume i valori $0, 1, 2, \dots, n$. La sua funzione di probabilità è data da:

$$P(X = k) = \binom{n}{k} p^k (1 - p)^{n-k}$$

dove $\binom{n}{k}$ è il coefficiente binomiale, che rappresenta il numero di modi diversi in cui è possibile ottenere k successi in n prove. In altre parole, la variabile aleatoria binomiale conta il numero di volte che si ottiene un certo evento in una serie di prove ripetute indipendenti.

Il valore atteso di una variabile aleatoria binomiale è dato da:

$$E[X] = np$$

mentre la sua varianza è:

$$Var(X) = np(1 - p)$$

6.2.1 Calcolo esplicito della distribuzione binomiale

Supponendo X una binomiale di parametri (n, p) , la funzione di ripartizione è:

$$P(X \leq i) = \sum_{k=0}^i \binom{n}{k} p^k (1 - p)^{n-k} \quad (6.3)$$

Per calcolare la funzione di massa:

$$P(X = i) = \binom{n}{i} p^i (1 - p)^{n-i} \quad (6.4)$$

6.3 Variabile aleatoria di Poisson

La variabile aleatoria di Poisson è una variabile aleatoria discreta che descrive il numero di eventi rari che si verificano in un intervallo di tempo o di spazio, dato il tasso di occorrenza di tali eventi. Ad esempio, il numero di chiamate che un centralino telefonico riceve in un minuto, il numero di incidenti stradali in un'ora, il numero di particelle radioattive che decadono in un secondo.

La variabile aleatoria di Poisson è indicata con X e si esprime attraverso un parametro λ che rappresenta il tasso di occorrenza degli eventi. La sua funzione di probabilità è data dalla seguente formula:

$$P(X = k) = \frac{e^{-\lambda} \lambda^k}{k!}$$

Mentre la media e la varianza sono:

$$E(X) = \lambda$$

$$Var(X) = \lambda$$

Una caratteristica interessante della Poissoniana è che può approssimare una binomiale con parametri (n, p) quando n è molto grande e p molto piccolo, ponendo $\lambda = np$:

$$P(X = i) \approx \frac{\lambda^i}{i!} e^{-\lambda} \quad (6.5)$$

6.3.1 Calcolo esplicito della distribuzione di Poisson

Il calcolo esplicito della distribuzione di Poisson si basa sulla seguente formula per la funzione di probabilità:

$$P(X = k) = \frac{e^{-\lambda} \lambda^k}{k!}$$

dove k è il numero di eventi che si verificano in un certo intervallo di tempo o di spazio, λ è il parametro della distribuzione che rappresenta il numero medio di eventi in quell'intervallo.

Per calcolare la probabilità di ottenere un certo numero di eventi k , basta sostituire i valori di k e λ nella formula e fare i calcoli. Ad esempio, supponiamo di voler calcolare la probabilità di osservare esattamente 2 eventi in un intervallo di tempo di un'ora, sapendo che in media accadono 3 eventi in un'ora. In questo caso, la formula diventa:

$$P(X = 2) = \frac{e^{-3} 3^2}{2!} \approx 0.224$$

Quindi la probabilità di osservare esattamente 2 eventi in un'ora, con una media di 3 eventi in un'ora, è di circa il 22,4

6.4 Variabili aleatorie ipergeometriche

Le variabili aleatorie ipergeometriche sono utilizzate per modellare situazioni in cui si effettuano estrazioni senza reinserimento da un insieme finito di elementi che si distinguono in due categorie (ad esempio, in un lotto di componenti elettronici, quelli funzionanti e quelli difettosi).

In particolare, si consideri un insieme di N elementi, dei quali K appartengono alla categoria 1 e $N - K$ alla categoria 2. Si effettuano n estrazioni senza reinserimento e si vuole determinare la probabilità che k estrazioni diano un esito di categoria 1. La variabile aleatoria X che misura il numero di successi (cioè il numero di elementi estratti appartenenti alla categoria 1) è una variabile aleatoria ipergeometrica.

La funzione di probabilità di una variabile aleatoria ipergeometrica è data dalla seguente formula:

$$P(X = k) = \frac{\binom{K}{k} \binom{N-K}{n-k}}{\binom{N}{n}}$$

dove $\binom{a}{b}$ rappresenta il coefficiente binomiale che indica il numero di modi in cui si possono scegliere b elementi da un insieme di a elementi.

Il valore atteso di una variabile aleatoria ipergeometrica è:

$$E(X) = n \frac{K}{N}$$

mentre la varianza è:

$$Var(X) = n \frac{K}{N} \frac{N-K}{N} \frac{N-n}{N-1}$$

Le variabili aleatorie ipergeometriche sono utili, ad esempio, nella valutazione della qualità di un lotto di componenti elettronici, in cui si vuole stimare la proporzione di componenti funzionanti effettuando un campionamento senza reinserimento.

6.5 Variabili aleatorie uniformi

Una variabile aleatoria uniforme è una distribuzione di probabilità in cui ogni valore possibile dell'intervallo di una variabile casuale ha la stessa probabilità di essere scelto. In altre parole, la probabilità di ottenere un valore in un certo intervallo è proporzionale alla lunghezza dell'intervallo.

Esistono due tipi di variabili aleatorie uniformi: la variabile aleatoria uniforme continua e la variabile aleatoria uniforme discreta.

La variabile aleatoria uniforme continua è caratterizzata dalla funzione di densità di probabilità (pdf):

$$f(x) = \begin{cases} \frac{1}{b-a} & \text{per } a \leq x \leq b \\ 0 & \text{altrimenti} \end{cases}$$

dove a e b sono i limiti inferiori e superiori dell'intervallo di valori possibili per la variabile casuale. La funzione di distribuzione cumulativa (cdf) è data da:

$$F(x) = \begin{cases} 0 & \text{per } x < a \\ \frac{x-a}{b-a} & \text{per } a \leq x \leq b \\ 1 & \text{per } x > b \end{cases}$$

La variabile aleatoria uniforme discreta è caratterizzata dalla funzione di massa di probabilità (pmf):

$$P(X = k) = \begin{cases} \frac{1}{n} & \text{per } k = 1, 2, \dots, n \\ 0 & \text{altrimenti} \end{cases}$$

dove n è il numero di valori possibili per la variabile casuale. La funzione di distribuzione cumulativa (cdf) è data da:

$$F(k) = \begin{cases} 0 & \text{per } k < 1 \\ \frac{k-1}{n-1} & \text{per } 1 \leq k \leq n \\ 1 & \text{per } k > n \end{cases}$$

In entrambi i casi, il valore atteso della variabile casuale è dato dalla formula:

$$E[X] = \frac{a+b}{2} \text{ o } \frac{n+1}{2}$$

a seconda che si tratti di una variabile aleatoria uniforme continua o discreta. La varianza della variabile casuale è invece data da:

$$Var(X) = \frac{(b-a)^2}{12} \text{ o } \frac{n^2-1}{12}$$

a seconda che si tratti di una variabile aleatoria uniforme continua o discreta.

6.6 Variabili aleatorie normali

Le variabili aleatorie normali (o gaussiane) sono molto importanti nella teoria della probabilità e nelle applicazioni pratiche, poiché molti fenomeni naturali e sociali seguono una distribuzione normale.

Una variabile aleatoria normale, indicata con X , è caratterizzata da due parametri: la media μ e la deviazione standard σ . La sua funzione di densità di probabilità (pdf) è data dalla seguente formula:

$$f(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

dove e è il numero di Nepero (costante matematica approssimativamente pari a 2,71828), μ è il rapporto tra la circonferenza e il diametro di un cerchio (approssimativamente pari a 3,14159), μ rappresenta il valore atteso della variabile e σ la sua deviazione standard.

La funzione di distribuzione cumulativa (cdf) di X è invece data dalla seguente formula:

$$F(x) = \int_{-\infty}^x f(t)dt = \frac{1}{\sqrt{2\pi\sigma^2}} \int_{-\infty}^x e^{-\frac{(t-\mu)^2}{2\sigma^2}} dt$$

La curva di una distribuzione normale è a forma di campana, simmetrica rispetto alla media μ . Il suo valore atteso, la mediana e la moda coincidono, ovvero:

$$E[X] = \text{mediana}(X) = \text{moda}(X) = \mu$$

Inoltre, il 68% dei dati si trova entro un intervallo di una deviazione standard dalla media, il 95% entro due deviazioni standard e il 99,7% entro tre deviazioni standard. Questa proprietà è nota come la regola empirica o regola del 68-95-99,7.

Le variabili aleatorie normali standardizzate, indicate con Z , sono ottenute dalla trasformazione:

$$Z = \frac{X - \mu}{\sigma}$$

e hanno media zero e deviazione standard pari a uno. La funzione di distribuzione cumulativa di Z è indicata con $\Phi(z)$ e viene chiamata funzione di distribuzione cumulativa standard normale. Esistono diverse tabelle o calcolatori online che permettono di calcolare le probabilità associate a $\Phi(z)$ per diversi valori di z .

6.7 Variabili aleatorie esponenziali

Le variabili aleatorie esponenziali sono un tipo di distribuzione di probabilità continua che descrive il tempo tra due eventi successivi in un processo di Poisson. Ad esempio, se si stanno monitorando gli arrivi di clienti in un negozio e si vuole sapere quanto tempo passerà tra un cliente e il successivo, la distribuzione esponenziale può essere usata per modellare il tempo di attesa.

La funzione densità di probabilità (pdf) della distribuzione esponenziale è:

$$f(x; \lambda) = \begin{cases} \lambda e^{-\lambda x}, & x \geq 0 \\ 0, & x < 0 \end{cases}$$

dove λ è il parametro di scala della distribuzione. Il valore atteso della distribuzione esponenziale è dato da:

$$E[X] = \frac{1}{\lambda}$$

e la deviazione standard è:

$$\text{Var}(X) = \frac{1}{\lambda^2}$$

Inoltre, la distribuzione esponenziale soddisfa la proprietà di mancanza di memoria, che significa che la probabilità che un evento si verifichi dopo un certo intervallo di tempo dipende solo dal tempo trascorso e non dal tempo trascorso finora. Questa proprietà la rende utile per modellare eventi rari e imprevedibili, come guasti a macchinari o incidenti.

6.7.1 Processo di Poisson

Il processo di Poisson è un processo stocastico che descrive il numero di eventi che si verificano in un intervallo di tempo, assumendo che questi eventi si verifichino in modo casuale e indipendente nel tempo.

In forma matematica, il processo di Poisson è definito come:

$$N(t) \sim \text{Pois}(\lambda t)$$

dove $N(t)$ è la variabile aleatoria che rappresenta il numero di eventi che si verificano nell'intervallo di tempo $[0, t]$ e λ è il parametro di intensità del processo, che rappresenta il numero medio di eventi che si verificano in unità di tempo.

La funzione di probabilità di $N(t)$ è data dalla distribuzione di Poisson:

$$P(N(t) = k) = e^{-\lambda t} \frac{(\lambda t)^k}{k!}$$

dove k è il numero di eventi che si verificano nell'intervallo di tempo $[0, t]$.

Il processo di Poisson gode delle proprietà di stazionarietà e indipendenza dei tratti, il che significa che il numero di eventi che si verificano in un intervallo di tempo dipende solo dalla lunghezza dell'intervallo e non dalla posizione dell'intervallo nel tempo, e che gli eventi che si verificano in diversi intervalli di tempo sono indipendenti tra loro.

6.8 Variabili aleatorie di tipo gamma

La variabile aleatoria di tipo gamma è una generalizzazione della distribuzione esponenziale e della distribuzione di Poisson. Una variabile aleatoria X si dice distribuita secondo una legge di tipo gamma con parametri α e λ se la sua funzione di densità di probabilità è data da:

$$f(x) = \begin{cases} \frac{\lambda^\alpha x^{\alpha-1} e^{-\lambda x}}{\Gamma(\alpha)}, & x > 0 \\ 0, & x \leq 0 \end{cases}$$

dove $\Gamma(\alpha)$ è la funzione Gamma, definita come:

$$\Gamma(\alpha) = (n-1)!$$

Il parametro α è un intero positivo e rappresenta il numero di eventi che si verificano in un certo intervallo di tempo, mentre il parametro λ rappresenta la frequenza con cui questi eventi si verificano.

La variabile aleatoria di tipo gamma è particolarmente utile per modellare processi di conteggio con tassi di eventi variabili nel tempo. In particolare, se X_i è la durata del i -esimo intervallo di tempo tra due eventi consecutivi, allora la somma $Y_n = X_1 + X_2 + \dots + X_n$ segue una distribuzione di tipo gamma con parametri $\alpha = n$ e λ uguale alla frequenza media degli eventi.

Il valore atteso di una variabile aleatoria di tipo gamma è pari a $\frac{\alpha}{\lambda}$, mentre la varianza è pari a $\frac{\alpha}{\lambda^2}$.

6.8.1 Relazione tra chi-quadro e gamma

La relazione tra la distribuzione chi-quadrato e la distribuzione gamma è che la distribuzione chi-quadrato è un caso particolare della distribuzione gamma.

In particolare, se una variabile aleatoria Z segue una distribuzione normale standard, allora la somma dei quadrati di k campioni estratti da Z segue una distribuzione chi-quadrato con k gradi di libertà.

Inoltre, se X_1, X_2, \dots, X_n sono variabili aleatorie indipendenti e identicamente distribuite con distribuzione normale standard, allora la somma dei loro quadrati segue una distribuzione chi-quadrato con n gradi di libertà.

La densità di probabilità della distribuzione gamma è data da:

$$f(x) = \frac{\beta^\alpha}{\Gamma(\alpha)} x^{\alpha-1} e^{-\beta x} \quad x \geq 0$$

dove α e β sono i parametri della distribuzione, e $\Gamma(\alpha)$ è la funzione Gamma.

La distribuzione chi-quadrato con k gradi di libertà è una distribuzione gamma con $\alpha = k/2$ e $\beta = 1/2$. In particolare, la densità di probabilità della distribuzione chi-quadrato con k gradi di libertà è:

$$f(x) = \frac{1}{2^{k/2} \Gamma(k/2)} x^{k/2-1} e^{-x/2} \quad x \geq 0$$

Quindi la distribuzione chi-quadrato è un caso particolare della distribuzione gamma, dove i parametri della distribuzione gamma sono determinati dal numero di gradi di libertà della distribuzione chi-quadrato.

6.9 Le distribuzioni T

Le distribuzioni T sono una famiglia di distribuzioni di probabilità continue che hanno una forma simile alla distribuzione normale standard, ma che dipendono anche dal parametro chiamato "gradi di libertà".

In particolare, la distribuzione T con k gradi di libertà è definita come la distribuzione della variabile casuale:

$$T = \frac{Z}{\sqrt{V/k}}$$

dove Z è una variabile casuale standard normale, V è una variabile casuale chi-quadrato con k gradi di libertà e Z e V sono indipendenti.

La distribuzione T è spesso utilizzata quando la deviazione standard della popolazione non è nota e deve essere stimata dalla deviazione standard campionaria. In questo caso, la statistica del test T è definita come:

$$T = \frac{\bar{X} - \mu}{S/\sqrt{n}}$$

dove \bar{X} è la media campionaria, μ è la media della popolazione, S è la deviazione standard campionaria e n è la dimensione del campione. La distribuzione T con $n-1$ gradi di libertà viene utilizzata per calcolare la probabilità di ottenere una statistica T data una determinata ipotesi sulla media della popolazione.

Le distribuzioni T sono utilizzate in diverse applicazioni statistiche, come ad esempio nell'analisi dei test t, nella regressione lineare, nella valutazione delle differenze tra gruppi e in altri tipi di analisi.

Il valore atteso di una distribuzione T non è zero in generale, ma dipende dai gradi di libertà della distribuzione. In particolare, se la distribuzione T ha n gradi di libertà, il valore atteso è zero se $n > 1$, mentre se $n = 1$ il valore atteso non esiste.

La varianza di una distribuzione T con n gradi di libertà è pari a $\frac{n}{n-2}$ se $n > 2$, mentre se $n \leq 2$ la varianza non esiste.

6.10 Le distribuzioni F

La distribuzione F è una distribuzione di probabilità continua che compare comunemente nell'ambito dell'analisi della varianza (ANOVA) e dei test statistici basati sul rapporto di due varianze.

La distribuzione F ha due gradi di libertà, uno per il numeratore e uno per il denominatore. La distribuzione F è quindi definita da due parametri, noti come gradi di libertà del numeratore e del denominatore.

Il valore atteso della distribuzione F dipende dai gradi di libertà del numeratore e del denominatore ed è definito come:

$$E(X) = \begin{cases} \frac{d_2}{d_2 - 2} & \text{se } d_2 > 2 \\ \text{undefined} & \text{se } d_2 \leq 2 \end{cases}$$

dove d_2 rappresenta i gradi di libertà del denominatore.

La varianza della distribuzione F dipende anch'essa dai gradi di libertà del numeratore e del denominatore ed è definita come:

$$Var(X) = \begin{cases} \frac{2d_2^2(d_1 + d_2 - 2)}{d_1(d_2 - 2)^2(d_2 - 4)} & \text{se } d_2 > 4 \\ \text{undefined} & \text{se } d_2 \leq 4 \end{cases}$$

dove d_1 rappresenta i gradi di libertà del numeratore.

La distribuzione F è simmetrica rispetto al valore 1 e assume valori positivi. La sua forma dipende dai gradi di libertà del numeratore e del denominatore e ha una coda lunga sulla destra quando il denominatore è piccolo rispetto al numeratore. La distribuzione F viene solitamente utilizzata per confrontare la varianza di due campioni indipendenti.

6.11 Distribuzione logistica

La distribuzione logistica è una distribuzione di probabilità continua utilizzata spesso in statistica e modellistica per descrivere l'evoluzione di un processo nel tempo.

La funzione di densità di probabilità (PDF) della distribuzione logistica è data da:

$$f(x) = \frac{e^{-(x-\mu)/s}}{s(1 + e^{-(x-\mu)/s})^2}$$

dove μ è il parametro di posizione che indica il valore atteso della distribuzione e s è il parametro di scala. La funzione di distribuzione cumulativa (CDF) può essere ottenuta integrando la PDF:

$$F(x) = \frac{1}{1 + e^{-(x-\mu)/s}}$$

Il parametro di scala s controlla l'inclinazione della curva della distribuzione logistica. Più grande è il valore di s , più stretta sarà la curva e più concentrati saranno i dati attorno al valore atteso μ .

La distribuzione logistica può essere utilizzata per modellare diverse variabili aleatorie, come ad esempio la distribuzione dei punteggi dei test standardizzati o la distribuzione della crescita di popolazioni biologiche.