

Probabilità

24 novembre 2025

Indice

1	Sommatorie e produttorie	11
1.1	Sommatorie	11
1.1.1	Sommatoria singola	11
1.1.1.1	Definizione	11
1.1.1.2	Tecniche utili	12
1.1.1.3	Proprietà	13
1.1.1.4	Applicazioni	15
1.1.2	Sommatorie doppie	17
1.1.2.1	Definizioni	17
1.1.2.2	Proprietà	18
1.2	Produttorie	22
1.2.1	Produttoria singola	22
1.2.1.1	Proprietà	22
1.3	Esercizi	24
2	Calcolo combinatorio	25
2.1	Introduzione	25
2.2	Casistica principale	26
2.2.1	Permutazioni	26
2.2.2	Disposizioni	27
2.2.3	Combinazioni	28
2.2.3.1	Combinazioni semplici	28
2.2.3.2	Combinazioni con ripetizione	28
2.3	Coefficiente binomiale e multinomiale	29
2.3.1	Coefficiente binomiale	29
2.3.1.1	Definizione	29
2.3.1.2	Proprietà	30
2.3.1.3	Origine del nome	31
2.3.2	Il coefficiente multinomiale	31
2.3.2.1	Definizione	31
2.3.2.2	Origine del nome	32
2.4	Calcolo combinatorio e funzioni	32
2.4.1	Principio dell'overcounting	33
2.4.2	Funzioni (disposizioni con ripetizione)	33
2.4.3	Funzioni iniettive (disposizioni semplici)	33
2.4.4	Permutazioni di un insieme (permutazioni semplici)	33
2.4.5	Funzioni caratteristiche (coefficiente binomiale)	33
2.5	Esercizi	34

3	Introduction	37
3.1	Probability space	37
3.1.1	Sample space, events	37
3.1.1.1	Events algebra	38
3.1.1.2	Relationship between events	39
3.1.2	σ -algebra \mathcal{A} (or σ -field \mathcal{F})	40
3.1.3	Probability measure \mathbb{P}	42
3.2	Probability	42
3.2.1	Immediate or useful general results	42
3.2.2	Finite equiprobable Ω and probability evaluation	45
3.2.3	Conditional probability	47
3.2.3.1	Introduction/definition	47
3.2.3.2	Probability of intersection	48
3.2.3.3	Law of total probability	48
3.2.3.4	Bayes formula	50
3.3	Independent events	52
3.3.1	Two events	52
3.3.2	n events	54
3.3.3	Other stuff	55
3.4	Further topics	55
3.4.1	Odds ratio	55
3.4.2	Conditional probability 2	56
3.4.2.1	È una probabilità	56
3.4.2.2	Risultati	57
3.4.2.3	Condizionare su più eventi	59
3.4.2.4	Indipendenza condizionata, aggiornamento delle stime	59
3.4.3	Schema operativo per il calcolo	61
3.5	Esercizi rigo	62
4	Random variables	67
4.1	Intro	67
4.1.1	Random variables linking probability spaces	67
4.1.2	Discrete and continuous rvs	68
4.2	Distribution (and other) functions	69
4.2.1	Types of RVs	70
4.2.2	Discrete rvs	71
4.2.3	Singular continuous rvs	71
4.2.4	Absolutely continuous rvs	72
4.3	OLD: Functions of random variables	74
4.3.1	Discrete rvs: PMF, CDF	74
4.3.2	Continuous rvs: PDF, CDF	75
4.3.3	Other useful rv functions	77
4.3.3.1	Support indicator	77
4.3.3.2	Survival and hazard function	78
4.4	Transformation	79
4.4.1	Discrete rv transform	79
4.4.2	Continuous rvs transform (linear case)	80
4.4.3	Continuous rvs (monotonic) transform	80
4.5	Independence	85

4.5.1	Independence	85
4.5.2	IID RVs	86
4.5.3	Conditional independence	87
4.6	Moments	87
4.6.1	Expected value	87
4.6.2	Variance	92
4.6.3	Asymmetry/skewness and kurtosis	94
4.6.3.1	Asymmetry/Skewness	95
4.6.3.2	Kurtosis	96
4.6.4	Other indicators	96
4.6.4.1	Mediana	96
4.6.4.2	Moda	96
4.7	Random vectors	96
4.7.1	Random vectors and their distribution	96
4.7.2	Type of random vectors	97
4.7.3	Marginals	99
4.7.4	Independence	99
4.8	Relationship between RVs	100
4.8.1	Covariance	100
4.8.2	Correlation coefficient	103
4.9	Exercises	104
4.9.1	Functions and moments	104
4.9.2	Random vectors	105
5	Discrete random variables	107
5.1	Recap notazione	107
5.2	Dirac	109
5.2.1	Defintion	109
5.2.2	Functions	109
5.2.3	Moments	109
5.2.4	Shape	110
5.2.5	Extras	110
5.3	Bernoulli	110
5.3.1	Definition	110
5.3.2	Functions	110
5.3.3	Moments	111
5.3.4	Shape	111
5.3.5	Extras	111
5.3.5.1	Indicator RV	111
5.3.5.2	Applications: probability	112
5.3.5.3	Applications: expected value evaluation	114
5.4	Binomial	115
5.4.1	Definition	115
5.4.2	Functions	115
5.4.3	Moments	116
5.4.4	Shape	117
5.4.5	Extras	119
5.4.5.1	Generazione mediante somma di bernoulliane	119
5.4.5.2	Somma di binomiali	120
5.4.5.3	Variabili derivate	120

5.5	Hypergeometric	121
5.5.1	Definition	121
5.5.2	Functions	121
5.5.3	Moments	122
5.5.4	Extras	123
5.5.4.1	Esperimenti assimilabili	123
5.5.4.2	Connessioni con la binomiale	123
5.6	Geometric (n. of trials)	126
5.6.1	Definition	127
5.6.2	Functions	127
5.6.3	Moments	127
5.6.4	Shape	128
5.6.5	Extras	128
5.6.5.1	Assenza di memoria	128
5.7	Geometric (n. of failures)	128
5.7.1	Functions	129
5.7.2	Moments	130
5.7.3	Extras	131
5.7.3.1	Conversione tra definizioni	131
5.7.3.2	Assenza di memoria	132
5.8	Negative binomial (n. of trials)	132
5.8.1	Definition	132
5.8.2	Functions	133
5.8.3	Moments	133
5.8.4	Shape	133
5.9	Negative binomial (n. of failures)	134
5.9.1	Definition	134
5.9.2	Functions	134
5.9.3	Moments	135
5.9.4	Shape	135
5.9.5	Extras	135
5.9.5.1	Generazione mediante somma di geometriche	135
5.10	Poisson	136
5.10.1	Definition	136
5.10.2	Functions	137
5.10.3	Moments	137
5.10.4	Shape	139
5.10.5	Extras	139
5.10.5.1	Origine e approssimazione	139
5.10.5.2	Legami con la binomiale	141
5.10.5.3	Somma di Poisson indipendenti	144
5.11	Discrete uniform	145
5.11.1	Definition	145
5.11.2	Functions	145
5.11.3	Moments	145
5.12	Exercises	146

6	Absolute continuous random variables	149
6.1	Uniforme continua	149
6.1.1	Definition	149
6.1.2	Functions	149
6.1.3	Moments	150
6.1.4	Shape	151
6.2	Esponenziale	151
6.2.1	Definition	151
6.2.2	Functions	151
6.2.3	Moments	152
6.2.4	Shape	152
6.2.5	Extras	153
6.3	Normale/Gaussiana	154
6.3.1	Definition	154
6.3.2	Functions	154
6.3.3	Shape	155
6.3.4	Normale standardizzata	155
6.3.5	Moments	156
6.3.6	Extras	156
6.4	Gamma	156
6.4.1	Definition	156
6.4.2	Functions	157
6.4.3	Moments	157
6.4.4	Shape	158
6.4.5	Extras	158
6.5	Chi-quadrato	159
6.5.1	Definition	159
6.5.2	Functions	159
6.5.3	Moments	160
6.5.4	Shape	160
6.5.5	Extras	160
6.6	Beta	161
6.6.1	Definition	161
6.6.2	Functions	161
6.6.3	Moments	161
6.6.4	Shape	162
6.6.5	Extras	162
6.7	T di Student	162
6.7.1	Definition	162
6.7.2	Functions	163
6.7.3	Moments	163
6.7.4	Shape	163
6.7.5	Extras	163
6.8	F di Fisher	164
6.8.1	Definition	164
6.8.2	Functions	164
6.8.3	Moments	165
6.8.4	Shape	165
6.8.5	Extras	165
6.9	Logistica	166

6.9.1	Definition	166
6.9.2	Functions	166
6.9.3	Moments	166
6.9.4	Shape	167
6.10	Lognormale	167
6.10.1	Definition	167
6.10.2	Functions	167
6.10.3	Moments	168
6.10.4	Shape	168
6.10.5	Extras	169
6.11	Weibull	169
6.11.1	Definition	169
6.11.2	Functions	169
6.11.3	Moments	169
6.11.4	Shape	169
6.11.5	Extras	170
6.12	Pareto	171
6.12.1	Definition	171
6.12.2	Functions	171
6.12.3	Moments	171
6.12.4	Shape	171
6.13	Exercises	172
7	Misc topics	177
7.1	Quantili	177
7.2	Order statistics	178
7.2.1	Minimum	179
7.2.2	Maximum	180
7.2.3	Generalized $X_{(i)}$	181
7.3	Inequalities	184
7.3.1	Tchebychev (Rigo)	184
7.3.2	Jensen (Rigo)	185
7.3.3	Markov (Viroli)	187
7.3.4	Tchebychev (Viroli)	188
7.4	Characteristic and moment generating function	190
7.4.1	Characteristic function	190
7.4.2	Moment generating function	194
7.5	Conditional distribution	204
7.5.1	Definition and examples	204
7.5.2	Formula to calculate it?	208
7.6	Multivariate normal	211
8	Convergences and related topics	215
8.1	Convergence	215
8.2	Laws of large numbers	220
8.2.1	Strong laws	221
8.2.2	Examples and consequences	222
8.2.3	A weak law	223
8.3	Central limit theorem	226
8.3.1	CLT	226

8.3.2	Examples	229
8.3.3	Berry-Esseen theorem	234
8.4	Additional topics	235
8.4.1	Borel-Cantelli lemma	235
8.4.2	Stable rvs	240
8.4.3	Infinite divisible rvs	242
8.4.4	Examples	244
9	Convergence	247
9.1	Convergence in probability	247
9.1.1	Definition	247
9.1.2	Weak consistence	248
9.1.3	Theorem: weak law of large numbers	251
9.2	Convergence in law/distribution	252
9.2.1	Theorem: central limit theorem	255
9.3	Convergence in mean of order k	257
9.3.1	Definition	257
9.3.2	Strong consistence	257
9.3.3	Theorem: strong law of large numbers	260
9.4	Almost sure convergence	261
9.5	Convergences properties	264
9.6	Delta method	265
10	Simulation	273
10.1	Sampling values from rvs	273
10.1.1	Inversion method	273
10.1.2	Accept-reject method	274
10.2	R exercises	278
10.2.1	CLT	278
10.2.2	Inversion method	278
10.2.3	Accept-reject	280
10.3	Other simulation based stuff	282
10.3.1	Definizioni	282
10.3.2	Metodo di monte carlo	283
10.3.2.1	Calcolo di campioni	283
10.3.2.2	Validazione stimatori	283
10.3.3	Test di permutazione/randomizzazione	285
10.3.3.1	Test di permutazione	285
10.3.3.2	Test di randomizzazione	286
10.3.4	Bootstrap	286

Capitolo 1

Sommatorie e produttorie

1.1 Sommatorie

1.1.1 Sommatoria singola

1.1.1.1 Definizione

Definition 1.1.1 (Sommatoria singola). Se $(a_j)_{j \in J}$, $a : J \rightarrow \mathbb{C}$ è una famiglia *finita* di numeri complessi (ossia l'insieme degli indici J è finito), è definita così la somma di tutti i numeri a_j per $j \in J$ e si indica con

$$\sum_{j \in J} a_j \quad (1.1)$$

Important remark 1. Se $J = \emptyset$ si pone per definizione $\sum_{j \in J} a_j = 0$.

Remark 1. È importante osservare che il simbolo $\sum_{j \in J} a_j$ non dipende da j ma solo dall'intero insieme J e dalla funzione $a : J \rightarrow \mathbb{C}$; la variabile j si dice *muta*, si ha cioè

$$\sum_{j \in J} a_j = \sum_{k \in J} a_k = \sum_{\lambda \in J} a_\lambda$$

Remark 2. Laddove si possa riescere ad esprimere il generico a_j come una $f(j)$ dipendente dall'indice la sommatoria degli elementi può essere usata anche per la somma di valori assunti di funzione che utilizza l'indice come input,

$$\sum_{j \in J} f(j)$$

Example 1.1.1. Se $a_j = 1/j$, allora $\sum_{j \in J} \frac{1}{j}$

Proposition 1.1.1 (Biezione e cambio di indici). *In generale se si ha una funzione $\varphi : K \rightarrow J$ biettiva allora:*

$$\sum_{j \in J} a_j = \sum_{k \in K} a_{\varphi(k)}$$

Remark 3. Ossia possiamo anche utilizzare un altro set di indici K posto che, per garantire l'uguaglianza, vi sia una biezione che ci garantisce che questi vadano a puntare agli stessi elementi.

Remark 4 (Indici comuni). Spesso l'insieme J degli indici è $I_n = \{1, \dots, n\}$ e si scrive allora anche

$$\sum_{j=1}^n a_j \quad \text{oppure} \quad \sum_{1 \leq j \leq n} a_j \quad \text{intendendo} \quad \sum_{j \in I_n} a_j$$

e per esteso si intende:

$$\sum_{j \in I_n} a_j = a_1 + \dots + a_n$$

Definition 1.1.2 (Sommatoria di sottofamiglia). Si intende la sommatoria di un pezzo, ossia di una sottofamiglia di successione $a : \mathbb{N} \rightarrow \mathbb{C}$ compresa tra due indici m, n , con $m \leq n$:

$$\sum_{j=m}^n a_j = \sum_{m \leq j \leq n} a_j = a_m + \dots + a_n$$

1.1.1.2 Tecniche utili

Remark 5. La traslazione di indici consiste nel cambiare gli indici senza cambiare gli oggetti puntati.

Proposition 1.1.2 (Traslazione di indici). *Per effettuarla occorre sostituire $j + \text{offset}$ al posto di j negli indici della sommatoria e sostituendo $j - \text{offset}$ nei termini indicati (sia offset un termine positivo o negativo)*

$$\sum_{j=m}^n a_j = \sum_{j=m-p}^{n-p} a_{j+p} = \sum_{j=m+p}^{n+p} a_{j-p} \quad (1.2)$$

Dimostrazione. È una applicazione di 1.1.1. □

Remark 6. In sostanza per garantire l'uguaglianza delle sommatorie, basta che alla fine l'indice punti allo stesso elemento poi la formula può essere cambiata a piacere.

Proposition 1.1.3 (Riflessione di indici). *Mediante questa tecnica si sommano gli stessi elementi, posti però in ordine inverso (si somma dall'indice originario più alto al più basso):*

$$\sum_{i=1}^n a_i = \sum_{i=1}^n a_{n-i+1} = \sum_{i=0}^{n-1} a_{n-i} \quad (1.3)$$

Dimostrazione. È una permutazione su indici quindi funzione biettiva e si applica 1.1.1. L'ultima eguaglianza si giustifica mediante una traslazione di indici (sostituendo con $i - 1$ negli indici della sommatoria e $i + 1$ nei termini della stessa). □

Remark 7. Il cambiamento di indice può tornare utile nel caso di sommatoria di funzione laddove si vogliano normalizzare un po' gli indici

Proposition 1.1.4 (Cambiamento di indice). *Sia $\sum_{i \in I} f(i)$ la sommatoria di nostro interesse. Supponendo che vi sia una funzione biettiva $\varphi : I \rightarrow J$ che esprima gli indici in un nuovo insieme e che sia $J = \varphi(I)$; esisterà anche $\varphi^{-1} : J \rightarrow I$. Dato che un singolo $j = \varphi(i)$ si potranno applicare φ e φ^{-1} rispettivamente a indici ed elementi della sommatoria, ottenendo lo stesso risultato poiché, per definizione*

$$\sum_{i \in I} f(i) = \sum_{j = \varphi(i) \in J} f(\varphi^{-1}(j)) \quad (1.4)$$

Remark 8. L'equazione di sopra ci dice che dobbiamo applicare la biezione trovata agli indici e la sua inversa all'argomento della sommatoria

Example 1.1.2. Ipotizziamo di avere

$$\sum_{i=-10}^{-8} \frac{1}{i+1} = -\frac{1}{9} - \frac{1}{8} - \frac{1}{7}$$

Al fine di semplificare gli indici della sommatoria applichiamo a questi $\varphi : I \rightarrow J$ definita come $j = i + 10$ (si vede che φ è biettiva: è una retta), si applica poi $\varphi^{-1} : J \rightarrow I$ definita come $i = j - 10$ agli elementi della sommatoria

$$\sum_{i=-10}^{-8} \frac{1}{i+1} = \sum_{j=0}^2 \frac{1}{j-9} = -\frac{1}{9} - \frac{1}{8} - \frac{1}{7}$$

Se si desidera, possiamo tornare all'indice iniziale con la sostituzione $i = j$:

$$\sum_{j=0}^2 \frac{1}{j-9} = \sum_{i=0}^2 \frac{1}{i-9}$$

Quindi:

$$\sum_{i=-10}^{-8} \frac{1}{i+1} = \sum_{i=0}^2 \frac{1}{i-9}$$

1.1.1.3 Proprietà

Remark 9. Valgono le seguenti *proprietà* (che possono essere utili lette sia da sinistra a destra che viceversa).

Proposition 1.1.5 (Sommatoria di costante). *Se k è una costante che non dipende dall'indice i , allora:*

$$\sum_{i=1}^n k = nk \quad (1.5)$$

Dimostrazione. Si pone convenzionalmente $a_i = k$, per cui:

$$\sum_{i=1}^n a_i = a_1 + a_2 + \dots + a_n = \underbrace{k + k + \dots + k}_{n \text{ volte}} = kn$$

□

Proposition 1.1.6 (Sommatoria di prodotto per costante). *Se k è una costante che non dipende dall'indice i , allora:*

$$\sum_{i=1}^n k a_i = k \sum_{i=1}^n a_i \quad (1.6)$$

Dimostrazione. Infatti

$$\sum_{i=1}^n k a_i = k a_1 + k a_2 + \dots + k a_n = k(a_1 + a_2 + \dots + a_n) = k \sum_{i=1}^n a_i$$

□

Proposition 1.1.7 (Scomposizione/somme su sottoinsiemi). *Se $m > n$, allora:*

$$\sum_{i=1}^n a_i + \sum_{i=n+1}^m a_i = \sum_{i=1}^m a_i \quad (1.7)$$

Dimostrazione. Infatti

$$\sum_{i=1}^n a_i + \sum_{i=n+1}^m a_i = (a_1 + \dots + a_n) + (a_{n+1} + \dots + a_m) = \sum_{i=1}^m a_i$$

□

Important remark 2. Generalizzando, se Λ è un insieme di indici, $a : \Lambda \rightarrow \mathbb{C}$ una famiglia di complessi, e J, K sottoinsiemi finiti *disgiunti* di Λ si ha:

$$\sum_{\lambda \in J \cup K} a_\lambda = \sum_{\lambda \in J} a_\lambda + \sum_{\lambda \in K} a_\lambda \quad (1.8)$$

Proposition 1.1.8 (Sommatoria di somme/additività rispetto alle famiglie). *Si ha che:*

$$\sum_{i=1}^n (a_i + b_i) = \sum_{i=1}^n a_i + \sum_{i=1}^n b_i \quad (1.9)$$

Dimostrazione. Infatti:

$$\begin{aligned} \sum_{i=1}^n (a_i + b_i) &= (a_1 + b_1) + (a_2 + b_2) + \dots + (a_n + b_n) \\ &= (a_1 + a_2 + \dots + a_n) + (b_1 + b_2 + \dots + b_n) \\ &= \sum_{i=1}^n a_i + \sum_{i=1}^n b_i \end{aligned}$$

□

Important remark 3. Generalizzando, se Λ è un insieme di indici, $a : \Lambda \rightarrow \mathbb{C}$ una famiglia di complessi e $b : \Lambda \rightarrow \mathbb{C}$ è un'altra famiglia di complessi si può definire la somma puntuale $a+b : \Lambda \rightarrow \mathbb{C}$ delle due famiglie ponendo $(a+b)(\lambda) = a_\lambda + b_\lambda$ per ogni $\lambda \in \Lambda$. Si ha anche che per ogni sottoinsieme finito J di Λ :

$$\sum_{j \in J} (a_j + b_j) = \sum_{j \in J} a_j + \sum_{j \in J} b_j \quad (1.10)$$

Proposition 1.1.9 (Sommatoria di termini lineari). *Se k e c sono costanti che non dipendono dall'indice i ,*

$$\sum_{i=1}^n (ka_i + c) = nc + k \sum_{i=1}^n a_i \quad (1.11)$$

Dimostrazione. Alla luce delle proprietà precedentemente viste:

$$\sum_{i=1}^n (ka_i + c) = \sum_{i=1}^n ka_i + \sum_{i=1}^n c = nc + k \sum_{i=1}^n a_i$$

□

Remark 10. Si noti che prima abbiamo preposto nc alla sommatoria per evitare confusione; un altro modo sarebbe $k(\sum_{i=1}^n a_i) + nc$

1.1.1.4 Applicazioni

Proposition 1.1.10 (Prodotti di sommatorie). *Se $(a_j)_{j \in J}$ è una famiglia finita di numeri complessi e $(b_k)_{k \in K}$ è un'altra famiglia finita di numeri complessi si ha:*

$$\left(\sum_{j \in J} a_j \right) \cdot \left(\sum_{k \in K} b_k \right) = \sum_{j \in J, k \in K} a_j b_k = \sum_{(j,k) \in J \times K} a_j b_k \quad (1.12)$$

Dimostrazione. Accettiamo il fatto (che si può dimostrare per induzione sul numero di elementi di K) e verificare nei casi più semplici, es $(a+b)(c+d) = ac + ad + bc + bd$. □

Prodotti di sommatorie aventi medesimo insieme di indici Nel caso gli elementi siano indicati dal medesimo set, es $J = I_n$, possiamo iniziare a pensare il relativo prodotto cartesiano $J \times J$ della precedente come una matrice quadrata:

$$\begin{aligned} \left(\sum_{i=1}^n a_i \right) \left(\sum_{i=1}^n b_i \right) &= (a_1 + a_2 + \dots + a_n)(b_1 + b_2 + \dots + b_n) \\ &= a_1 b_1 + a_1 b_2 + \dots + a_1 b_n + \\ &\quad a_2 b_1 + a_2 b_2 + \dots + a_2 b_n + \\ &\quad \dots + \\ &\quad a_n b_1 + a_n b_2 + \dots + a_n b_n \\ &= \sum_{i=1}^n a_i b_i + \sum_{i \neq j} a_i b_j \end{aligned}$$

In altre parole abbiamo scomposto la sommatoria in due pezzi; quella degli elementi residenti sulla diagonale principale (primo termine) e i rimanenti (secondo termine).

Quadrato di sommatoria Nel caso particolare di quadrato di sommatoria degli elementi $(a_j)_{j \in J}$, si ha:

$$\left(\sum_{j \in J} a_j \right)^2 = \left(\sum_{j \in J} a_j \right) \cdot \left(\sum_{j \in J} a_k \right) = \sum_{(j,k) \in J \times J} a_j a_k \quad (1.13)$$

Per ritrovare l'usuale espressione del quadrato di una somma spezziamo indici $J \times J$ (e relative sommatorie) nella diagonale $\Delta = \{(j, j) : j \in J\}$ e nel suo complementare $J \times J \setminus \Delta$. Si ha

$$\sum_{(j,k) \in J \times J} a_j a_k = \sum_{(j,k) \in \Delta} a_j a_k + \sum_{(j,k) \in J \times J \setminus \Delta} a_j a_k$$

e dato che essendo $j = k$ per $(j, k) \in \Delta$ possiamo riscrivere il primo termine come

$$\sum_{(j,k) \in \Delta} a_j a_k = \sum_{j \in J} a_j a_j = \sum_{j \in J} a_j^2$$

mentre il secondo termine, che dipende dal set di indici $J \times J \setminus \Delta$, può essere diviso in due parti disgiunte, $S = \{(j, k) : j < k\}$ e $T = \{(j, k) : j > k\}$ (si pensi al triangolo superiore e inferiore della matrice che rappresenta il prodotto cartesiano $J \times J$):

$$\sum_{(j,k) \in J \times J \setminus \Delta} a_j a_k = \sum_{(j,k) \in S} a_j a_k + \sum_{(j,k) \in T} a_j a_k$$

e poiché $(j, k) \rightarrow (k, j)$ è una biezione dell'insieme S su T si ha

$$\sum_{(j,k) \in T} a_j a_k = \sum_{(k,j) \in S} a_j a_k = \sum_{(r,s) \in S} a_s a_r$$

dove nell'ultimo passaggio abbiamo effettuato un mero cambio di indici muti. Se ne effettuiamo uno lievemente simile nell'altro termine

$$\sum_{(j,k) \in S} a_j a_k = \sum_{(r,s) \in S} a_r a_s$$

possiamo tornare a

$$\begin{aligned} \sum_{(j,k) \in S} a_j a_k + \sum_{(j,k) \in T} a_j a_k &= \sum_{(r,s) \in S} a_r a_s + \sum_{(r,s) \in S} a_s a_r \\ &= \sum_{(r,s) \in S} (a_r a_s + a_s a_r) \\ &= \sum_{(r,s) \in S} 2a_r a_s \end{aligned}$$

e si conclude con infine a

$$\left(\sum_{j \in J} a_j \right)^2 = \sum_{j \in J} a_j^2 + \sum_{(j,k) \in J \times J : j < k} 2a_j a_k$$

cioè il quadrato di una somma è la somma dei quadrati di tutti i termini, più la somma di tutti i doppi prodotti dei termini stessi.

1.1.2 Sommatorie doppie

1.1.2.1 Definizioni

Remark 11. Date più quantità dipendenti da due indici, es:

$$\begin{array}{cccc} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \dots & \dots & \dots & \dots \\ a_{m1} & a_{m2} & \dots & a_{mn} \end{array}$$

la loro somma si può scrivere utilizzando la notazione di sommatoria:

$$\begin{aligned} &= (a_{11} + a_{12} + \dots a_{1n}) + \\ &\quad + (a_{21} + a_{22} + \dots a_{2n}) + \\ &\quad + \dots + \\ &\quad + (a_{m1} + a_{m2} + \dots a_{mn}) \\ &= \sum_{i=1}^n a_{1i} + \sum_{i=1}^n a_{2i} + \dots + \sum_{i=1}^n a_{mi} \end{aligned}$$

Ponendo $\sum_{i=1}^n a_{1i} = S_1, \sum_{i=1}^n a_{2i} = S_2, \dots, \sum_{i=1}^n a_{mi} = S_m$, la somma degli $m \times n$ elementi a diviene

$$S_1 + S_2 + \dots + S_m = \sum_{j=1}^m S_j = \sum_{j=1}^m \sum_{i=1}^n a_{ji}$$

e si legge “sommatoria doppia delle a_{ji} con j che varia da 1 a m ed i che varia da 1 a n ”, essendo a_{ji} il termine generico che compare nella somma.

Remark 12. Si noti il caso particolare $\sum_{j=c}^c \sum_{i=k}^k a_{ji} = a_{ck}$.

Remark 13. Le due sommatorie si possono invertirsi (con l'effetto che prima di sommare una riga e poi passare alla successiva, prima si somma una colonna per passare poi alla susseguente; il quale ovviamente non ha riverbero sui risultati)

Proposition 1.1.11 (Inversione delle sommatorie).

$$\sum_{j=1}^m \sum_{i=1}^n a_{ji} = \sum_{i=1}^n \sum_{j=1}^m a_{ji}$$

Dimostrazione. Si ha

$$\begin{aligned} \sum_{j=1}^m \sum_{i=1}^n a_{ji} &= (a_{11} + a_{12} + \dots a_{1n}) + \\ &\quad + (a_{21} + a_{22} + \dots a_{2n}) + \\ &\quad + \dots + \\ &\quad + (a_{m1} + a_{m2} + \dots a_{mn}) \\ &= (a_{11} + a_{21} + \dots + a_{m1}) + \\ &\quad + (a_{12} + a_{22} + \dots + a_{m2}) + \\ &\quad + \dots + \\ &\quad + (a_{1n} + a_{2n} + \dots + a_{mn}) = \\ &= \sum_{j=1}^m a_{j1} + \sum_{j=1}^m a_{j2} + \dots + \sum_{j=1}^m a_{jn} \end{aligned}$$

Ponendo $\sum_{j=1}^m a_{j1} = Z_1, \sum_{j=1}^m a_{j2} = Z_2, \dots, \sum_{j=1}^m a_{jn} = Z_n$ si ha

$$\sum_{j=1}^m \sum_{i=1}^n a_{ji} = Z_1 + Z_2 + \dots + Z_n = \sum_{i=1}^n Z_i = \sum_{i=1}^n \sum_{j=1}^m a_{ji}$$

□

Remark 14. Anche in questo caso le lettere j e i , indici del termine generico, possono essere sostituite da qualsiasi altre lettere. Talvolta si può trovare $\sum_{j=1}^m \sum_{i=1}^n a_{ji}$ espresso omettendo gli estremi del campo di variazione della i e della j (se ciò non crea confusione o equivoci), mediante $\sum_j \sum_i a_{ji}$ o anche $\sum_{j,i} a_{ji}$. Talvolta si può trovare la scrittura $\sum \sum a_{ji}$ che è bene evitare perché è sempre meglio indicare gli indici variabili (nel nostro caso j e i) rispetto ai quali si esegue la somma.

1.1.2.2 Proprietà

Proposition 1.1.12 (Sommatoria di costante). *Se k è una costante che non dipende dagli indici j e i :*

$$\sum_{j=1}^m \sum_{i=1}^n k = kmn \quad (1.14)$$

Dimostrazione. Infatti è una sommatoria doppia in cui il termine generico $a_{ji} = k$:

$$\sum_{j=1}^m \sum_{i=1}^n k = \sum_{j=1}^m kn = n \sum_{j=1}^m k = kmn$$

□

Proposition 1.1.13 (Sommatoria di prodotto per costante). *Se k è una costante che non dipende dagli indici j e i :*

$$\sum_{j=1}^m \sum_{i=1}^n ka_{ji} = k \sum_{j=1}^m \sum_{i=1}^n a_{ji} \quad (1.15)$$

Dimostrazione. Infatti

$$\begin{aligned}
 \sum_{j=1}^m \sum_{i=1}^n ka_{ji} &= ka_{11} + ka_{12} + \dots + ka_{1n} + \\
 &\quad + ka_{21} + ka_{22} + \dots + ka_{2n} \\
 &\quad + \dots + \\
 &\quad + ka_{m1} + ka_{m2} + \dots + ka_{mn} \\
 &= \sum_{i=1}^n ka_{1i} + \sum_{i=1}^n ka_{2i} + \dots + \sum_{i=1}^n ka_{mi} = \\
 &= k \sum_{i=1}^n a_{1i} + k \sum_{i=1}^n a_{2i} + \dots + k \sum_{i=1}^n a_{mi} = \\
 &= k \left(\sum_{i=1}^n a_{1i} + \sum_{i=1}^n a_{2i} + \dots + \sum_{i=1}^n a_{mi} \right) = \\
 &= k \sum_{j=1}^m \sum_{i=1}^n a_{ji}
 \end{aligned}$$

□

Proposition 1.1.14 (Scomposizione/somme su sottoinsiemi). *Si ha che:*

$$\sum_{j=1}^m \sum_{i=1}^{n_1} a_{ji} + \sum_{j=1}^m \sum_{i=n_1+1}^n a_{ji} = \sum_{j=1}^m \sum_{i=1}^n a_{ji} \quad (1.16)$$

Dimostrazione. Infatti

$$\begin{aligned}
 \sum_{j=1}^m \sum_{i=1}^{n_1} a_{ji} + \sum_{j=1}^m \sum_{i=n_1+1}^n a_{ji} &= \sum_{j=1}^m \left(\sum_{i=1}^{n_1} a_{ji} + \sum_{i=n_1+1}^n a_{ji} \right) = \\
 &= \sum_{j=1}^m \sum_{i=1}^n a_{ji} \\
 \sum_{j=1}^{m_1} \sum_{i=1}^n a_{ji} + \sum_{j=m_1+1}^m \sum_{i=1}^n a_{ji} &= \sum_{i=1}^n \sum_{j=1}^{m_1} a_{ji} + \sum_{i=1}^n \sum_{j=m_1+1}^m a_{ji} = \\
 &= \sum_{j=1}^m \sum_{i=1}^n a_{ji}
 \end{aligned}$$

□

Remark 15. Per visualizzare le operazioni di cui sopra si pensi ad una somma degli elementi di una matrice che procede attraverso le colonne (sommatoria interna) e poi passa alla prossima riga (ciclo sulla sommatoria esterna); nel primo caso qui sopra abbiamo aggiunto delle colonne ad una matrice, mentre nel secondo abbiamo aggiunto delle righe ad un'altra matrice.

Proposition 1.1.15 (Sommatoria di somme). *Vale la:*

$$\sum_{j=1}^m \sum_{i=1}^n (a_{ji} + b_{ji}) = \sum_{j=1}^m \sum_{i=1}^n a_{ji} + \sum_{j=1}^m \sum_{i=1}^n b_{ji} \quad (1.17)$$

Dimostrazione. Infatti

$$\begin{aligned}
 \sum_{j=1}^m \sum_{i=1}^n (a_{ji} + b_{ji}) &= \sum_{j=1}^m \left[\sum_{i=1}^n (a_{ji} + b_{ji}) \right] \\
 &= \sum_{j=1}^m \left[\sum_{i=1}^n a_{ji} + \sum_{i=1}^n b_{ji} \right] \\
 &= \sum_{j=1}^m \sum_{i=1}^n a_{ji} + \sum_{j=1}^m \sum_{i=1}^n b_{ji}
 \end{aligned}$$

□

Proposition 1.1.16 (Sommatoria di termini lineari). *Se k e c sono costanti che non dipendono dagli indici j e i , vale:*

$$\sum_{j=1}^m \sum_{i=1}^n (ka_{ji} + c) = mnc + k \sum_{j=1}^m \sum_{i=1}^n a_{ji} \quad (1.18)$$

Dimostrazione. Infatti:

$$\begin{aligned}
 \sum_{j=1}^m \sum_{i=1}^n (ka_{ji} + c) &= \sum_{j=1}^m \sum_{i=1}^n ka_{ji} + \sum_{j=1}^m \sum_{i=1}^n c \\
 &= k \sum_{j=1}^m \sum_{i=1}^n a_{ji} + c \sum_{j=1}^m \sum_{i=1}^n 1 \\
 &= cmn + k \sum_{j=1}^m \sum_{i=1}^n a_{ji}
 \end{aligned}$$

□

Proposition 1.1.17 (Portar fuori sommatoria). *È lecito estrarre da ogni sommatoria i termini che non dipendono dall'indice della sommatoria:*

$$\sum_{j=1}^m \sum_{i=1}^n a_j b_i = \sum_{j=1}^m a_j \sum_{i=1}^n b_i \quad (1.19)$$

Cioè dalla seconda sommatoria, fatta secondo l'indice i , si può estrarre il termine a_j che da i non dipende.

Dimostrazione. Infatti

$$\begin{aligned}
 \sum_{j=1}^m \sum_{i=1}^n a_j b_i &= a_1 b_1 + a_1 b_2 + \dots + a_1 b_n + \\
 &\quad a_2 b_1 + a_2 b_2 + \dots + a_2 b_n + \\
 &\quad \vdots \\
 &\quad a_n b_1 + a_n b_2 + \dots + a_n b_n = \\
 &= a_1 \sum_{i=1}^n b_i + a_2 \sum_{i=1}^n b_i + \dots + a_m \sum_{i=1}^n b_i \\
 &= (a_1 + a_2 + \dots + a_m) \sum_{i=1}^n b_i \\
 &= \sum_{j=1}^m a_j \sum_{i=1}^n b_i
 \end{aligned}$$

□

Lemma 1.1.18. *Da ciò deriva ad esempio che si può scrivere*

$$\left(\sum_{i=1}^n a_i \right)^2 = \sum_{j=1}^n \sum_{i=1}^n a_i a_j$$

Dimostrazione. Infatti

$$\left(\sum_{i=1}^n a_i \right)^2 = \sum_{i=1}^n a_i \cdot \sum_{i=1}^n a_i = \sum_{j=1}^n a_j \cdot \sum_{i=1}^n a_i = \sum_{j=1}^n \sum_{i=1}^n a_j a_i$$

dove abbiamo posto j al posto di i in una delle due sommatorie per evitare confusioni. □

Lemma 1.1.19. *È lecito anche scrivere:*

$$\begin{aligned}
 \sum_{j=1}^m \sum_{i=1}^n a_j &= \sum_{j=1}^m a_j \sum_{i=1}^n 1 = n \sum_{j=1}^m a_j \\
 \sum_{j=1}^m \sum_{i=1}^n b_i &= \sum_{i=1}^n b_i \sum_{j=1}^m 1 = m \sum_{i=1}^n b_i
 \end{aligned}$$

Lemma 1.1.20. *È corretto effettuare la seguente posizione:*

$$\sum_{j=1}^m \sum_{i=1}^n a_j b_{ji} = \sum_{j=1}^m a_j \sum_{i=1}^n b_{ji}$$

cioè estrarre a_j dalla seconda sommatoria da cui non dipende, perché quest'ultima è fatta rispetto all'indice i .

Important remark 4. Si osservi che è scorretto scrivere

$$\sum_{i=1}^n b_{ji} \sum_{j=1}^m a_j$$

cioè non è possibile estrarre b_{ji} da alcuna sommatoria perché dipende da entrambi gli indici e quindi da entrambe le sommatorie.

1.2 Produttorie

1.2.1 Produttoria singola

Definition 1.2.1 (Produttoria). Se $(a_j)_{j \in J}$, $a : J \rightarrow \mathbb{C}$ è una famiglia *finita*, il prodotto di tutti i numeri a_j per $j \in J$ si indica con:

$$\prod_{j \in J} a_j \quad (1.20)$$

Remark 16. Si pone per convenzione

$$\prod_{j \in \emptyset} a_j = 1 \quad (1.21)$$

1.2.1.1 Proprietà

Remark 17. Analogamente al caso delle sommatorie valgono le seguenti *proprietà* (che possono essere utili lette sia da sinistra a destra che viceversa).

Proposition 1.2.1 (Produttoria di costante). *Se k è una costante che non dipende dall'indice i :*

$$\prod_{i=1}^n k = k^n \quad (1.22)$$

Dimostrazione. Infatti è una produttoria in cui il termine generico $a_i = k$

$$\prod_{i=1}^n a_i = a_1 a_2 \dots a_n = k \cdot k \cdot \dots \cdot k = k^n$$

□

Proposition 1.2.2 (Produttoria di prodotto per costante). *Se k è una costante che non dipende dall'indice i :*

$$\prod_{i=1}^n k a_i = k^n \prod_{i=1}^n a_i \quad (1.23)$$

Dimostrazione. Infatti

$$\prod_{i=1}^n k a_i = k a_1 \cdot k a_2 \cdot \dots \cdot k a_n = k^n (a_1 a_2 \dots a_n) = k^n \prod_{i=1}^n a_i$$

□

Scomposizione in sottoinsiemi**Proposition 1.2.3.** *Vale la seguente:*

$$\prod_{i=1}^m a_i \prod_{i=m+1}^n a_i = \prod_{i=1}^n a_i \quad (1.24)$$

Dimostrazione. Infatti

$$\prod_{i=1}^m a_i \prod_{i=m+1}^n a_i = (a_1 a_2 \dots a_m)(a_{m+1} a_{m+2} \dots a_n) = \prod_{i=1}^n a_i$$

□

Remark 18. Generalizzando, se Λ è un insieme di indici, $a : \Lambda \rightarrow \mathbb{C}$ una famiglia di complessi, e J, K sottoinsiemi finiti *disgiunti* di Λ si ha:

$$\prod_{\lambda \in J \cup K} a_\lambda = \prod_{\lambda \in J} a_\lambda \cdot \prod_{\lambda \in K} a_\lambda \quad (1.25)$$

Proposition 1.2.4 (Scomposizione: produttoria di prodotti). *Vale la seguente:*

$$\prod_{i=1}^n a_i b_i = \prod_{i=1}^n a_i \prod_{i=1}^n b_i \quad (1.26)$$

Dimostrazione.

$$\begin{aligned} \prod_{i=1}^n a_i b_i &= a_1 b_1 \cdot a_2 b_2 \cdot \dots \cdot a_n b_n = \\ &= (a_1 a_2 \dots a_n)(b_1 b_2 \dots b_n) \\ &= \prod_{i=1}^n a_i \prod_{i=1}^n b_i \end{aligned}$$

□

Remark 19. Generalizzando, se Λ è un insieme di indici, $a : \Lambda \rightarrow \mathbb{C}$ una famiglia di complessi e $b : \Lambda \rightarrow \mathbb{C}$ è un'altra famiglia di complessi si può definire il prodotto $a \cdot b : \Lambda \rightarrow \mathbb{C}$ delle due famiglie ponendo $(a \cdot b)(\lambda) = a_\lambda \cdot b_\lambda$ per ogni $\lambda \in \Lambda$. Si ha anche che per ogni sottoinsieme finito J di Λ :

$$\prod_{j \in J} (a_j \cdot b_j) = \prod_{j \in J} a_j \cdot \prod_{j \in J} b_j \quad (1.27)$$

Proposition 1.2.5 (Logaritmi e sommatorie). *Vale la*

$$\log \prod_{i=1}^n a_i = \sum_{i=1}^n \log a_i \quad (1.28)$$

Dimostrazione. Infatti:

$$\log \prod_{i=1}^n a_i = \log(a_1 a_2 \dots a_n) = \log a_1 + \log a_2 + \dots + \log a_n = \sum_{i=1}^n \log a_i$$

□

1.3 Esercizi

Exercise 1.3.1 (es 10 pag 34 bps1). Ricavare la formula per la somma dei primi n numeri pari

$$\sum_{k=1}^n (2k)$$

e dimostrarla per induzione

Soluzione. Elaboriamola intanto

$$\sum_{k=1}^n 2k = 2 \sum_{k=1}^n k = 2 \frac{n(n+1)}{2} = n(n+1)$$

Dimostriamo per induzione (anche se non ci sarebbe bisogno, essendo che è moltiplicare per 2 entrambi i membri dell'equazione $\sum_{k=1}^n k = \frac{n(n+1)}{2}$):

- per il passo base

$$\sum_{k=1}^1 2k = 2$$

$$1(1+1) = 2$$

sono uguali quindi il passo base è ok

- per il passo induttivo

$$\sum_{k=1}^{n+1} 2k = \left(\sum_{k=1}^n 2k \right) + n(n+1) = (n+1)(n+2)$$

Quest'ultima è proprio $n(n+1)$ con sostituzione $n \rightarrow n+1$, quindi anche il passo induttivo è ok

Capitolo 2

Calcolo combinatorio

2.1 Introduzione

Definition 2.1.1 (Calcolo combinatorio). Studio di come quantificare raggruppamenti aventi determinate caratteristiche degli elementi di un insieme finito di oggetti.

Remark 20. È fondamentale per il calcolo delle probabilità in quanto spesso la probabilità di un evento è calcolabile come il numero di modi in cui detto evento può verificarsi in rapporto al numero di casi possibili.

Definition 2.1.2 (Principio fondamentale del calcolo combinatorio). Se si realizzano due esperimenti:

- in cui il primo ha m esiti possibili;
- e per ognuno di questi il secondo ha n esiti possibili;
- e l'ordinamento conta per qualificare un esito (ossia sequenze diverse dei singoli esiti dei due esperimenti producono esiti finali distinti):

allora i due esperimenti (considerati congiuntamente) hanno $m \cdot n$ esiti possibili.

Remark 21. Generalizzato, con r esperimenti nel quale il primo abbia n_1 esiti possibili, per ciascuno di questi il secondo ne abbia $n_2 \dots$ per ogni esito dei primi due $r - 1$ l' r -esimo n_r esiti possibili e l'ordinamento conta, allora gli esperimenti hanno in tutto $\prod_{i=1}^r n_i$ esiti possibili.

Definition 2.1.3 (Funzione fattoriale). Il fattoriale di n , indicato con $n!$ è una funzione $f : \mathbb{N} \rightarrow \mathbb{N}$ è definito come il prodotto dei primi n numeri interi:

$$n! = n \cdot (n - 1) \cdot (n - 2) \cdot \dots \cdot 1 \quad (2.1)$$

Si conviene che $0! = 1$.

Remark 22 (Definizione ricorsiva). Dato che $(n - 1) \cdot (n - 2) \cdot \dots \cdot 1 = (n - 1)!$ il fattoriale può esser definito anche come:

$$n! = \begin{cases} 1 & n \in \mathbb{N}, n = 0 \\ n \cdot (n - 1)! & n \in \mathbb{N}, n \neq 0 \end{cases} \quad (2.2)$$

Remark 23 (Una semplificazione utile). Se $0 < k < n$, si ha:

$$\frac{n!}{(n-k)!} = n \cdot (n-1) \cdot (n-2) \cdot \dots \cdot (n-k+1) \quad (2.3)$$

2.2 Casistica principale

Supponendo di voler costruire sottoinsiemi contenenti k elementi scelti tra gli n elementi di un insieme U :

- nel caso in cui l'*ordine* abbia importanza (configurazioni con gli stessi elementi posti in ordine diverso danno origine ad esiti diversi) abbiamo a che fare con:
 - **permutazioni**: disponiamo di $k = n$ slot ed n elementi ($\in U$) da utilizzare per riempirli. Ci interessa sapere in quanti modi si possono ordinare gli n oggetti: ognuno di questi ordinamenti si chiama *permutazione*. Possiamo avere due casi:
 1. permutazioni *semplici*: gli n elementi da ordinare sono unici (ad esempio gli anagrammi della parola “AMORE”);
 2. permutazioni *con ripetizione*: ammettono che un elemento si presenti più volte tra gli n dai quali si può pescare (ad esempio gli anagrammi della parola “PEPPER”).
 - **disposizioni** (che costituiscono una versione generalizzata della permutazioni): gli slot sono in numero $k \leq n$ inferiore (o uguale) rispetto agli elementi n con il quale possiamo riempirli. Di fatto qua si considera che gli n elementi siano tutti distinti/diversi. Abbiamo:
 1. disposizioni *semplici*: i k elementi sono pescati da un insieme di n elementi distinti e una volta che l'elemento è stato scelto esce dal pool degli utilizzabili;
 2. disposizioni *con ripetizione*: ciascun elemento dei n può essere estratto più volte
- se viceversa l'*ordine non ha rilevanza*, ossia sottoinsiemi composti da medesimi elementi posti in ordine differente sono considerati uguali (ad esempio quando si vogliono contare insiemi nell'accezione matematica del termine) si ha a che fare con le **combinazioni**. Le combinazioni semplici sono le più utilizzate e si hanno quando il pool dal quale si pesca è composto da oggetti diversi/distinti tra loro.

2.2.1 Permutazioni

Proposition 2.2.1 (Permutazioni semplici). *Il numero di permutazioni di n elementi distinti in n slot è:*

$$P_n = n \cdot (n-1) \cdot (n-2) \cdot \dots \cdot 1 = n! \quad (2.4)$$

Dimostrazione. Nella prima posizione possiamo porre n alternative, nella seconda $n-1$ (visto che una è già andata nella prima), e così via; arrivando così all'ultima posizione rimane un solo oggetto possibile degli n iniziali. Pertanto per il principio fondamentale del calcolo combinatorio si conclude. \square

Remark 24. Nel caso in cui vi siano elementi ripetuti/uguali dai quali pescare (ad esempio se vogliamo permutare le lettere di “PEPPER”) vogliamo che il numero di esiti complessivi diminuisca (evitando di contare come differenti due configurazioni con elementi uguali permutati tra loro)

Proposition 2.2.2 (Permutazioni con ripetizione). *Tra gli n dai quali pescare vi siano $i = 1, 2 \dots r$ elementi univoci che si possono ripetere, aventi numerosità rispettivamente $k_1, k_2 \dots k_r$ (ossia si ha $\sum_{i=1}^r i \cdot k_i = n$). Le permutazioni uniche (non ripetute) sono:*

$$P_n^{k_1, k_2, \dots, k_r} = \frac{n!}{k_1! \cdot k_2! \cdot \dots \cdot k_r!} \quad (2.5)$$

Dimostrazione. Si parte dal numero di permutazioni degli n oggetti al numeratore. Applicando il principio fondamentale del calcolo combinatorio al contrario, si tratta di dividere queste per il numero delle $k_1!$ permutazioni uguali fra loro (dovute al “girare” di uno stesso elemento), poi per le $k_2!$ permutazioni del secondo elemento multiplo, e così via. \square

Example 2.2.1. Considerando le permutazioni PEPPER ad ogni sequenza univoca (ad esempio REPPEP) corrisponderanno $3!2!$ sequenze che sono di fatto uguali. Pertanto il numero di permutazioni univoche (con ripetizione) di PEPPER saranno $6!/(3! \cdot 2!)$.

Remark 25. La formula delle permutazioni è una generalizzazione e vale in realtà per qualsiasi permutazione, anche senza ripetizioni di elementi. Infatti, se abbiamo elementi univoci, ossia $k_1 = k_2 = \dots = k_r = 1$, otteniamo esattamente la formula delle permutazioni semplici in quanto:

$$P_n^{k_1, k_2, \dots, k_r} = \frac{n!}{k_1! \cdot k_2! \cdot \dots \cdot k_r!} = \frac{n!}{1! \cdot 1! \cdot \dots \cdot 1!} = n! \quad (2.6)$$

2.2.2 Disposizioni

Definition 2.2.1 (Disposizioni semplici). Se il numero degli slot disponibili è inferiore (o uguale) al numero di elementi dai quali si pesca, gli elementi dai quali si pesca sono distinti tra loro e non vengono reinseriti nel pool dove pescare si hanno le disposizioni semplici.

Sono quello che in statistica si chiama *campionamento senza ripetizione*.

Proposition 2.2.3 (Numero di disposizioni semplici). *Il numero $D_{n,k}$ di disposizioni semplici di $k \leq n$ oggetti estratti da un insieme di n oggetti differenti è:*

$$D_{n,k} = n \cdot (n-1) \cdot \dots \cdot (n-k+1) = \frac{n!}{(n-k)!} \quad (2.7)$$

Dimostrazione. Il primo componente di una tale sequenza può essere scelto in n modi diversi, il secondo in $(n-1)$ e così via, sino al k -esimo che può essere scelto in $(n-k+1)$ modi diversi. \square

Remark 26. Le permutazioni semplici (quando $k = n$) sono casi particolari delle disposizioni semplici (quando $k \leq n$):

$$P_n = D_{n,n} = \frac{n!}{(n-n)!} = \frac{n!}{0!} = \frac{n!}{1} = n! \quad (2.8)$$

Definition 2.2.2 (Disposizioni con ripetizione). Le disposizioni con ripetizione sono caratterizzate dal fatto che ciascuno degli n elementi possa essere estratto più volte per riempire i k slot.

Sono quello che in statistica si chiama *campionamento con ripetizione*.

Proposition 2.2.4 (Numero di disposizioni con ripetizione). *Il numero di disposizioni con ripetizione di n elementi in k slot:*

$$D'_{n,k} = \underbrace{n \cdot n \cdot \dots \cdot n}_{k \text{ volte}} = n^k \quad (2.9)$$

Dimostrazione. Si hanno n possibilità per scegliere il primo componente, n per il secondo, altrettante per il terzo e così via, sino al k -esimo; si conclude per il principio fondamentale del calcolo combinatorio. \square

2.2.3 Combinazioni

2.2.3.1 Combinazioni semplici

Remark 27. Gli n elementi dai quali si pesca sono univoci: si pescano k elementi, l'ordine/disposizione di questi non è rilevante a qualificare un esito differente. Si hanno le combinazioni semplici che conteggiano il numero di sottoinsiemi di ampiezza definita di un determinato insieme base.

Proposition 2.2.5 (Combinazioni semplici). *Il numero delle combinazioni semplici di n elementi di lunghezza k , indicato con $C_{n,k}$ è:*

$$C_{n,k} = \frac{D_{n,k}}{P_k} = \frac{n \cdot (n-1) \cdot \dots \cdot (n-k+1)}{k!} = \frac{n!}{k! \cdot (n-k)!} = \binom{n}{k} \quad (2.10)$$

Dimostrazione. Analogamente alle disposizioni semplici sceglieremo k elementi da n : si inizierà avendo n possibilità per il primo, sino a $n-k+1$ per il k -esimo. Tuttavia all'interno dei gruppi così determinati ci saranno combinazioni che sono formate dagli stessi elementi di altre, anche se in ordine inverso. Per non contare tali gruppi più volte (dato che l'ordine non interessa), sempre applicando il principio fondamentale del calcolo combinatorio, occorrerà dividere le disposizioni per il numero di permutazioni dei k elementi estratti ($k!$). \square

2.2.3.2 Combinazioni con ripetizione

Remark 28. Nelle combinazioni semplici non è ammesso pescare lo stesso elemento più volte. Una volta estratto non rimane negli oggetti estraibili.

Nelle combinazioni con ripetizione invece vogliamo determinare quanti modi vi sono di scegliere k volte da un insieme di n oggetti diversi tra loro, ammettendo che però uno stesso oggetto possa essere pescato più volte.

L'ordine continua a non essere importante (ci interessa sono quante volte ogni oggetto è stato scelto, non l'ordine con cui esso appare).

Le combinazioni con ripetizione contano i *multiset* (insiemi che ammettono ripetizioni) sottoinsieme di un insieme dato.

Proposition 2.2.6. *Il numero di combinazioni con ripetizione di k oggetti scelti tra n è*

$$C_{n,k}^* = \binom{n+k-1}{k} \quad (2.11)$$

Dimostrazione. Se l'ordine contasse il numero di combinazioni sarebbe n^k , ma questo non è il caso. Per dimostrare la formula risolviamo narrativamente un problema isomorfo (stesso problema con setup differente).

Il problema può essere posto come: porre k palline identiche in n scatole differenti: quello che conta è solamente il numero di palline in ciascuna scatola. Una qualsiasi configurazione può essere rappresentata come una sequenza di $|$ per rappresentare i lati di una scatola e o per rappresentare le palline in essa. Ad esempio ipotizzando di avere $k = 7$ palline e $n = 4$ scatole, per rappresentare una pallina nella prima scatola, due nella seconda, tre nella terza e una nella quarta:

$$|o|oo|ooo|o|$$

Per essere valida ciascuna sequenza deve iniziare e finire con $|$: pertanto si tratta solo di contare il modo in cui si possono riarrangiare i termini rimanenti al suo interno (varie configurazioni di scatole). I termini all'interno dei bordi numero $n+k-1$: di questi k (le palline) ed $((n+k-1)-k) = n-1$ anche (i bordi rimanenti utili per formare le n scatole, una volta che due sono stati impiegati per i lati). La soluzione è pertanto

$$\frac{(n+k-1)!}{k! \cdot (n-1)!} = \binom{n+k-1}{k}$$

□

2.3 Coefficiente binomiale e multinomiale

2.3.1 Coefficiente binomiale

2.3.1.1 Definizione

Remark 29. Approfondiamo il coefficiente che risulta dal calcolo del numero di combinazioni semplici di k elementi presi da n .

Definition 2.3.1 (Coefficiente binomiale). Indicato con $\binom{n}{k}$ e pronunciato “n su k” si definisce come

$$\binom{n}{k} = \frac{n \cdot (n-1) \cdot \dots \cdot (n-k+1)}{k!} = \frac{n!}{k! \cdot (n-k)!}$$

se $k \leq n$. Se $n < k$ si pone $\binom{n}{k} = 0$.

Remark 30. Per quanto riguarda il calcolo a mano, spesso è più utile/veloce la prima definizione, mentre la seconda è più compatta ed utilizzabile nelle parti teoriche.

2.3.1.2 Proprietà

Proposition 2.3.1. *Si ha che:*

$$\boxed{\binom{n}{k} = \binom{n}{n-k}} \quad (2.12)$$

Dimostrazione.

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

□

Remark 31. Una intuizione sul significato di 2.12: per scegliere un comitato di k persone tra n sappiamo che ci sono $\binom{n}{k}$ modi. Un'altro modo di scegliere il comitato è specificare quali $n-k$ non ne faranno parte; specificare chi è nel comitato determina chi non vi è e viceversa. Pertanto i due lati sono uguali dato che sono due modi di contare la stessa cosa.

Remark 32. Esempi notevoli/utili della 2.12 sono:

$$\binom{n}{0} = \binom{n}{n} = 1, \quad \binom{n}{1} = \binom{n}{n-1} = n \quad (2.13)$$

Proposition 2.3.2.

$$\boxed{\binom{n}{k} = \binom{n-1}{k-1} + \binom{n-1}{k}} \quad (2.14)$$

Dimostrazione.

$$\begin{aligned} \binom{n-1}{k-1} + \binom{n-1}{k} &= \frac{(n-1)!}{(k-1)! \cdot (n-k)!} + \frac{(n-1)!}{k! \cdot (n-k-1)!} \\ &= \frac{(n-1)! \cdot k}{k! \cdot (n-k)!} + \frac{(n-1)! \cdot (n-k)}{k! \cdot (n-k)!} \\ &= \frac{(n-1)! \cdot n}{k! \cdot (n-k)!} \\ &= \binom{n}{k} \end{aligned}$$

□

Remark 33. Per il significato di 2.14: se ho un insieme di n oggetti $I_n = \{1, \dots, n\}$ isolando un oggetto (diciamo l' n -esimo) posso dividere i sottoinsiemi di I_n che hanno k oggetti in quelli che non contengono l' n -esimo (che sono $\binom{n-1}{k}$), essendo esattamente i sottoinsiemi di I_{n-1} a k oggetti) ed in quelli che lo contengono, i quali si ottengono aggiungendo n ad un insieme di $k-1$ oggetti di I_{n-1} e quindi sono in numero di $\binom{n-1}{k-1}$ ¹; questi due gruppi di sottoinsiemi di I_n sono evidentemente disgiunte, quindi l'unione ha la somma come cardinale, e quindi si ha la formula.

¹Sarebbero $\binom{n-1}{k-1} \cdot 1$ poiché vi è un solo modo di aggiungere l' n -esimo ad un insieme di $k-1$ elementi già formati (scelti tra $n-1$ elementi disponibili)

Proposition 2.3.3 (Identità di Vandermonde).

$$\binom{m+n}{k} = \sum_{j=0}^k \binom{m}{j} \binom{n}{k-j} \quad (2.15)$$

Dimostrazione. La prova mediante espansione dei termini e forza bruta ce la si può evitare. Una dimostrazione narrativa sul perché l'uguaglianza valga è comunque efficace.

Considerando un gruppo di m uomini ed n donne dal quale un comitato di k persone verrà scelto: ci sono $\binom{m+n}{k}$ per farlo. Se vi sono j uomini nel comitato, allora vi debbono essere $k-j$ donne. Il lato destro dell'uguaglianza somma per il numero j di uomini. \square

Proposition 2.3.4 (Squadra con capitano). Per $k, n \in \mathbb{N}$ con $k \leq n$ si ha

$$n \binom{n-1}{k-1} = k \binom{n}{k} \quad (2.16)$$

Dimostrazione. Una dimostrazione narrativa: consideriamo un gruppo di n persone dal quale una squadra di k verrà scelta; uno di queste sarà capitano. Il numero possibile di team così formati può derivare da (lato sinistro) prima scegliere il capitano tra gli n e poi scegliere i $k-1$ rimanenti tra gli $n-1$ disponibili. Oppure ed equivalentemente scegliendo gli $\binom{n}{k}$ componenti e tra questi sceglierne uno dei k come capitano. \square

2.3.1.3 Origine del nome

Remark 34. Il coefficiente binomiale prende nome dal fatto che determina i coefficienti dello sviluppo della potenza del binomio $(x+y)^n$

Proposition 2.3.5 (Teorema binomiale).

$$(x+y)^n = \sum_{k=0}^n \binom{n}{k} x^k y^{n-k} \quad (2.17)$$

Dimostrazione. Per provare il teorema espandiamo il prodotto:

$$(x+y)^n = \underbrace{(x+y) \cdot (x+y) \cdot \dots \cdot (x+y)}_{n \text{ fattori}}$$

I termini del prodotto $(x+y)^n$ sono ottenuti scegliendo la x o la y da ognuno dei fattori. Vi sono $\binom{n}{k}$ modi per scegliere esattamente k volte x (scegliendo y nei $n-k$ rimanenti): in questi casi si ottiene il termine $x^k y^{n-k}$. Il teorema si ottiene facendo variare il numero k di x scelti e sommando i termini risultati. \square

2.3.2 Il coefficiente multinomiale

2.3.2.1 Definizione

Proposition 2.3.6. Il numero di modi in cui è possibile distribuire n oggetti distinti in r scatole distinte in modo che queste contengano, nell'ordine, n_1, n_2, \dots, n_r oggetti ($\sum_{i=1}^r n_i = n$) è:

$$\boxed{\binom{n}{n_1, n_2, \dots, n_r} = \frac{n!}{n_1! \cdot n_2! \cdot \dots \cdot n_r!}} \quad (2.18)$$

Dimostrazione. Vi sono $\binom{n}{n_1}$ possibili scelte per gli oggetti della prima scatola; per ogni tale scelta vi sono $\binom{n-n_1}{n_2}$ scelte per la seconda; per ogni scelta effettuata nelle prime due vi sono $\binom{n-n_1-n_2}{n_3}$ nella terza e così via. Dal principio fondamentale del calcolo combinatorio discende che il risultato cercato è:

$$\binom{n}{n_1} \cdot \binom{n-n_1}{n_2} \cdot \dots \cdot \binom{n-n_1-\dots-n_{r-1}}{n_r} \quad (2.19)$$

Sviluppando si ha

$$\frac{n!}{(n-n_1)!n_1!} \cdot \frac{(n-n_1)!}{(n-n_1-n_2)!n_2!} \cdot \dots \cdot \frac{(n-n_1-n_2-\dots-n_{r-1})!}{0!n_r!}$$

dalla quale, in seguito alle semplificazioni, si ottiene il coefficiente. \square

Remark 35. Costituisce una generalizzazione del coefficiente binomiale (che si ottiene considerando due scatole).

Remark 36. Il coefficiente multinomiale è la formula che viene utilizzato nelle permutazioni con ripetizione (utile ad esempio per il numero di permutazioni di una parola con lettere ripetute).

2.3.2.2 Origine del nome

Remark 37. La formula del coefficiente multinomiale determina i coefficienti dello sviluppo di un polinomio di r termini

Proposition 2.3.7 (Teorema multinomiale).

$$(x_1 + x_2 + \dots + x_r)^n = \sum_{\substack{(n_1, n_2, \dots, n_r): \\ n_1 + n_2 + \dots + n_r = n}} \binom{n}{n_1, n_2, \dots, n_r} x_1^{n_1} \cdot x_2^{n_2} \cdot \dots \cdot x_r^{n_r}$$

Dimostrazione. Analoga al caso binomiale. \square

Example 2.3.1. Nello sviluppo del cubo di un trinomio potremmo procedere manualmente:

$$(a + b + c)^3 = a^3 + b^3 + c^3 + 3a^2b + 3a^2c + 3b^2a + 3b^2c + 3c^2a + 3c^2b + 6abc$$

o calcolare più velocemente, ad esempio che:

- il termine $a^2b^0c^1$ presenta come coefficiente:

$$\binom{3}{2, 0, 1} = \frac{3!}{2! \cdot 0! \cdot 1!} = \frac{6}{2 \cdot 1 \cdot 1} = 3$$

- il termine $a^1b^1c^1$ ha invece coefficiente pari a:

$$\binom{3}{1, 1, 1} = \frac{3!}{1! \cdot 1! \cdot 1!} = \frac{6}{1 \cdot 1 \cdot 1} = 6$$

2.4 Calcolo combinatorio e funzioni

Il calcolo combinatorio può essere applicato per contare le funzioni aventi determinate caratteristiche tra due insiemi finiti. Vediamo innanzitutto un criterio utile per contare e poi alcune applicazioni al conteggio delle funzioni.

2.4.1 Principio dell'overcounting

Sia $f : X \rightarrow Y$ suriettiva; si ha allora

$$\text{Card}(X) = \sum_{y \in Y} \text{Card}(f^{-1}(\{y\})) \quad (2.20)$$

In particolare se tutte le fibre $f^{-1}(\{y\})$ hanno una stessa cardinalità, ossia $\text{Card}(f^{-1}(\{y\})) = \alpha$, si ha:

$$\text{Card}(X) = \alpha \text{Card}(Y) \quad (2.21)$$

essendo X una unione disgiunta delle fibre $f^{-1}(\{y\})$ al variare di $y \in Y$.

Anche detto principio del pastore, questo torna utile quando conosciamo la cardinalità di uno dei due insiemi (ad esempio pecore) e desideriamo ricavare quella dell'altro (numero di zampe).

2.4.2 Funzioni (disposizioni con ripetizione)

Si indica con X^{I_p} l'insieme di tutte le funzioni $f : I_p \rightarrow X$ con $\text{Card}(I_p) = p$ e $\text{Card}(X) = m$. Il numero di tutte le funzioni possibili tra i due insiemi è

$$\text{Card}(X^{I_p}) = \underbrace{m \cdot m \cdot \dots \cdot m}_{p \text{ volte}} = m^p \quad (2.22)$$

e corrisponde alle disposizioni con ripetizione, a p a p degli m oggetti di X .

2.4.3 Funzioni iniettive (disposizioni semplici)

Siamo interessati a quantificare la cardinalità del sottoinsieme delle funzioni iniettive $\Lambda(n, p) \subset I_n^{I_p}$ del tipo $f : I_p \rightarrow I_n$. Si ha che

$$\text{Card}(\Lambda(n, p)) = n \cdot (n-1) \cdot \dots \cdot (n-(p-1)) \quad (2.23)$$

vedendo che all'ultimo elemento di I_p ho già fatto $p-1$ collegamenti, quindi me ne rimangono possibili $n-(p-1)$.

2.4.4 Permutazioni di un insieme (permutazioni semplici)

In particolare se $p = n$ si hanno le biiezioni di un insieme I_n in se stesso, ossia le permutazioni dell'insieme, che sono in numero $n!$

2.4.5 Funzioni caratteristiche (coefficiente binomiale)

Il calcolo del numero di sottoinsiemi a p elementi di un insieme di n oggetti equivale a quantificare la cardinalità delle funzioni caratteristiche che scelgono p elementi tra un insieme di n (ossia tali che $\sum \chi(I_n) = p$).

Indicando con $C(n, p)$ l'insieme dei sottoinsiemi di I_n che hanno cardinale p si ha una funzione suriettiva:

$$s : \Lambda(n, p) \rightarrow C(n, p) \quad (2.24)$$

Il dominio $\Lambda(n, p)$ è un insieme di funzioni mentre il codominio $C(n, p)$ è un insieme di insiemi: la funzione suriettiva è quella che ad ogni funzione iniettiva

$f : I_p \rightarrow I_n$ (con $f \in \Lambda(n, p)$) associa l'immagine $f(I_p) \in C(n, p)$, sottoinsieme a p oggetti di I_n .

Essendo che due funzioni iniettive facenti parte del dominio $f, g \in \Lambda(n, p)$ hanno la stessa immagine se e solo se differiscono per una permutazione sul proprio dominio, le fibre di s hanno tutte cardinale $p!$ (ossia ciascun insieme di p elementi si presenta in $p!$ ordini possibili), segue dal principio dell'overcounting che $\text{Card}(\Lambda(n, p)) = p! \text{Card}(C(n, p))$, quindi:

$$\text{Card}(C(n, p)) = \frac{n \cdot (n-1) \cdot \dots \cdot (n-(p-1))}{p!} = \binom{n}{p} = \frac{n!}{p!(n-p)!} \quad (2.25)$$

2.5 Esercizi

Exercise 2.5.1 (Es 3.4 pg 49 de marco). Sia $p \geq 0$ naturale fissato. Mostrare che per ogni naturale $n \geq p$ si ha

$$\sum_{p \leq k \leq n} \binom{k}{p} = \binom{n+1}{p+1}$$

Soluzione. Si ha che

- se $n = p$ l'eguaglianza è verificata in quanto

$$\sum_{p=k=n=a} \binom{a}{a} = \binom{a+1}{a+1} = 1$$

- supponendo sia vera per $n \geq p$ si ha che

$$\begin{aligned} \sum_{p \leq k \leq n+1} \binom{k}{p} &= \sum_{p \leq k \leq n} \binom{k}{p} + \binom{n+1}{p} = \frac{(n+1)!}{(p+1)!(n-p)!} + \frac{(n+1)!}{p!(n-p+1)!} \\ &= \frac{(n+1)!}{(p+1)p!(n-p)!} + \frac{(n+1)!}{p!(n-p+1)(n-p)!} \\ &= \frac{(n-p+1)(n+1)! + (p+1)(n+1)!}{(p+1)p!(n-p)!(n-p+1)} = \dots \\ &= \frac{n+2}{(p+1)!(n-p+1)} \\ &= \binom{n+2}{p+1} = \binom{(n+1)+1}{p+1} \end{aligned}$$

Dove avviene la sostituzione $n \rightarrow n+1$

Exercise 2.5.2 (Es 1.26.1 pag 33 giusti1). Dimostrare per induzione che $n^n \geq n!$.

Soluzione. Per l'induzione:

- se $n = 1$ si ha che $1 \geq 1$
- ipotizzando che valga per il generico $n \geq 1$, multiplico per $n+1 > 0$ entrambi i membri ottenendo

$$n^n(n+1) \geq n!(n+1) = (n+1)!$$

ora notiamo che $n^n \cdot (n+1)$ sono $n+1$ termini e si ha che

$$(n+1)^{n+1} \geq n^n(n+1)$$

perché per i primi n termini di entrambe si ha che $n+1 > n$. Pertanto considerando le due precedenti equazioni

$$(n+1)^{n+1} \geq n^n(n+1) \geq (n+1)!$$

si conclude guardando al primo e terzo membro

Exercise 2.5.3 (Es 1.26.3 pag 33 giusti1). Dimostrare per induzione che $2 \cdot 4 \cdot \dots \cdot 2n = 2^n n!$

Soluzione. Si vuole dimostrare che

$$\prod_{i=1}^n 2i = 2^n n!$$

Per induzione:

- se $i = 1$ si ha che $2 = 2^1 \cdot 1! = 2$ quindi ok
- per il passo induttivo moltiplichiamo entrambi i termini per $2(n+1)$

$$\begin{aligned} \left(\prod_{i=1}^n 2i \right) \cdot (2(n+1)) &= \prod_{i=1}^{n+1} 2i = 2^n n! \cdot (2(n+1)) \\ &= 2^{n+1} \cdot (n+1)! \end{aligned}$$

ed è ok.

Exercise 2.5.4 (Es 1.26.3 pag 33 giusti1). Dimostrare per induzione che $\forall n \geq 4, n! > 2^n$

Soluzione. Si ha:

- per il passo base, per $n = 4$ si ha

$$4! > 2^4 \iff 24 > 16$$

che è verificato

- per il passo induttivo moltiplichiamo entrambi i termini della disequazione generica per $(n+1)$

$$(n+1)! > 2^n(n+1)$$

ora si ha che $2^2(n+1) > 2^{n+1}$ dato che $(n+1) > 2$. Pertanto

$$(n+1)! > 2^n(n+1) > 2^{n+1}$$

e si conclude guardando primo e ultimo termine

Capitolo 3

Introduction

3.1 Probability space

Definition 3.1.1 (Probability space). Considering an experiment, it's a triplet $(\Omega, \mathcal{F}, \mathbb{P})$, used to describe the experiment in mathematical way, composed by a set called sample space Ω , a σ -algebra \mathcal{A} on it (or, same σ -field \mathcal{F}) and a probability function \mathbb{P} .

3.1.1 Sample space, events

Definition 3.1.2 (Sample space, Ω). The (non-null) set of possible outcomes of an experiment, $\Omega = \{\omega_1, \omega_2, \dots\}$, of which *only one will occur*.

Remark 38. The assumption is that before executing the experiment we can know all the possible outcomes.

Example 3.1.1 (Sample space of coin toss). Here $\Omega = \{h, t\}$ while h is one possible outcome. We could be interested in the events outcome is head $\{h\}$ (singleton), outcome is either head or tail, outcome is not a head etc.

Example 3.1.2 (Sample space of two dice throwing). $\Omega = \{(1, 1), (2, 1), \dots, (6, 6)\}$. The event $E = \text{first is one} = \{(1, 1), \dots, (1, 6)\}$

Example 3.1.3 (Sample space of arrival order). In arrival order of a race with 7 numbered horses $\Omega = \{7! \text{ permutations of } (1, 2, 3, 4, 5, 6, 7)\}$.

Example 3.1.4 (Sample space of cars counted at a crossroad during a minute). $\Omega = \{0, 1, 2, \dots\}$

Example 3.1.5 (Sample space of bulb lifetime). Will be a positive real number so $\Omega = \{x \in \mathbb{R}^+ | x \geq 0\}$.

Definition 3.1.3 (Sample space cardinality). Sample spaces of experiments can be *finite* (eg 3.1.1, 3.1.2) *countable* (in bijection with \mathbb{N} , eg 3.1.4) or *non countable* (bijection with \mathbb{R} , eg 3.1.5)

Definition 3.1.4 (Outcome, ω). One possible result of the experiment: $\omega \in \Omega$.

Definition 3.1.5 (Event (E or A)). Any subset of the sample space Ω .

Definition 3.1.6 (Occurred event). E occurred if it contains the result of the experiment.

Remark 39. Since an event is any subset of Ω the following are valid.

Definition 3.1.7 (True event (Ω)). Always occurs, since at least an element of the Ω occurs during the event.

Definition 3.1.8 (Impossible event (\emptyset)). Never occurs.

Definition 3.1.9 (Singleton event (eventi elementari), $\{\omega\}$). Events composed by a single experiment outcome.

Remark 40 (Plotting). With Venn diagram Ω is given by a rectangle, while events are represented by circles.

3.1.1.1 Events algebra

Remark 41. Rules that applies to create new events; inherits from set theory being the events a set.

Definition 3.1.10 (Events union: $A \cup B$). Event that occurs if occurs one of A or B .

Remark 42. The outcomes composing the event are given by union of the outcomes of starting events.

Remark 43. Union can be extended to a numerable infinite number of events

$$E_1 \cup E_2 \cup \dots \cup E_n \cup \dots = \bigcup_{i=1}^{\infty} E_i \quad (3.1)$$

and verifies if at least one of E_i happens.

Definition 3.1.11 (Events intersection $A \cap B$ (A, B or AB)). Event that occurs if occur both A and B .

Remark 44. The outcome composing the event are given by intersection of the outcomes of starting events.

Remark 45. Similarly intersection event can be extended to a numerable infinite set of events

$$E_1 \cap E_2 \cap \dots \cap E_n \cap \dots = \bigcap_{i=1}^{\infty} E_i \quad (3.2)$$

Definition 3.1.12 (Event complement/negation). The negation of the event A , typed \bar{A} or A^c , is the event that happens if A does not: $A^c = \Omega \setminus A$.

Definition 3.1.13 (Events difference $A \setminus B$). Events that occurs when A occurs but not B : $A \setminus B = A \cap \bar{B}$.

Remark 46. The outcome composing the event are given by the set difference $A \setminus B$ outcomes of starting events.

Definition 3.1.14 (Events symmetric difference $A \triangle B$ (xor)). Events that occur if A or B occurs, but not both

Remark 47. The outcome composing the event are given by $(A \cup B) \setminus (A \cap B)$.

Property	Union	Intersection
Idempotenza	$A \cup A = A$	$A \cap A = A$
Elemento neutro	$A \cup \emptyset = A$	$A \cap \Omega = A$
Commutativa	$A \cup B = B \cup A$	$A \cap B = B \cap A$
Associativa	$(A \cup B) \cup C = A \cup (B \cup C)$	$(A \cap B) \cap C = A \cap (B \cap C)$
Distributiva	$A \cup (B \cap C) = (A \cup B) \cap (A \cup C)$	$A \cap (B \cup C) = (A \cap B) \cup (A \cap C)$

Tabella 3.1: Proprietà di unione ed intersezione

Operation properties

Important remark 5 (Events operation properties). Operation properties are the same as set properties and summarized in tab 3.1; same for DeMorgan Laws.

Proposition 3.1.1 (DeMorgan laws). *With two events*

$$\overline{A \cap B} = \overline{A} \cup \overline{B} \quad (3.3)$$

$$\overline{A \cup B} = \overline{A} \cap \overline{B} \quad (3.4)$$

while in the general form

$$\overline{\bigcap_i E_i} = \bigcup_i \overline{E_i} \quad (3.5)$$

$$\overline{\bigcup_i E_i} = \bigcap_i \overline{E_i} \quad (3.6)$$

3.1.1.2 Relationship between events

Definition 3.1.15 (Event inclusion, $A \subseteq B$). Event A is included in B , $A \subseteq B$ if each time A happens, B happens as well.

Example 3.1.6. $E_1 = \{1, 2\}$ (“dice below 3”) is included in $E_2 = \{1, 2, 3\}$ (“dice below 4”)

Definition 3.1.16 (Monotone increasing sequence of events). A sequence of events E_1, E_2, \dots where $E_1 \subseteq E_2 \subseteq \dots$

Definition 3.1.17 (Monotone decreasing sequence of events). A sequence of events E_1, E_2, \dots where $E_1 \supseteq E_2 \supseteq \dots$

Definition 3.1.18 (Incompatibility/disjointness, $A \cap B = \emptyset$). A and B are incompatible (or disjoint) if they can’t verify together, that is, $A \cap B = \emptyset$.

Example 3.1.7 (Incompatible events with two dice throwing). If $A = \{(1, 6), (2, 5), (3, 4), (4, 3), (5, 2), (6, 1)\}$ (two dice sum to 7) and $B = \{(1, 5), (2, 4), (3, 3), (4, 2), (5, 1)\}$ (sum to 6) are incompatible because $A \cap B = \emptyset$.

Remark 48. In Venn diagrams, two disjoint events are represented by non overlapping areas.

Definition 3.1.19 (Pairwise disjointness/incompatibility). Given a collection of events E_i , $1 \leq i \leq \infty$, they are pairwise disjoint if

$$E_i \cap E_j = \emptyset \quad \forall i \neq j$$

Important remark 6. The same can be defined for 3-folded incompatibility or n -folded. Clearly pairwise disjointness implies higher level disjointness (eg 3-folded, etc); viceversa does not happens.

Definition 3.1.20 (Jointly exhaustive events (eventi necessari), $A \cup B = \Omega$). A and B are jointly exhaustive if at least one event occurs, that is $A \cup B = \Omega$.

Remark 49. Same applies for a collection: $E_i, 1 \leq i \leq \infty$ is jointly exhaustive if at least one event occurs $\bigcup_{i=1}^{\infty} E_i = \Omega$

Definition 3.1.21 (Ω partition). It's a set of events $\{E_i\}_{i \in I}, E_i \subseteq \Omega$ which are both disjoint and jointly exhaustive:

$$E_i \cap E_j = \emptyset \quad i \neq j, \quad \bigcup_{i=1}^{\infty} E_i = \Omega$$

Remark 50. If the set of events E_i is finite, countable or uncountable (eg idem the set of index I), the partition of Ω will respectively be called finite, countable or uncountable.

Remark 51. On Venn diagrams it's a set of non overlapping shapes that sum up to Ω .

Example 3.1.8. Suppose $\Omega = \mathbb{R}$, collection of all $\{x\}$ with $x \in \mathbb{R}$ is a partition (not finite nor countable, it's uncountable).

3.1.2 σ -algebra \mathcal{A} (or σ -field \mathcal{F})

Definition 3.1.22 (σ -algebra \mathcal{A} (or σ -field \mathcal{F})). Set of all the possible events of interest, $\mathcal{A} \subseteq \mathcal{P}(\Omega)$ having the following properties

1. $\Omega \in \mathcal{A}$
2. \mathcal{A} is closed under complements: $A \in \mathcal{A} \implies A^c \in \mathcal{A}$
3. \mathcal{A} is closed under *finite* or *countable* unions (and intersection as well): if $E_1, E_2, \dots \in \mathcal{A}$ is a finite or countable set of events then $\bigcup_{i=1}^{\infty} E_i \in \mathcal{A}$

Lemma 3.1.2. Thus we have that $\emptyset \in \mathcal{A}$ and \mathcal{A} is closed under finite or countable intersection as well:

$$\emptyset = \Omega^c \in \mathcal{A}$$

$$E_1, E_2, \dots \in \mathcal{A} \implies \bigcap_{i=1}^{+\infty} E_i = \left(\bigcup_{i=1}^{+\infty} E_i^c \right)^c \in \mathcal{A}$$

the last by applying proprieties 2, 3 of the definition and DeMorgan's laws.

Important remark 7 (The idea). Events are subset of Ω but it's not needed all the subsets of Ω , elements of $\mathcal{P}(\Omega)$, to be events (for technical complex reasons). It suffices for us to think of the collection of events as a subcollection $\mathcal{A} \subseteq \mathcal{P}(\Omega)$ of the power set of the sample space, having certain reasonable/minimal properties. The idea is that:

- \mathcal{A} can be thought as the set of all possible events that are relevant regarding the considered experiment (probabilistic meaning of \mathcal{A})

- if I make some operations of interest between events (unions, intersections, complement), I can be confident of being inside the σ -algebra.
- if the set of possible events \mathcal{E} of our interest is not a σ -algebra, then we set $\mathcal{A} = \sigma(\mathcal{E})$ as the minimum σ -algebra containing \mathcal{E} , and “work” with this one.

Example 3.1.9. $\mathcal{A} = \{\emptyset, \Omega\}$ is the least possible (più piccola) σ -algebra

Example 3.1.10. $\mathcal{A} = \{\emptyset, \Omega, A, A^c\}$ is the least possible σ -algebra including A .

Example 3.1.11 (Power set (insieme delle parti) as \mathcal{A}). $\mathcal{A} = \mathcal{P}(\Omega)$ is the most possible sigma field; no other \mathcal{A} can be bigger (in terms of cardinality). If:

- Ω is finite, it can be $\mathcal{A} = \mathcal{P}(\Omega)$.
- Ω is countable (eg \mathbb{N}), its power set can be a σ -algebra (see here).
- Ω is *non countable* (eg $\Omega = \mathbb{R}$), its power set is a too large collection for probabilities to be assigned reasonably (eg all being non negative and singleton events probabilities summing up to 1) to all its members

Important remark 8. In case of $\Omega = \mathbb{R}$ or $\Omega = \mathbb{R}^n$ we consider a particular case of σ -field/algebra called Borel σ -field/algebra

Definition 3.1.23 (Intervals of \mathbb{R}). The intervals of \mathbb{R} are (a, b) , $[a, b]$, $(a, b]$, $[a, b)$ $(-\infty, b]$, $(-\infty, b)$, (a, ∞) , $[a, \infty)$, and \mathbb{R} as well.

Definition 3.1.24 (Borel σ -field on \mathbb{R}). The borel sigma-field on \mathbb{R} , denoted by $\beta(\mathbb{R})$, is the least possible sigma-field including all the \mathbb{R} intervals.

Remark 52. Some remarks:

- if $\Omega = \mathbb{R}$ and \mathcal{E} is a set of intervals of \mathbb{R} but *not* a σ -algebra (eg like borel), by definition it could happen that $(-1, 5) \cup [7, 8] \notin \mathcal{E}$; so the property/definition of borel seem reasonable/desiderable;
- $\beta(\mathbb{R})$ includes all sets which can be obtained, starting from intervals, by a countable numbers of unions, intersections and complements;
- note that $\exists A \subset \mathbb{R}$ such as that $A \notin \beta(\mathbb{R})$; in other terms $\beta(\mathbb{R})$ is *not* the power set of \mathbb{R} .

Example 3.1.12 (singleton events and $\beta(\mathbb{R})$). Singleton events are contained in $\beta(\mathbb{R})$ since can be written as intersection between intervals $x = (x - 1, x] \cap [x, x + 1) \in \beta(\mathbb{R}) \forall x \in \mathbb{R}$.

Definition 3.1.25 (Borel σ -field on \mathbb{R}^n). In the same way, if $\Omega = \mathbb{R}^n$, $\beta(\mathbb{R}^n)$ equals to the least σ -field on \mathbb{R}^n including all sets of the form $I_1 \times I_2 \times \dots \times I_n$, where each I_i is an interval of \mathbb{R} .

Example 3.1.13. Graphically think as set of rectangles in the space, eg if $n = 2$ a set of rectangles $I_1 \times I_2$

3.1.3 Probability measure \mathbb{P}

Remark 53. In our construction the third element is the probability function \mathbb{P} , defined according to three Kolmogorov axioms that specifies basic features of any probability function.

Definition 3.1.26 (Measure). A measure, generally speaking, is a function:

1. assigning a positive number to each set
2. for which measure of union of disjoint set is sum of measure of the sets.

Definition 3.1.27 (Probability function, \mathbb{P}). It's a measure characterized¹ by $\mathbb{P}(\Omega) = 1$, so it's a function $\mathbb{P} : \mathcal{A} \rightarrow [0, 1]$ such that:

$$\mathbb{P}(A) \geq 0, \quad \forall A \in \mathcal{A} \quad (3.7)$$

$$\mathbb{P}(\Omega) = 1 \quad (3.8)$$

$$A_i \cap A_j = \emptyset, \forall i \neq j \implies \mathbb{P}\left(\bigcup_i A_i\right) = \sum_i \mathbb{P}(A_i) \quad (3.9)$$

For the latter one (called σ -additivity), set $\{A_1, A_2, \dots\}$ is a *finite* or *countable* set of incompatible events.

Example 3.1.14. A coin, possibly biased is tossed once. We have $\Omega = \{h, t\}$, $\mathcal{A} = \{\emptyset, \{h\}, \{t\}, \Omega\}$ and a *possible* probability measure (it fullfill the requirements) $\mathbb{P} : \mathcal{A} \rightarrow [0, 1]$ is given by

$$\mathbb{P}(\emptyset) = 0, \quad \mathbb{P}(\{h\}) = p, \quad \mathbb{P}(\{t\}) = 1 - p, \quad \mathbb{P}(\Omega = \{h, t\}) = 1$$

where p is a fixed real number in the interval $[0, 1]$. If $p = \frac{1}{2}$ then we say the coin is *fair* or unbiased.

Definition 3.1.28 (Null event). Events A such as $\mathbb{P}(A) = 0$.

Definition 3.1.29 (Almost sure event). Event A such as $\mathbb{P}(A) = 1$.

Important remark 9 (Null vs impossible events, true vs almost surely events). Null events should not be confused with the impossible event \emptyset : null events are happening all around us, even though they have zero probability (eg what's the chance that a dart strikes any given point of the target at which it's thrown). That is: the impossible event is null, but null events need not to be impossible. Specular considerations for Ω with events A such as $\mathbb{P}(A) = 1$.

3.2 Probability

3.2.1 Immediate or useful general results

Remark 54. Let's see some properties following directly from the definition; in what follows we consider generic events $A, B \subseteq \Omega$.

¹For a measure (in general), it may be that $P(\Omega) = 0$ or $P(\Omega) = +\infty$ as well; but not for probability, for which $\mathbb{P}(\Omega) = 1$.

Proposition 3.2.1.

$$\boxed{\mathbb{P}(\overline{A}) = 1 - \mathbb{P}(A)} \quad (3.10)$$

Dimostrazione.

$$\begin{aligned} \Omega &= A \cup \overline{A} \\ \mathbb{P}(\Omega) &= \mathbb{P}(A \cup \overline{A}) \\ 1 &= \mathbb{P}(A) + \mathbb{P}(\overline{A}) \end{aligned}$$

□

Example 3.2.1. If the probability of having head with coin is $\frac{3}{8}$ then probability of tail have to be $\frac{5}{8}$.

Proposition 3.2.2.

$$\boxed{\mathbb{P}(\emptyset) = 0} \quad (3.11)$$

Dimostrazione. Setting $A = \Omega$ in 3.10,

$$\begin{aligned} \mathbb{P}(\overline{\Omega}) &= 1 - \mathbb{P}(\Omega) \\ \mathbb{P}(\emptyset) &= 1 - 1 \end{aligned}$$

□

Proposition 3.2.3.

$$\boxed{A \subseteq B \implies \mathbb{P}(A) \leq \mathbb{P}(B)} \quad (3.12)$$

Dimostrazione. If $A \subseteq B$, B can be written as union of two incompatible events A and $(B \setminus A)$; applying third axiom

$$\begin{aligned} B &= A \cup (B \setminus A) \\ \mathbb{P}(B) &= \mathbb{P}(A) + \mathbb{P}(B \setminus A) \end{aligned}$$

since $\mathbb{P}(B \setminus A) \geq 0$ by axioms, then $\mathbb{P}(B) \geq \mathbb{P}(A)$,

□

Proposition 3.2.4 (Probability that A occurs but not B).

$$\boxed{\mathbb{P}(A \setminus B) = \mathbb{P}(A \cap \overline{B}) = \mathbb{P}(A) - \mathbb{P}(A \cap B)} \quad (3.13)$$

Dimostrazione. Looking at A as union of incompatible events (think using Venn diagram):

$$\begin{aligned} A &= (A \cap B) \cup (A \cap \overline{B}) \\ \mathbb{P}(A) &= \mathbb{P}(A \cap B) + \mathbb{P}(A \cap \overline{B}) \end{aligned}$$

then we conclude as in proposition.

□

Proposition 3.2.5 (Probability of union).

$$\boxed{\mathbb{P}(A \cup B) = \mathbb{P}(A) + \mathbb{P}(B) - \mathbb{P}(A \cap B)} \quad (3.14)$$

Dimostrazione. Writing $A \cup B$ as union of two incompatible events, we apply axioms and 3.13:

$$\begin{aligned} A \cup B &= A \cup (B \cap \bar{A}) \\ \mathbb{P}(A \cup B) &= \mathbb{P}(A) + \mathbb{P}(B \cap \bar{A}) \\ \mathbb{P}(A \cup B) &= \mathbb{P}(A) + \mathbb{P}(B) - \mathbb{P}(A \cap B) \end{aligned}$$

□

Proposition 3.2.6 (Inclusion/exclusion formula). *Considering a finite union of events, probability of their union is calculated according to the following:*

$$\mathbb{P}\left(\bigcup_{i=1}^n E_i\right) = \sum_{r=1}^n (-1)^{r+1} \sum_{i_1 < \dots < i_r} \mathbb{P}(E_{i_1} \cap E_{i_2} \cap \dots \cap E_{i_r}) \quad (3.15)$$

$$\begin{aligned} &= \sum_i \mathbb{P}(E_i) - \sum_{i < j} \mathbb{P}(E_i \cap E_j) + \sum_{i < j < k} \mathbb{P}(E_i \cap E_j \cap E_k) - \dots \\ &\dots + (-1)^{n+1} \mathbb{P}(E_1 \cap \dots \cap E_n) \end{aligned} \quad (3.16)$$

Dimostrazione. Can be proved by induction, as we'll see in 5.3.5.2. □

Example 3.2.2. In case of three events, E, F, G :

$$\begin{aligned} \mathbb{P}(E \cup F \cup G) &= \mathbb{P}(E) + \mathbb{P}(F) + \mathbb{P}(G) - \mathbb{P}(E \cap F) \dots \\ &\quad - \mathbb{P}(E \cap G) - \mathbb{P}(F \cap G) + \mathbb{P}(E \cap F \cap G) \end{aligned}$$

Proposition 3.2.7 (Boole inequality (on union)).

$$\mathbb{P}(E_1 \cup E_2 \cup \dots \cup E_n) \leq \sum_{i=1}^n \mathbb{P}(E_i) \quad (3.17)$$

Dimostrazione. Done in the following section 5.3.5.2. □

Proposition 3.2.8 (Bonferroni inequality (on intersection)).

$$\mathbb{P}(E_1 \cap E_2 \cap \dots \cap E_n) \geq 1 - \sum_{i=1}^n \mathbb{P}(\bar{E}_i) \quad (3.18)$$

Dimostrazione. In section 5.3.5.2. □

Proposition 3.2.9. *If A_1, A_2, \dots is an increasing sequence of events, so that $A_1 \subseteq A_2 \subseteq \dots$ and we set A as the limit of the union:*

$$A = \bigcup_{i=1}^{+\infty} A_i = \lim_{i \rightarrow +\infty} A_i$$

then it follows that

$$\mathbb{P}(A) = \lim_{i \rightarrow +\infty} \mathbb{P}(A_i) \quad (3.19)$$

Proposition 3.2.10. *Similarly if B_1, B_2, \dots is decreasing sequence of events $B_1 \supseteq B_2 \supseteq \dots$ and we set as B the limit of the intersection:*

$$B = \bigcap_{i=1}^{+\infty} B_i = \lim_{i \rightarrow +\infty} B_i$$

then

$$\mathbb{P}(B) = \lim_{i \rightarrow +\infty} \mathbb{P}(B_i) \quad (3.20)$$

Dimostrazione. We prove only the first; we have that A can be seen as an union of a disjoint family of events

$$A = A_1 \cup (A_2 \setminus A_1) \cup (A_3 \setminus A_2) \cup \dots$$

Thus by definition of the probability function its probability is a sum of the disjoint events (again think with Venn, these are enclosing circles)

$$\begin{aligned} \mathbb{P}(A) &= \mathbb{P}(A_1) + \sum_{i=1}^{+\infty} \mathbb{P}(A_{i+1} \setminus A_i) \\ &= \mathbb{P}(A_1) + \lim_{n \rightarrow +\infty} \sum_{i=1}^{n-1} [\mathbb{P}(A_{i+1}) - \mathbb{P}(A_i)] \\ &= \lim_{n \rightarrow +\infty} \mathbb{P}(A_n) \end{aligned}$$

The last passage involve simplification/elision. For the second results on B , take complements and use the first part. \square

3.2.2 Finite equiprobable Ω and probability evaluation

Remark 55. In previous section we never evaluated a probability. In this one we show how it's done for the particular case where Ω is finite with every $\omega \in \Omega$ having the same probability of occurring.

It's a reasonable assumption in several cases (eg balanced dice, coins etc)

Proposition 3.2.11 (Probability of singleton a event). *If Ω is finite, $\Omega = \{1, 2, \dots, n\}$, and $\mathbb{P}(1) = \mathbb{P}(2) = \dots = \mathbb{P}(n) = p$, being the singleton events disjoint and the probability of their union summing to 1 ($p \cdot n = 1$), we'll have*

$$p = \frac{1}{n}$$

Proposition 3.2.12 (Probability of general event). *Given a generic event E , its probability will be*

$$\mathbb{P}(E) = \frac{\# \text{ of outcomes composing } E}{\# \text{ possible outcomes}} = \frac{|E|}{|\Omega|}$$

Remark 56. In words, number of favorable outcome of event E out of possible outcomes of Ω . Often, count of numerator/denominator uses combinatorics.

Remark 57. Suppose a partition E_1, E_2, \dots of Ω is *finite* or *countable* and we want to assign the same probability to all E_i . Is it possible?

Proposition 3.2.13. *It's possible to assign to element/events of a finite partition of Ω the same probability; if the partition is countable this is no more possible.*

Dimostrazione. If the partition is *finite* in n events E_i , it suffices to assign $\mathbb{P}(E_i) = \frac{1}{n}$, so that $\mathbb{P}(\Omega) = \mathbb{P}(\cup_{i=1}^n E_i) = 1$.
If the partition is countable this is impossible: let's prove it by absurd/contradiction. Suppose be $\mathbb{P}(E_i) = c \geq 0, \forall i$. Then

$$1 = \mathbb{P}(\Omega) = \mathbb{P}(\cup_{i=1}^{\infty} E_i) = \sum_{i=1}^{\infty} \mathbb{P}(E_i) = \sum_{i=1}^{\infty} c = \begin{cases} 0 & \text{if } c = 0 \\ +\infty & \text{if } c > 0 \end{cases}$$

Therefore we have a contraddiction: 1 can't be equal to 0 or $+\infty$ \square

Example 3.2.3 (Concordanza estrazione trial). We have an urn with n numbered balls from 1 to n , we draw without replacement. Let's define C_i = "concordance at trial i " as the selected ball at draw i is numbered i . We are interested in evaluating $\mathbb{P}(E)$ where E = no concordance in n draws. By applying the previous properties:

$$\begin{aligned} \mathbb{P}(E) &= 1 - \mathbb{P}(\text{at least one concordance}) = 1 - \mathbb{P}\left(\bigcup_{i=1}^n C_i\right) \\ &= 1 - \left\{ \sum_i \mathbb{P}(C_i) - \sum_{i < j} \mathbb{P}(C_i \cap C_j) + \sum_{i < j < k} \mathbb{P}(C_i \cap C_j \cap C_k) \dots + (-1)^{n+1} \mathbb{P}(C_1 \cap \dots \cap C_n) \right\} \end{aligned}$$

Now

$$\mathbb{P}(C_i) = \frac{(n-1)!}{n!} = \frac{1}{n}$$

we have n slots, the sequences of balls can be $n!$, while the sequence where i ball is at the i -th place are $(n-1)!$ (fix i in its place and then permute the remaining balls). Furthermore for similar reasons

$$\begin{aligned} \mathbb{P}(C_i \cap C_j) &= \frac{(n-2)!}{n!} \\ \mathbb{P}(C_i \cap C_j \cap C_k) &= \frac{(n-3)!}{n!} \\ &\dots \\ \mathbb{P}(C_1 \cap \dots \cap C_n) &= \frac{1}{n!} \end{aligned}$$

Therefore

$$\mathbb{P}(E) = 1 - \left\{ n \cdot \frac{1}{n} - \binom{n}{2} \frac{(n-2)!}{n!} + \binom{n}{3} \frac{(n-3)!}{n!} \dots + (-1)^{n+1} \frac{1}{n!} \right\}$$

Example 3.2.4 (Birthday problem). Ci sono k persone in una stanza. Assumendo che siano nate in uno dei 365 giorni dell'anno con probabilità uguale per ciascun giorno (escludiamo anni bisestili) e che i compleanni siano indipendenti (es non vi sono gemelli nella stanza), quale è la probabilità che due o più persone

nel gruppo compiano gli anni lo stesso giorno?

La calcoliamo come complemento della probabilità che nessuno faccia compleanno lo stesso giorno: questa è data da casi favorevoli (numero di modi possibili per avere compleanni in date diverse) fratto casi possibili (numero di possibili configurazioni di compleanni). Si ha:

$$\mathbb{P}(k \text{ compleanni diversi}) = \frac{365 \cdot \dots \cdot (365 - k + 1)}{365^k}$$

da cui

$$\mathbb{P}(\text{Almeno due uguali tra } k) = 1 - \frac{365 \cdot \dots \cdot (365 - k + 1)}{365^k}$$

Eseguendo i conti si nota come si supera la probabilità del 50% già con $k = 23$ persone (ossia in un gruppo di 23 persone c'è poco più del 50% di probabilità di averne due o più che fanno gli anni lo stesso giorno) mentre a $k = 57$ la probabilità è già oltre il 99%.

```
prob_birthday <- function(k){
  # vectorized for several k
  num <- unlist(lapply(k, function(k2) prod(seq(365, 365 - k2 + 1))))
  den <- 365^k
  1 - num/den
}
k <- 1:60
round(prob_birthday(k = k), 4)

## [1] 0.0000 0.0027 0.0082 0.0164 0.0271 0.0405 0.0562 0.0743 0.0946 0.1169
## [11] 0.1411 0.1670 0.1944 0.2231 0.2529 0.2836 0.3150 0.3469 0.3791 0.4114
## [21] 0.4437 0.4757 0.5073 0.5383 0.5687 0.5982 0.6269 0.6545 0.6810 0.7063
## [31] 0.7305 0.7533 0.7750 0.7953 0.8144 0.8322 0.8487 0.8641 0.8782 0.8912
## [41] 0.9032 0.9140 0.9239 0.9329 0.9410 0.9483 0.9548 0.9606 0.9658 0.9704
## [51] 0.9744 0.9780 0.9811 0.9839 0.9863 0.9883 0.9901 0.9917 0.9930 0.9941
```

3.2.3 Conditional probability

3.2.3.1 Introduction/definition

Remark 58 (Idea). Often is needed to compute probability of an event in case another happened; or it's easier to compute a probability of event A conditioning on information of another event B .

Definition 3.2.1 (Conditioned probability of A given B). If $\mathbb{P}(B) > 0$ it's defined as

$$\mathbb{P}(A|B) = \frac{\mathbb{P}(A \cap B)}{\mathbb{P}(B)} \quad (3.21)$$

Important remark 10. $\mathbb{P}(A|B) \neq \mathbb{P}(B|A)$; denominators are different.

Remark 59. Limit/extreme cases:

$$\begin{aligned} A \cap B = \emptyset &\implies \mathbb{P}(A|B) = 0 \\ A \subseteq B &\implies \mathbb{P}(A|B) = 1 \end{aligned}$$

3.2.3.2 Probability of intersection

Proposition 3.2.14 (For two events, $\mathbb{P}(A \cap B)$). If $\mathbb{P}(B) > 0$:

$$\boxed{\mathbb{P}(A \cap B) = \mathbb{P}(B) \mathbb{P}(A|B)} \quad (3.22)$$

Symmetrically, if $\mathbb{P}(A) > 0$:

$$\boxed{\mathbb{P}(A \cap B) = \mathbb{P}(A) \mathbb{P}(B|A)} \quad (3.23)$$

Dimostrazione. Algebraic manipulation of 3.21. \square

Proposition 3.2.15 (n events (product rule)). Given $E_1, \dots, E_n \in \mathcal{A}$ if $\mathbb{P}(E_1 \cap E_2 \cap \dots \cap E_{n-1}) > 0$, then:

$$\mathbb{P}\left(\bigcap_{i=1}^n E_i\right) = \mathbb{P}(E_1) \cdot \mathbb{P}(E_2|E_1) \cdot \mathbb{P}(E_3|E_1 \cap E_2) \cdot \dots \cdot \mathbb{P}(E_n|E_1 \cap E_2 \cap \dots \cap E_{n-1})$$

Dimostrazione. To verify it we apply recursively the definition 3.23 to the second member:

$$\mathbb{P}(E_1) \cdot \frac{\mathbb{P}(E_1 \cap E_2)}{\mathbb{P}(E_1)} \cdot \frac{\mathbb{P}(E_1 \cap E_2 \cap E_3)}{\mathbb{P}(E_1 \cap E_2)} \cdot \dots \cdot \frac{\mathbb{P}(E_1 \cap E_2 \cap \dots \cap E_n)}{\mathbb{P}(E_1 \cap E_2 \cap \dots \cap E_{n-1})} \quad (3.24)$$

and after simplifying it remains $\mathbb{P}(E_1 \cap E_2 \cap \dots \cap E_n) = \mathbb{P}\left(\bigcap_{i=1}^n E_i\right)$.

Note that denominators in 3.24 are strictly positive thanks to the hypothesis $\mathbb{P}(E_1 \cap E_2 \cap \dots \cap E_{n-1}) > 0$: since intersection on $n-1$ events is not null, even the intersection of less events will be. \square

Remark 60. In practice we can handle/manipulate events as we prefer, eg:

$$\begin{aligned} \mathbb{P}(E_1 \cap E_2 \cap E_3) &= \mathbb{P}(E_1) \cdot \mathbb{P}(E_2|E_1) \cdot \mathbb{P}(E_3|E_1 \cap E_2) \\ &= \mathbb{P}(E_3) \cdot \mathbb{P}(E_2|E_3) \cdot \mathbb{P}(E_1|E_3 \cap E_2) \end{aligned}$$

3.2.3.3 Law of total probability

Remark 61 (Conditioning for problem solving). Sometimes is difficult to calculate $\mathbb{P}(E)$; this can become easier if we can condition on C (and \bar{C}), and summing up applying the previous formula. It's common practice to condition on hypothesis/hypothetical situation or, in sequential experiment, conditioning on previous steps.

Definition 3.2.2 (LTP with a single event (and its complement)). If E and C are two events we have (with E of interest for probability evaluation) then:

$$\mathbb{P}(E) = \mathbb{P}(C) \mathbb{P}(E|C) + \mathbb{P}(\bar{C}) \mathbb{P}(E|\bar{C}) \quad (3.25)$$

Dimostrazione. We can split E in disjoint union as follows:

$$E = (E \cap C) \cup (E \cap \bar{C})$$

Being disjoint:

$$\begin{aligned}\mathbb{P}(E) &= \mathbb{P}((E \cap C) \cup (E \cap \overline{C})) \\ &= \mathbb{P}(E \cap C) + \mathbb{P}(E \cap \overline{C}) \\ &= \mathbb{P}(C) \mathbb{P}(E|C) + \mathbb{P}(\overline{C}) \mathbb{P}(E|\overline{C})\end{aligned}$$

□

Example 3.2.5. Domani potrebbe o piovere o nevicare, ma i due eventi non si possono verificare contemporaneamente. La probabilità che piova è $2/5$, mentre la probabilità che nevichi è $3/5$. Se pioverà, la probabilità che io faccia tardi a lezione è di $1/5$, mentre la probabilità corrispondente nel caso in cui nevichi è di $3/5$. Calcolare la probabilità che io sia in ritardo.

Si ha P = piove, $N = P^c$ = nevica, R = ritardo; avendo a che fare con una partizione

$$\mathbb{P}(R) = \mathbb{P}(P) \mathbb{P}(R|P) + \mathbb{P}(N) \mathbb{P}(R|N) = \frac{2}{5} \frac{1}{5} + \frac{3}{5} \frac{3}{5} = \frac{11}{25}$$

Theorem 3.2.16 (LTP with a partition). *If C_1, C_2, \dots is a finite or countable partition of Ω , the probability of a generic event E can be written as (disintegrability):*

$$\boxed{\mathbb{P}(E) = \sum_i \mathbb{P}(C_i) \mathbb{P}(E|C_i)} \quad (3.26)$$

Important remark 11. Looking at the formula, here it's not a problem if $\mathbb{P}(C_i) = 0$ (which is at the denominator of $\mathbb{P}(E|C_i)$, which would be undefined); undefined multiplied by zero is not considered in the sum.

Dimostrazione. If C_1, C_2, \dots, C_n is a partition of Ω , we can split E in disjoint pieces by intersection with C_i

$$E = \Omega \cap E = \left(\bigcup_{i=1}^n C_i \right) \cap E = (C_1 \cap E) \cup (C_2 \cap E) \cup \dots \cup (C_n \cap E)$$

Being $(C_i \cap E)$ disjoint probability is the sum:

$$\mathbb{P}(E) = \sum_{i=1}^n \mathbb{P}(C_i \cap E) = \sum_{i=1}^n \mathbb{P}(C_i) \mathbb{P}(E|C_i) \quad (3.27)$$

and in the last passage we substituted 3.23. □

Example 3.2.6 (Esempio Rigo). Having an urn with n_w white and n_b black balls, we draw without replacement. We are interested in $\mathbb{P}(W_2)$ where W_2 = “white ball at second draw”: it is not trivial without formula, since we don't know the result of the first trial. We however can calculate it conditioning on first draw results.

Let's set W_1 = “white at first draw” and B_1 = “black at first draw”; since $\{W_1, B_1\}$ is a finite partition of the sample space of the first trial, we can apply the law of total probabilities:

$$\mathbb{P}(W_2) = \mathbb{P}(W_1) \mathbb{P}(W_2|W_1) + \mathbb{P}(B_1) \mathbb{P}(W_2|B_1)$$

Given that we have $n = n_w + n_b$ balls and we draw without replacement

$$\mathbb{P}(W_1) = \frac{n_w}{n}, \mathbb{P}(B_1) = \frac{n_b}{n}, \mathbb{P}(W_2|W_1) = \frac{n_w - 1}{n - 1}, \mathbb{P}(W_2|B_1) = \frac{n_w}{n - 1},$$

Therefore, overall

$$\mathbb{P}(W_2) = \frac{n_w}{n} \cdot \frac{n_w - 1}{n - 1} + \frac{n_b}{n} \cdot \frac{n_w}{n - 1} = \dots = \frac{n_w}{n}$$

This is a counterintuitive result, since it's the same as drawing *with* replacement. In general if $W_j = \text{white at draw } j$, $\mathbb{P}(W_j)$ is still $\frac{n_w}{n}$. In this case we have to condition on the partition of the first $j - 1$ trials.

For example, considering “ $W_3 = \text{white at draw } 3$ ” the first two draws will have $\Omega = \{ww, wb, bw, bb\}$, so

$$\begin{aligned} \mathbb{P}(W_3) &= \mathbb{P}(ww) \mathbb{P}(W_3|ww) + \mathbb{P}(wb) \mathbb{P}(W_3|wb) + \mathbb{P}(bw) \mathbb{P}(W_3|bw) + \mathbb{P}(bb) \mathbb{P}(W_3|bb) \\ &= \dots = \frac{n_w}{n} \end{aligned}$$

Eg in this case $\mathbb{P}(W_3|ww) = \frac{n_w - 2}{n - 2}$

3.2.3.4 Bayes formula

Theorem 3.2.17 (Bayes formula). *If A, B are two events, with $P(B) > 0$ then*

$$\boxed{\mathbb{P}(A|B) = \frac{\mathbb{P}(A) \cdot \mathbb{P}(B|A)}{\mathbb{P}(B)}} \quad (3.28)$$

Dimostrazione. Substitute 3.23 in 3.21. □

Remark 62 (Decision making and knowledge update). When performing a test to verify an hypothesis, bayes formula is used like this: let H be “my hypothesis is true”, and T “positive test”; then:

$$\mathbb{P}(H|T) = \frac{\mathbb{P}(H) \cdot \mathbb{P}(T|H)}{\mathbb{P}(T)}$$

in this case $\mathbb{P}(H)$ is called *a priori probability* $\mathbb{P}(T|H)$ *likelihood* and $\mathbb{P}(H|T)$ *posterior probability* (the denominator is merely a normalizing constant).

Remark 63 (Bayes in diagnostic: PPV and NPV). If D is “being diseased” and T è “being positive to diagnostic test”, $\mathbb{P}(D|T)$ (applying bayes formula) is Positive predictive value while $\mathbb{P}(\overline{D}|\overline{T})$ is negative predictive value..

Corollary 3.2.18. *Let E be a generic event and C_1, C_2, \dots a finite or countable partition of Ω ; the conditional probability of C_i given E is:*

$$\boxed{\mathbb{P}(C_i|E) = \frac{\mathbb{P}(C_i) \mathbb{P}(E|C_i)}{\sum_i \mathbb{P}(C_i) \mathbb{P}(E|C_i)}}$$

Dimostrazione. We started from $\mathbb{P}(C_i|E)$ defined using Bayes law and then substituted the denominator using the law of total probability:

$$\mathbb{P}(C_i|E) = \frac{\mathbb{P}(C_i) \mathbb{P}(E|C_i)}{\mathbb{P}(E)} = \frac{\mathbb{P}(C_i) \mathbb{P}(E|C_i)}{\sum_i \mathbb{P}(C_i) \mathbb{P}(E|C_i)}$$

□

Remark 64. For example in Bayesian statistics C_1, C_2, \dots are the possible values of a random parameter while E is the observed sample.

Remark 65 (Interpretation). E can be thought as an occurred event/effect that is due to only one of n causes C_i (disjoint, exhaustive: that is one and only one of them surely happened) each one of the cause has probability $\mathbb{P}(C_i)$ to happen.

The theorem allows us to evaluate $\mathbb{P}(C_i|E)$, that is probability that having observed E , this has been caused by C_i . In the process we use prior probability $\mathbb{P}(C_i)$ and likelihood $\mathbb{P}(E|C_i)$ at numerator (denominator is a normalizing constant):

- when prior probability is not known, if the partition is *finite* (see 3.2.13), one can assign a common probability $\mathbb{P}(C_i) = 1/n, \forall i$;
- likelihood is generally easier to know/evaluate;
- we conclude C_i as the most reasonable cause if its $\mathbb{P}(C_i|E)$ is higher than the others;
- the final result depends only on the numerator, being the denominator a normalizing constant common for all C_i (and making posteriors $\mathbb{P}(C_i|E)$ to sum up to 1). For this reason we can write

$$\mathbb{P}(C_i|E) \propto \mathbb{P}(C_i) \mathbb{P}(E|C_i)$$

that is posterior probability is proportional to the prior time likelihood

Important remark 12. It's often useful the simpler version of (where the partition of Ω composed by two events, only one of which is of interest, the other is the complement) reported here:

$$\mathbb{P}(H|T) = \frac{\mathbb{P}(H) \cdot \mathbb{P}(T|H)}{\mathbb{P}(H) \cdot \mathbb{P}(T|H) + \mathbb{P}(\overline{H}) \cdot \mathbb{P}(T|\overline{H})} \quad (3.29)$$

Example 3.2.7 (Moneta bilanciata). Abbiamo una moneta bilanciata e una sbilanciata che cade su testa con probabilità $3/4$. Si sceglie una moneta a caso e la si lancia tre volte; restituisce testa tutte e tre le volte. Quale è la probabilità che la moneta scelta sia quella bilanciata?

Se H è l'evento "testa tre volte" e B è l'evento "scelta la moneta bilanciata"; siamo interessati alla probabilità $\mathbb{P}(B|H)$. Ci risulta tuttavia più semplice trovare $\mathbb{P}(H|B)$ e $\mathbb{P}(H|\overline{B})$ dato che aiuta sapere quale moneta consideriamo per calcolare la probabilità di tre teste. Questo suggerisce l'utilizzo del teorema di Bayes e della legge delle probabilità totali. Si ha

$$\begin{aligned} \mathbb{P}(B|H) &= \frac{\mathbb{P}(B) \cdot \mathbb{P}(H|B)}{\mathbb{P}(B) \cdot \mathbb{P}(H|B) + \mathbb{P}(\overline{B}) \cdot \mathbb{P}(H|\overline{B})} \\ &= \frac{(1/2) \cdot (1/2)^3}{(1/2) \cdot (1/2)^3 + (1/2) \cdot (3/4)^3} \\ &\approx 0.23 \end{aligned}$$

Example 3.2.8 (Test di una malattia rara). Un paziente è testato per una malattia che colpisce l'1% della popolazione. Sia D l'evento che "il paziente ha la

malattia” e T il test è positivo (ossia suggerisce che il paziente abbia la malattia). Il paziente sottoposto al test risulta effettivamente positivo. Supponendo che il test sia accurato al 95%, ossia che $\mathbb{P}(T|D) = 0.95$ (la sensibilità) ma anche che $\mathbb{P}(\bar{T}|\bar{D}) = 0.95$ (la specificità), qual è la probabilità che il paziente abbia effettivamente la malattia data la positività del test?

Applicando la formula di Bayes:

$$\begin{aligned}\mathbb{P}(D|T) &= \frac{\mathbb{P}(D) \mathbb{P}(T|D)}{\mathbb{P}(T)} \\ &= \frac{0.01 \cdot 0.95}{\mathbb{P}(T)}\end{aligned}$$

$\mathbb{P}(T)$ non è così facile da ottenere (necessiterebbe di provare il test su tutta la popolazione), ma il teorema delle probabilità totali ci viene in soccorso:

$$\begin{aligned}\mathbb{P}(D|T) &= \frac{0.01 \cdot 0.95}{\mathbb{P}(D) \mathbb{P}(T|D) + \mathbb{P}(\bar{D}) \mathbb{P}(T|\bar{D})} \\ &= \frac{0.01 \cdot 0.95}{0.01 \cdot 0.95 + 0.99 \cdot 0.05} \\ &\approx 0.16\end{aligned}$$

Pertanto vi è il 16% di probabilità che il paziente sia malato, anche se il test è positivo e lo strumento è affidabile: il fatto è che la malattia è estremamente rara e potrebbe essere un falso positivo, ossia un errore del test applicato (nella maggioranza dei casi) ad individui negativi.

3.3 Independent events

NB: Per rigo potrebbe essere un esercizio verificare indipendenza

Definition 3.3.1 (Independence of a collection (even infinite) of events). In general, a collection (even infinite) of events $\mathcal{E} = \{E_1, E_2, \dots\} \subset \mathcal{A}$ is composed by independent events if, for *every finite* subset of the collection $\{E_1, \dots, E_n\} \subset \mathcal{E}$, we have that

$$\mathbb{P}(E_1 \cap \dots \cap E_n) = \mathbb{P}(E_1) \cdot \dots \cdot \mathbb{P}(E_n)$$

Remark 66. Things become easier when events are independent but in reality this is rarely happening.

3.3.1 Two events

Example 3.3.1 (2 independent events, $A \perp\!\!\!\perp B$). Applying definition 3.3.1 to a collection $\mathcal{E} = \{A, B\}$ of two events we say that A, B are independent if:

$$\boxed{\mathbb{P}(A \cap B) = \mathbb{P}(A) \mathbb{P}(B)} \quad (3.30)$$

Example 3.3.2. Tossing a fair coin two times we have $\Omega = \{ht, hh, th, tt\}$ each outcome with probability $1/4$. Defining $H_i =$ “i-th toss is a head”, we have $H_1 = \{ht, hh\}$, $H_2 = \{th, hh\}$; each has probability $\frac{1}{2}$. We have that $H_1 \cap H_2 = \{hh\}$ and since that

$$\mathbb{P}(H_1 \cap H_2) = \frac{1}{4} = \mathbb{P}(H_1) \cdot \mathbb{P}(H_2) = \frac{1}{2} \cdot \frac{1}{2}$$

the two events are independent: $H_1 \perp\!\!\!\perp H_2$. It makes sense since the result of the first outcome does not affect the next.

Important remark 13 (Conditional probability of independent events). If A and B are independent and at the same time we have that $\mathbb{P}(B) > 0$, then we can redefine conditional probability as

$$\mathbb{P}(A|B) = \frac{\mathbb{P}(A \cap B)}{\mathbb{P}(B)} = \frac{\mathbb{P}(A) \cdot \mathbb{P}(B)}{\mathbb{P}(B)} = \mathbb{P}(A) \quad (3.31)$$

Thus under these conditions $\mathbb{P}(A|B) = \mathbb{P}(A)$.

Remark 67. Think independence in this latter way ($\mathbb{P}(A|B) = \mathbb{P}(A)$) may be clearer (knowing that B occurs or not it's the same, we don't need it), but we can't define independence $\mathbb{P}(A|B) = \mathbb{P}(A)$ out of the box because we're assuming $\mathbb{P}(B) > 0$

Proposition 3.3.1. *If $\mathbb{P}(B) = 0 \vee \mathbb{P}(B) = 1$, then A is independent of B , $\forall A$.*

Dimostrazione.

$$\begin{aligned} \mathbb{P}(B) = 0 &\implies \mathbb{P}(A \cap B) = 0 = 0 \cdot \mathbb{P}(A) = \mathbb{P}(B) \cdot \mathbb{P}(A) \\ \mathbb{P}(B) = 1 &\implies \mathbb{P}(A \cap B) = \mathbb{P}(A) = 1 \cdot \mathbb{P}(A) = \mathbb{P}(B) \cdot \mathbb{P}(A) \end{aligned}$$

□

Important remark 14. The previous results applies even if the two events seems to be somewhat connected. Eg suppose $\mathbb{P}(B) = 0$ and $B \subseteq A$. According to intuition these seems not to be independent because if B happens A happens as well. However logic and math definition/point of view can be different in practice.

Proposition 3.3.2 (Independence and complements). *If A and B are independent then the following couples of events are independent as well: A and \overline{B} , \overline{A} and B , \overline{A} e \overline{B} .*

Dimostrazione. Showing the first; suppose A, B are independent so $\mathbb{P}(A \cap B) = \mathbb{P}(A)\mathbb{P}(B)$. We want to prove

$$\mathbb{P}(A \cap \overline{B}) = \mathbb{P}(A)\mathbb{P}(\overline{B})$$

We split $A = (A \cap B) \cup (A \cap \overline{B})$ in a disjoint union and sum its component probability:

$$\mathbb{P}(A) = \mathbb{P}(A \cap B) + \mathbb{P}(A \cap \overline{B})$$

therefore

$$\begin{aligned} \mathbb{P}(A \cap \overline{B}) &= \mathbb{P}(A) - \mathbb{P}(A \cap B) \\ &= \mathbb{P}(A) - \mathbb{P}(A)\mathbb{P}(B) \\ &= \mathbb{P}(A)[1 - \mathbb{P}(B)] \\ &= \mathbb{P}(A)\mathbb{P}(\overline{B}) \end{aligned}$$

Regarding \overline{A} e B independence (and \overline{A} e \overline{B}) it suffices to swap roles by negation/complement. □

3.3.2 n events

Example 3.3.3 (Independence of n events (finite set)). Again, applying definition 3.3.1 to a finite set of n events $A_1, \dots, A_n \subset \Omega$, we have independence if for any subgroup of m events, $1 < m \leq n$, we have:

$$\mathbb{P}\left(\bigcap_{i=1}^m A_i\right) = \prod_{i=1}^m \mathbb{P}(A_i) \quad (3.32)$$

Example 3.3.4 (Independence and pairwise independence of 3 events). E, F, G are independent if:

$$\begin{aligned} \mathbb{P}(E \cap F) &= \mathbb{P}(E) \mathbb{P}(F) \\ \mathbb{P}(E \cap G) &= \mathbb{P}(E) \mathbb{P}(G) \\ \mathbb{P}(F \cap G) &= \mathbb{P}(F) \mathbb{P}(G) \\ \mathbb{P}(E \cap F \cap G) &= \mathbb{P}(E) \mathbb{P}(F) \mathbb{P}(G) \end{aligned}$$

E, F, G are *pairwise* independent if the first three equation above holds.

Important remark 15. Generally speaking, n -wise independence implies $(n-1)$ -wise of its components but viceversa does not hold; eg having *pairwise independence* of the above three events is not enough to prove their *independence*.

Important remark 16. In general given *any* collection $\mathcal{E} = \{E_1, E_2, \dots\} \subset \mathcal{A}$ of events, it may be that $\mathbb{P}(E_i \cap E_j) = \mathbb{P}(E_i) \cdot \mathbb{P}(E_j)$, $\forall i \neq j$, but \mathcal{E} is not independent. An example follows with three events.

NB: Altro esempio, volendo, rigo lez 2023-09-21.

Example 3.3.5. Throwing two coins ha $\Omega = \{tt, tc, ct, cc\}$. Following events are pairwise independent but not independent:

- $A = \text{"first tail"} = \{th, tt\}$
- $B = \text{"second tail"} = \{ht, tt\}$
- $C = \text{"same result"} = \{hh, tt\}$

Infatti

$$\begin{aligned} \mathbb{P}(A) &= \mathbb{P}(B) = \mathbb{P}(C) = \frac{2}{4} = \frac{1}{2} \\ \mathbb{P}(A \cap B) &= \mathbb{P}(\{tt\}) = \frac{1}{4} = \mathbb{P}(A) \mathbb{P}(B) \\ \mathbb{P}(A \cap C) &= \mathbb{P}(\{tt\}) = \frac{1}{4} = \mathbb{P}(A) \mathbb{P}(C) \\ \mathbb{P}(B \cap C) &= \mathbb{P}(\{tt\}) = \frac{1}{4} = \mathbb{P}(B) \mathbb{P}(C) \\ \mathbb{P}(A \cap B \cap C) &= \mathbb{P}(\{tt\}) = \frac{1}{4} \neq \mathbb{P}(A) \mathbb{P}(B) \mathbb{P}(C) = \frac{1}{8} \end{aligned}$$

Point is: knowing what happened with A and B gives us complete information on C .

Important remark 17 (Independence and complements). Similar to the two events case, for three events, if E, F, G are independent, then E is independent from any event formed by union/intersection/complement of F e G .

Example 3.3.6. E is independent from $F \cup G$ being:

$$\begin{aligned}\mathbb{P}(E \cap (F \cup G)) &= \mathbb{P}((E \cap F) + (E \cap G)) \\ &= \mathbb{P}(E \cap F) + \mathbb{P}(E \cap G) - \mathbb{P}(E \cap F \cap G) \\ &= \mathbb{P}(E)\mathbb{P}(F) + \mathbb{P}(E)\mathbb{P}(G) - \mathbb{P}(E)\mathbb{P}(F \cap G) \\ &= \mathbb{P}(E)[\mathbb{P}(F) + \mathbb{P}(G) - \mathbb{P}(F \cap G)] \\ &= \mathbb{P}(E)\mathbb{P}(F \cup G)\end{aligned}$$

3.3.3 Other stuff

Important remark 18 (Independence and disjointness). These are two different concepts, often confused:

- disjointness is a *relation between events*, depicted on Venn diagrams as non overlapping areas;
- independence is a *relation between probability of events*; since on Venn diagrams probability are not depicted, it's not graphically representable

In general, disjointness and independence have no relation, except in the following case

Proposition 3.3.3. *Let be A, B events with positive probability; if they are disjoint/incompatible then they cannot be independent.*

Dimostrazione. If A, B are disjoint/incompatible it must be:

$$\mathbb{P}(A \cap B) = \mathbb{P}(\emptyset) = 0$$

If they were also independent it should be:

$$\mathbb{P}(A \cap B) = \mathbb{P}(A)\mathbb{P}(B)$$

but since we hypothesized $\mathbb{P}(A), \mathbb{P}(B) > 0$, then $\mathbb{P}(A)\mathbb{P}(B) > 0$, which contradict the previous statement on disjointness. \square

3.4 Further topics

3.4.1 Odds ratio

Definition 3.4.1 (Odds ratio (rapporto a favore)). L'odds ratio di un evento A è definito come

$$\text{OR}(A) = \frac{\mathbb{P}(A)}{\mathbb{P}(\overline{A})} = \frac{\mathbb{P}(A)}{1 - \mathbb{P}(A)} \quad (3.33)$$

ed esprime quanto è più probabile che l'evento si realizzi rispetto al fatto che non si realizzi.

Remark 68. Per convertire da odds ratio a probabilità, come si può verificare sostituendo, si ha:

$$\mathbb{P}(A) = \frac{\text{OR}(A)}{1 + \text{OR}(A)} \quad (3.34)$$

Remark 69. Può essere di interesse la modifica della probabilità che una ipotesi H sia vera $\mathbb{P}(H)$ quando si dispone di informazioni su una prova E ; le probabilità condizionate dato E che H sia vera o meno

$$\begin{aligned}\mathbb{P}(H|E) &= \frac{\mathbb{P}(H \cap E)}{\mathbb{P}(E)} = \frac{\mathbb{P}(H) \mathbb{P}(E|H)}{\mathbb{P}(E)} \\ \mathbb{P}(\bar{H}|E) &= \frac{\mathbb{P}(\bar{H} \cap E)}{\mathbb{P}(E)} = \frac{\mathbb{P}(\bar{H}) \mathbb{P}(E|\bar{H})}{\mathbb{P}(E)}\end{aligned}$$

Definition 3.4.2 (Odds ratio condizionato). L'odds ratio dell'ipotesi H non è più $\frac{\mathbb{P}(H)}{\mathbb{P}(\bar{H})}$, ma a seguito delle conoscenze su (o nell'ipotesi di) E è dato da:

$$\frac{\mathbb{P}(H|E)}{\mathbb{P}(\bar{H}|E)} = \frac{\frac{\mathbb{P}(H) \mathbb{P}(E|H)}{\mathbb{P}(E)}}{\frac{\mathbb{P}(\bar{H}) \mathbb{P}(E|\bar{H})}{\mathbb{P}(E)}} = \frac{\mathbb{P}(H)}{\mathbb{P}(\bar{H})} \cdot \frac{\mathbb{P}(E|H)}{\mathbb{P}(E|\bar{H})} \quad (3.35)$$

Remark 70. A seguito dell'introduzione di una prova l'originale rapporto a favore $\frac{\mathbb{P}(H)}{\mathbb{P}(\bar{H})}$ viene moltiplicato per un secondo termine che ne determina l'eventuale variazione: il rapporto a favore finale (e quindi la probabilità di H) aumenta se E è più probabile quando H è vera che quando H è falsa (secondo termine del prodotto) e diminuisce in caso contrario.

Example 3.4.1. Con riferimento all'esempio un altro modo conveniente era utilizzare 3.4.2 per il calcolo dell'odds ratio (e poi la 3.34 per passare a probabilità), evitando di dover utilizzare il teorema delle probabilità totali:

$$\frac{\mathbb{P}(D|T)}{\mathbb{P}(\bar{D}|T)} = \frac{\mathbb{P}(D) \mathbb{P}(T|D)}{\mathbb{P}(\bar{D}) \mathbb{P}(T|\bar{D})} = \frac{0.01}{0.99} \cdot \frac{0.95}{0.05} \approx 0.19$$

da cui applicando la 3.34 si ha:

$$\mathbb{P}(D|T) \approx 0.19 / (1 + 0.19) \approx 0.16$$

3.4.2 Conditional probability 2

3.4.2.1 È una probabilità

Remark 71. Quando condizioniamo su un evento F , aggiorniamo la nostra idea per essere coerente con questa conoscenza, ponendoci in un universo dove sappiamo che F è accaduto.

Entro questo nuovo universo, tuttavia, le leggi della probabilità funzionano come in precedenza dato che le probabilità condizionate sono probabilità a tutti gli effetti.

Proposition 3.4.1. La probabilità condizionata è una valida funzione di probabilità a tutti gli effetti in quanto rispetta gli assiomi di Kolmogorov. Si ha:

$$0 \leq \mathbb{P}(E|F) \leq 1$$

$$\mathbb{P}(\Omega|F) = 1$$

$$\mathbb{P}\left(\bigcup_{i=1}^{\infty} E_i|F\right) = \sum_{i=1}^{\infty} \mathbb{P}(E_i|F) \quad \text{se } E_i \cap E_j = \emptyset, \forall i \neq j$$

Dimostrazione. Per la prima dobbiamo mostrare che:

$$0 \leq \frac{\mathbb{P}(E \cap F)}{\mathbb{P}(F)} \leq 1$$

La prima disuguaglianza è ovvia, mentre la seconda discende dal fatto che $(E \cap F) \subseteq F$, da cui $\mathbb{P}(E \cap F) \leq \mathbb{P}(F)$.

La seconda segue dalla:

$$\mathbb{P}(\Omega|F) = \frac{\mathbb{P}(\Omega \cap F)}{\mathbb{P}(F)} = \frac{\mathbb{P}(F)}{\mathbb{P}(F)} = 1$$

Per la terza

$$\begin{aligned} \mathbb{P}\left(\bigcup_{i=1}^{\infty} E_i|F\right) &= \frac{\mathbb{P}((\bigcup_{i=1}^{\infty} E_i) \cap F)}{\mathbb{P}(F)} \quad \text{applicata la def. di } \mathbb{P}(A|B); \text{ per la prop. distributiva, poi } \dots \\ &= \frac{\mathbb{P}(\bigcup_{i=1}^{\infty} (E_i \cap F))}{\mathbb{P}(F)} \quad \dots \text{ma dato che si tratta di unione eventi disgiunti} \\ &= \frac{\sum_{i=1}^{\infty} \mathbb{P}(E_i \cap F)}{\mathbb{P}(F)} \quad \text{e portando il denominatore sotto sommatoria} \\ &= \sum_{i=1}^{\infty} \mathbb{P}(E_i|F) \end{aligned}$$

□

Remark 72 (Notazione). A volte si vuole esprimere compattamente la probabilità condizionata di un evento E condizionata al verificarsi di un altro evento F . Per farlo definiamo

$$\tilde{\mathbb{P}}(E) = \mathbb{P}(E|F)$$

Remark 73. Pertanto si ha che ogni probabilità condizionata è una probabilità. Allo stesso modo *tutte le probabilità possono essere pensate come probabilità condizionate*. Vi è sempre qualche informazione di fondo sulla quale condizioniamo anche se non esplicitata. Quando scriviamo pertanto $\mathbb{P}(A)$ stiamo pensando a $\mathbb{P}(A|K)$ con K background knowledge.

3.4.2.2 Risultati

Remark 74. Il fatto che, in seguito a 3.4.1, la probabilità condizionata sia una funzione di probabilità a tutti gli effetti, fa sì che tutti i risultati sviluppati in precedenza (per la probabilità non condizionata) valgano anche per la probabilità condizionata.

Possiamo aggiornare tutti i risultati visti in precedenza aggiungendo F a destra della barra di condizionamento. Ne mostriamo alcuni.

Lemma 3.4.2.

$$\tilde{\mathbb{P}}(\bar{A}) = 1 - \tilde{\mathbb{P}}(A) \tag{3.36}$$

Dimostrazione. Infatti

$$\begin{aligned} 1 - \tilde{\mathbb{P}}(A) &= 1 - \mathbb{P}(A|F) = 1 - \frac{\mathbb{P}(A \cap F)}{\mathbb{P}(F)} = \frac{\mathbb{P}(F) - \mathbb{P}(A \cap F)}{\mathbb{P}(F)} = \frac{\mathbb{P}(\bar{A} \cap F)}{\mathbb{P}(F)} \\ &= \mathbb{P}(\bar{A}|F) = \tilde{\mathbb{P}}(\bar{A}) \end{aligned}$$

□

Lemma 3.4.3 (Probabilità dell'unione e principio di inclusione/esclusione). *Si ha*

$$\tilde{\mathbb{P}}(A \cup B) = \tilde{\mathbb{P}}(A) + \tilde{\mathbb{P}}(B) - \tilde{\mathbb{P}}(A \cap B)$$

o equivalentemente

$$\mathbb{P}(A \cup B|F) = \mathbb{P}(A|F) + \mathbb{P}(B|F) - \mathbb{P}(A \cap B|F)$$

Lemma 3.4.4 (Condizionamento ulteriore). *La probabilità condizionata $A|B$ dove B è un nuovo condizionamento e F è già presente/sottointeso si sviluppa come*

$$\tilde{\mathbb{P}}(A|B) = \frac{\tilde{\mathbb{P}}(A \cap B)}{\tilde{\mathbb{P}}(B)} = \frac{\mathbb{P}(A \cap B|F)}{\mathbb{P}(B|F)} = \frac{\frac{\mathbb{P}(A \cap B \cap F)}{\mathbb{P}(F)}}{\frac{\mathbb{P}(B \cap F)}{\mathbb{P}(F)}} = \mathbb{P}(A|B \cap F)$$

Lemma 3.4.5 (Regola di Bayes con condizionamento ulteriore). *A patto che $\mathbb{P}(A \cap F) > 0$ e $\mathbb{P}(B \cap F) > 0$ si ha*

$$\tilde{\mathbb{P}}(A|B) = \frac{\tilde{\mathbb{P}}(A) \cdot \tilde{\mathbb{P}}(B|A)}{\tilde{\mathbb{P}}(B)} = \frac{\mathbb{P}(A|F) \cdot \mathbb{P}(B|A \cap F)}{\mathbb{P}(B|F)}$$

Lemma 3.4.6 (Odds ratio con condizionamento ulteriore). *Si ha:*

$$\frac{\tilde{\mathbb{P}}(A|B)}{\tilde{\mathbb{P}}(\bar{A}|B)} = \frac{\mathbb{P}(A|B \cap F)}{\mathbb{P}(\bar{A}|B \cap F)} = \frac{\mathbb{P}(A|F) \cdot \mathbb{P}(B|A \cap F)}{\mathbb{P}(\bar{A}|F) \cdot \mathbb{P}(B|\bar{A} \cap F)} \quad (3.37)$$

Lemma 3.4.7 (Teorema delle probabilità totali 1). *La probabilità condizionata dell'evento E può essere spezzata come somma delle probabilità di eventi incompatibili, analogamente a quanto fatto in 3.25*

$$\tilde{\mathbb{P}}(E) = \tilde{\mathbb{P}}(C) \tilde{\mathbb{P}}(E|C) + \tilde{\mathbb{P}}(\bar{C}) \tilde{\mathbb{P}}(E|\bar{C})$$

ossia, equivalentemente

$$\mathbb{P}(E|F) = \mathbb{P}(C|F) \mathbb{P}(E|C \cap F) + \mathbb{P}(\bar{C}|F) \mathbb{P}(E|\bar{C} \cap F)$$

Lemma 3.4.8 (Teorema delle probabilità totali (versione generica)). *Se C_1, \dots, C_n è una partizione di Ω e nell'ipotesi che $\mathbb{P}(C_i \cap F) > 0$ per ogni i , allora analogamente a 3.26 si ha*

$$\tilde{\mathbb{P}}(E) = \mathbb{P}(E|F) = \sum_{i=1}^n \mathbb{P}(C_i|F) \cdot \mathbb{P}(E|C_i \cap F)$$

Example 3.4.2 (Moneta bilanciata 2). Riprendendo l'esempio 3.2.7, supponiamo di aver visto la moneta uscire testa tre volte. Se la rilanciamo quale è la probabilità che esca testa una volta ancora?

Sia H l'evento testa tre volte, e T esce testa anche la quarta volta. Siamo interessati a $\mathbb{P}(T|H)$; la legge delle probabilità totali ci permette di scriverla come media ponderata dei condizionamenti su B (scelta la moneta bilanciata)

$$\begin{aligned} \mathbb{P}(T|H) &= \mathbb{P}(B|H) \mathbb{P}(T|B \cap H) + \mathbb{P}(\bar{B}|H) \mathbb{P}(T|\bar{B} \cap H) \\ &= 0.23 \cdot \frac{1}{2} + (1 - 0.23) \cdot \frac{3}{4} \\ &\approx 0.69 \end{aligned}$$

con $\mathbb{P}(B|H) = 0.23$ come derivato in esempio 3.2.7.

3.4.2.3 Condizionare su più eventi

Spesso si vuole condizionare su più eventi/informazioni, ora abbiamo vari modi per farlo. Ipotizzando di essere interessati a $\mathbb{P}(A|B \cap C)$, ossia di voler condizionare a sia B che C :

- possiamo utilizzare la definizione di probabilità condizionata

$$\mathbb{P}(A|B \cap C) = \frac{\mathbb{P}(A \cap B \cap C)}{\mathbb{P}(B \cap C)}$$

- possiamo utilizzare la regola di Bayes condizionando ulteriormente su C (questo è l'approccio naturale se pensiamo che ogni evento nel nostro problema sia condizionato su C)

$$\mathbb{P}(A|B \cap C) = \frac{\mathbb{P}(A|C) \cdot \mathbb{P}(B|A \cap C)}{\mathbb{P}(B|C)}$$

- viceversa utilizzare la regola di Bayes condizionando ulteriormente su B (questo è l'approccio naturale se pensiamo che ogni evento nel nostro problema sia condizionato su B)

$$\mathbb{P}(A|B \cap C) = \frac{\mathbb{P}(A|B) \cdot \mathbb{P}(C|A \cap B)}{\mathbb{P}(C|B)}$$

3.4.2.4 Indipendenza condizionata, aggiornamento delle stime

Definition 3.4.3 (Indipendenza condizionata). Gli eventi A e B sono indipendenti condizionatamente dato l'evento F se

$$\mathbb{P}(A \cap B|F) = \mathbb{P}(A|F) \cdot \mathbb{P}(B|F) \quad (3.38)$$

Remark 75. Attenzione, due eventi:

- possono essere indipendenti condizionatamente (dato F), ma non indipendenti;
- possono essere indipendenti, ma non indipendenti condizionatamente (dato F);
- possono essere indipendenti condizionatamente dato F ma non dato \bar{F} .

Lo vediamo nei seguenti esempi.

Example 3.4.3 (Eventi indipendenti condizionatamente ma non indipendenti). Tornando al setup di esempio 3.2.7, sia F “ho scelto la moneta bilanciata”, A_1 “primo lancio da testa” e A_2 “secondo lancio da testa”. Condizionatamente a F , A_1 e A_2 sono indipendenti; ma A_1 e A_2 non sono indipendenti da soli perché A_1 fornisce informazioni su A_2 .

Example 3.4.4 (Eventi indipendenti ma non condizionatamente). Siano Alice e Bob sono le uniche due persone che mi telefonano; ogni giorno decidono indipendentemente se farlo e sia A “mi chiama Alice”, B “mi chiama Bob”. Questi

sono eventi indipendenti. Ma supponendo che R “il telefono squilla”, condizionatamente a questo A e B non sono più indipendenti, perché se non è Alice deve essere Bob, ossia

$$\mathbb{P}(B|R) < 1 = \mathbb{P}(B|\bar{A} \cap R)$$

per cui B e \bar{A} non sono condizionalmente indipendenti dato R (e allo stesso modo A e B)

Example 3.4.5. Supponendo che vi siano solo due tipi di classi: classi buone dove se si lavora tanto si prendono buoni voti e classi cattive dove il professore assegna voti a caso. Sia G “classe è buona”, W “si lavora tanto” e A “si prende un bel voto”. Allora W, A sono indipendenti condizionatamente a \bar{G} , ma non lo sono dato G .

Example 3.4.6 (Aggiornamento delle stime (e indipendenza condizionale)). Riprendendo l'esempio 3.4.1 sul test della malattia rara, ipotizziamo che il paziente decida di intraprendere un secondo test; questo è indipendente dal primo test effettuato (condizionatamente allo stato di malattia) e ha la stessa sensibilità e specificità. Il paziente risulta positivo per la seconda volta. Come si aggiorna la sua probabilità di essere effettivamente malato?

Siamo interessati a $\tilde{\mathbb{P}}(D|T_2)$, condizionata a T_1 , dove D è essere malato, T_1 è essere risultati positivi al primo test e T_2 al secondo. Utilizziamo la forma per l'odds ratio per ricondurci in secondo luogo alla probabilità; si ha

$$\begin{aligned} \frac{\tilde{\mathbb{P}}(D|T_2)}{\tilde{\mathbb{P}}(\bar{D}|T_2)} &= \frac{\mathbb{P}(D|T_1 \cap T_2)}{\mathbb{P}(\bar{D}|T_1 \cap T_2)} = \frac{\mathbb{P}(D) \cdot \mathbb{P}(T_1 \cap T_2|D)}{\mathbb{P}(\bar{D}) \cdot \mathbb{P}(T_1 \cap T_2|\bar{D})} \\ &= \frac{\mathbb{P}(D) \cdot \mathbb{P}(T_1|D) \cdot \mathbb{P}(T_2|D)}{\mathbb{P}(\bar{D}) \cdot \mathbb{P}(T_1|\bar{D}) \cdot \mathbb{P}(T_2|\bar{D})} = \boxed{\frac{\mathbb{P}(D|T_1)}{\mathbb{P}(\bar{D}|T_1)} \cdot \frac{\mathbb{P}(T_2|D)}{\mathbb{P}(T_2|\bar{D})}} \\ &= 0.19 \cdot \frac{0.95}{0.05} \approx 3.646 \end{aligned}$$

Di particolare interesse è la seconda riga dove, in contesto di indipendenza condizionale, si vede che aggiorniamo i risultati cui eravamo giunti in precedenza mediante le informazioni sul nuovo test. Passiamo alla probabilità seguendo la consueta formula

$$\mathbb{P}(D|T_1 \cap T_2) = \frac{3.646}{1 + 3.646} = 0.78$$

La probabilità di essere malati in seguito ad un secondo test positivo (indipendente condizionalmente) aumenta molto, da 0.16 a 0.78.

Example 3.4.7 (Calcolo diretto della probabilità). Volendo invece calcolare direttamente la probabilità in un colpo solo si applica Bayes e torna comodo il teorema delle probabilità totali condizionando su D :

$$\begin{aligned} \mathbb{P}(D|T_1 \cap T_2) &= \frac{\mathbb{P}(D) \cdot \mathbb{P}(T_1 \cap T_2|D)}{\mathbb{P}(T_1 \cap T_2)} \\ &= \frac{\mathbb{P}(D) \cdot \mathbb{P}(T_1 \cap T_2|D)}{\mathbb{P}(D) \cdot \mathbb{P}(T_1 \cap T_2|D) + \mathbb{P}(\bar{D}) \cdot \mathbb{P}(T_1 \cap T_2|\bar{D})} \\ &= \frac{\mathbb{P}(D) \cdot \mathbb{P}(T_1|D) \cdot \mathbb{P}(T_2|D)}{\mathbb{P}(D) \cdot \mathbb{P}(T_1|D) \cdot \mathbb{P}(T_2|D) + \mathbb{P}(\bar{D}) \cdot \mathbb{P}(T_1|\bar{D}) \cdot \mathbb{P}(T_2|\bar{D})} \\ &= \frac{0.01 \cdot 0.95 \cdot 0.95}{0.01 \cdot 0.95 \cdot 0.95 + 0.99 \cdot 0.05 \cdot 0.05} = 0.78 \end{aligned}$$

Soffermendoci un attimo sulla equazione prima del calcolo dell'ultima riga, se dividiamo algebricamente per $\mathbb{P}(T_1)$ sia numeratore che denominatore si ottiene:

$$\begin{aligned}\mathbb{P}(D|T_1 \cap T_2) &= \frac{\mathbb{P}(D|T_1) \cdot \mathbb{P}(T_2|D)}{\mathbb{P}(D|T_1) \cdot \mathbb{P}(T_2|D) + \mathbb{P}(\overline{D}|T_1) \cdot \mathbb{P}(T_2|\overline{D})} \\ &= \frac{0.16 \cdot 0.95}{0.16 \cdot 0.95 + 0.84 \cdot 0.05} \approx 0.78\end{aligned}$$

che equivale ad un normale teorema di Bayes dove al posto delle probabilità a priori secca $\mathbb{P}(D)$ che avevamo utilizzato in esempio 3.4.1, abbiamo sostituito i risultati disponibili alla fine del primo test, ossia $\mathbb{P}(D|T_1) = 0.16$ e $\mathbb{P}(\overline{D}|T_1) = 1 - 0.16 = 0.84$; come si nota l'unica cosa che cambia nella formula (anche perché T_1 e T_2 performano allo stesso modo), sono tali parti, evidenziate in rosso. Aggiorniamo dunque i risultati al termine del primo test con le informazioni del secondo test, per arrivare alla probabilità a posteriori $\mathbb{P}(D|T_1 \cap T_2)$. Seguendo questa impostazione, è facile generalizzare ad n test applicando ripetutamente il teorema.

3.4.3 Schema operativo per il calcolo

Il seguente schema può esser utile per la risoluzione dei problemi di calcolo delle probabilità, poiché il calcolo della probabilità di un evento E , anche complesso, ad una serie di operazioni semplici su eventi elementari:

1. individuare correttamente la prova e tutti gli eventi elementari di cui è composta, distinguendo tra prove semplici e prove composte
2. assegnare la probabilità agli eventi elementari distinguendo tra
 - partizioni finite di eventi
 - partizioni numerabili di eventi
 - partizioni non numerabili di eventi

Dedurre poi, a seconda dei casi, la misura della probabilità da ragioni di simmetria, dalla struttura fisica dell'esperimento, da argomentazioni geometriche o analitiche, da esperienze precedenti, da valutazioni soggettive

3. Controllare che l'assegnazione delle probabilità agli eventi elementari rispetti gli Assiomi di Kolmogorov
4. Esplicitare l'evento E di cui si desidera calcolare la probabilità mediante l'unione, la negazione, l'intersezione degli eventi elementari, in numero finito o numerabile.
5. Possono verificarsi tre casi:
 - se l'evento E è costituito dall'unione di eventi elementari, occorre chiedersi e essi siano incompatibili o meno, applicando le rispettive formule;
 - se l'evento E è costituito dall'intersezione di eventi elementari, occorre chiedersi se essi siano indipendenti o meno, applicando poi le rispettive formule

- se l'evento E è rappresentato dalla *negazione di eventi elementari* (o dalla loro unione e/o intersezione) occorre applicare le formule delle probabilità per l'evento negazione, ovvero le leggi di DeMorgan, se convenienti
- 6. Qualora l'evento derivi da più sottoprove, occorre verificarne la indipendenza oppure calcolare la probabilità condizionate corrispondenti ai singoli eventi, applicando il teorema delle probabilità totali oppure il teorema di Bayes
- 7. Controllare alla fine di ogni esercizio la coerenza dei risultati ottenuti, per esempio calcolando in modo indipendente la probabilità di altri eventi.

3.5 Esercizi rigo

Example 3.5.1 (Es rigo). Stai viaggiando su un treno con un amico. Nessuno di voi ha il biglietto e il controllore vi ha beccato. Il controllore è autorizzato a infliggervi una punizione molto particolare. Vi porge una scatola contenente 9 cioccolatini identici, 3 dei quali avvelenati. Vi costringe a sceglierne uno a testa, a turno, e mangiarlo immediatamente.

1. Se scegli prima del tuo amico, qual è la probabilità che tu sopravviva?
2. Se scegli per primo e sopravvivi, qual è la probabilità che anche il tuo amico sopravviva?
3. Se scegli per primo e muori, qual è la probabilità che il tuo amico sopravviva?
4. E' nel tuo interesse far scegliere prima al tuo amico?
5. Se scegli per primo, qual è la probabilità che tu sopravviva, tenendo conto del fatto che il tuo amico resti in vita?

Se A ="primo cioccolatino scelto è non avvelenato", e B ="secondo scelto non avvelenato"

1. $\mathbb{P}(A) = 6/9$
2. $\mathbb{P}(B|A) = 5/8$
3. $\mathbb{P}(B|A^c) = 6/8$
4. $\mathbb{P}(B) = \mathbb{P}(A)\mathbb{P}(B|A) + \mathbb{P}(A^c)\mathbb{P}(B|A^c) = \frac{6}{9}\frac{5}{8} + \frac{6}{9}\frac{6}{8} = \frac{6}{9}$ quindi non vi è vantaggio nello scegliere dopo il tuo amico
5. $\mathbb{P}(A|B) = \frac{\mathbb{P}(A)\mathbb{P}(B|A)}{\mathbb{P}(B)} = \dots = \frac{5}{8}$; notiamo che $\mathbb{P}(A|B) = \mathbb{P}(B|A)$ in accordo con l'osservazione precedente, ossia che l'ordine della scelta non influenzi le probabilità di sopravvivenza

Example 3.5.2 (Rs rigo). Un dado a sei facce non truccato viene lanciato due volte.

1. Scrivere lo spazio di probabilità dell'esperimento.

- Supponiamo che B sia l'evento corrispondente al fatto che il risultato del primo lancio sia un numero non maggiore di 3, e supponiamo anche che C sia l'evento corrispondente al fatto che la somma dei due numeri ottenuti nei due lanci sia uguale a 6. Determinare le probabilità di B e C , e le probabilità condizionali di C dato B , e di B dato C .

Lo spazio di probabilità in questo esperimento è la tripla $(\Omega, \mathcal{A}, \mathbb{P})$, dove:

- $\Omega = \{(1, 1), \dots, (6, 6)\}$
- $\mathcal{A} = \mathcal{P}(\Omega)$
- ciascun punto in Ω ha uguale probabilità di successo, ossia $\mathbb{P}((i, j)) = 1/36$

Per il secondo punto:

- $B = \text{primo lancio} \leq 3 = \{(1, 1), \dots, (1, 6), (2, 1), \dots, (2, 6), (3, 1), \dots, (3, 6)\}$
pertanto $\mathbb{P}(B) = \frac{18}{36}$
- $C = \text{somma} = 6 = \{(1, 5), (5, 1), (2, 4), (4, 2), (3, 3)\}$, $\mathbb{P}(C) = \frac{5}{36}$
- si ha che $C \cap B = \{(1, 5), (2, 4), (3, 3)\}$ quindi $\mathbb{P}(C|B) = \frac{\mathbb{P}(C \cap B)}{\mathbb{P}(B)} = \frac{3/36}{18/36} = \frac{1}{6}$
- $\mathbb{P}(B|C) = \frac{3/36}{5/36} = \frac{3}{5}$

Example 3.5.3. Una scatola contiene n palline di cui k bianche e $n - k$ nere, dove $1 \leq k \leq n - 1$. faccio due estrazioni senza reinserimento. Calcolare la probabilità che la prima sia bianca dato che la seconda estratta è nera. Si ha

$$\begin{aligned} \mathbb{P}(1b|2n) &= \frac{\mathbb{P}(1b \text{ e } 2n)}{\mathbb{P}(2n)} = \frac{\mathbb{P}(1b) \cdot \mathbb{P}(2n|1b)}{\mathbb{P}(1b) \cdot \mathbb{P}(2n|1b) + \mathbb{P}(1n) \cdot \mathbb{P}(2n|1n)} \\ &= \frac{\frac{k}{n} \frac{n-k}{n-1}}{\frac{k}{n} \frac{n-k}{n-1} + \frac{n-k}{n} \frac{n-k-1}{n-1}} = \frac{k(n-k)}{k(n-k) + (n-k)(n-k-1)} \\ &= \frac{k}{k+n-k-1} = \frac{k}{n-1} \end{aligned}$$

Example 3.5.4. Considerati 3 lanci di una moneta:

- costruire lo spazio di probabilità che descrive il numero di teste
- stabilire se gli eventi $A = \{\text{ottengo almeno una testa}\}$ $B = \{\text{ottengo almeno una croce}\}$ sono indipendenti
- calcolare $\mathbb{P}(A \cup B^c)$ e $\mathbb{P}(A|B^c)$

Si ha che

- $(\Omega, \mathcal{P}(\Omega), \mathbb{P})$ definito a partire da $\Omega = \{ttt, ttc, tct, ctt, tcc, ctc, cct, ccc\}$ e $(X(\Omega), \mathcal{P}(X(\Omega)), \nu)$ $X(\Omega) = \{0, 1, 2, 3\}$ e $\nu(E) = \mathbb{P}(X^{-1}(E))$ con, ad

esempio:

$$\nu(t0) = \mathbb{P}(\{ccc\}) = \frac{1}{8}$$

$$\nu(t1) = \mathbb{P}(\{ttc, tct, ctt\}) = \frac{3}{8}$$

$$\nu(t2) = \mathbb{P}(\{tcc, ctc, cct\}) = \frac{3}{8}$$

$$\nu(t3) = \mathbb{P}(\{ttt\}) = \frac{1}{8}$$

2. i due eventi sono indipendenti se

$$\mathbb{P}(A \wedge B) = \mathbb{P}(A) \cdot \mathbb{P}(B)$$

si ha che

$$\mathbb{P}(A \wedge B) = \mathbb{P}(\text{almeno una testa e almeno una croce}) = \mathbb{P}(\{ttc, tct, ctt, tcc, ctc, cct\}) = \frac{6}{8} = \frac{3}{4}$$

$$\mathbb{P}(A) = 1 - \mathbb{P}(\{ccc\}) = \frac{7}{8}$$

$$\mathbb{P}(B) = 1 - \mathbb{P}(\{ttt\}) = \frac{7}{8}$$

$$\frac{3}{4} \neq \frac{7}{8} \cdot \frac{7}{8} = \frac{49}{64}$$

ergo i due eventi non sono indipendenti

3. si ha che $B^c = \{ttt\}$ e $A \cap B^c = \{ttt\}$

$$\mathbb{P}(A \cup B^c) = \mathbb{P}(A) + \mathbb{P}(B^c) - \mathbb{P}(A \cap B^c) = \frac{7}{8} + \frac{1}{8} - \mathbb{P}(\{ttt\}) = \frac{7}{8} + \frac{1}{8} - \frac{1}{8}$$

$$\mathbb{P}(A|B^c)$$

$$= \frac{\mathbb{P}(A \cap B^c)}{\mathbb{P}(B^c)} = \frac{1/8}{1/8} = 1$$

Example 3.5.5. Si consideri $\Omega = \{1, 2, 3, 4, 5, 6, 7, 8, 9, 10\}$ con $\mathbb{P}(\{i\}) = \frac{i}{10}$ $\forall i \in \Omega$:

1. stabilire se gli eventi $A = \{\text{multipli di 2}\}$ e $B = \{\text{multipli di 3}\}$ sono indipendenti

2. dato $C = \{< 6\}$ calcolare $\mathbb{P}(A|C)$ e $\mathbb{P}(B|C)$

Si ha

1.

$$\mathbb{P}(A) = \frac{2}{55} + \frac{4}{55} + \frac{6}{55} + \frac{8}{55} + \frac{10}{55} = \frac{30}{55}$$

$$\mathbb{P}(B) = \frac{3}{55} + \frac{6}{55} + \frac{9}{55} = \frac{18}{55}$$

$$\mathbb{P}(A \cap B) = \mathbb{P}(6) = \frac{6}{55} \neq \mathbb{P}(A) \cdot \mathbb{P}(B)$$

quindi gli eventi non sono indipendenti

2.

$$\begin{aligned}\mathbb{P}(C) &= \frac{1+2+3+4+5}{55} = \frac{15}{55} \\ \mathbb{P}(A|C) &= \frac{\mathbb{P}(A \cap C)}{\mathbb{P}(C)} = \frac{\frac{2+4}{55}}{\frac{15}{55}} = \frac{6}{15} = \frac{2}{5} \\ \mathbb{P}(B|C) &= \frac{\mathbb{P}(B \cap C)}{\mathbb{P}(C)} = \frac{\frac{3}{55}}{\frac{15}{55}} = \frac{1}{5}\end{aligned}$$

Example 3.5.6. Una scatola contiene due palline bianche e una nera. Estraggo una pallina a caso: se bianca lancio un dado e registro il risultato ottenuto, se è nera lancio due dadi e registro il minore dei due. Calcolare la probabilità di ottenere 2 al termine dell'esperimento. Si ha

$$\begin{aligned}\mathbb{P}(2) &= \mathbb{P}(2|\text{bianca}) \cdot \mathbb{P}(\text{bianca}) + \mathbb{P}(2|\text{nera}) \cdot \mathbb{P}(\text{nera}) \\ &= \frac{2}{3} \cdot \mathbb{P}(\{2\}) + \frac{1}{3} \cdot \mathbb{P}(\{(2, 2), (2, 3), (2, 4), (2, 5), (2, 6), (3, 2), (4, 2), (5, 2), (6, 2)\}) \\ &= \frac{2}{3} \cdot \frac{1}{6} + \frac{1}{3} \cdot \frac{9}{36} = \frac{1}{9} + \frac{1}{12} \\ &= \frac{7}{36}\end{aligned}$$

Example 3.5.7 (Esercizio esame rigo). Da un'urna contenente 5 palline bianche e 4 nere effettuiamo estrazioni senza reinserimento. Si determini la probabilità di ottenere una pallina bianca alla terza prova.

$$\begin{aligned}\mathbb{P}(3b) &= \mathbb{P}(3b|1b \cap 2b) \cdot \mathbb{P}(1b \cap 2b) + \mathbb{P}(3b|1n \cap 2b) \cdot \mathbb{P}(1n \cap 2b) + \dots \\ &\quad \dots + \mathbb{P}(3b|1b \cap 2n) \cdot \mathbb{P}(1b \cap 2n) + \mathbb{P}(3b|1n \cap 2n) \cdot \mathbb{P}(1n \cap 2n) \\ \mathbb{P}(1b \cap 2b) &= \frac{5}{9} \cdot \frac{4}{8} \\ \mathbb{P}(1n \cap 2b) &= \frac{4}{9} \cdot \frac{5}{8} \\ \mathbb{P}(1b \cap 2n) &= \frac{5}{9} \cdot \frac{4}{8} \\ \mathbb{P}(1n \cap 2n) &= \frac{4}{9} \cdot \frac{3}{8} \\ \mathbb{P}(3b) &= \frac{5}{9} \cdot \frac{4}{8} \cdot \frac{3}{7} + \frac{4}{9} \cdot \frac{5}{8} \cdot \frac{4}{7} + \frac{5}{9} \cdot \frac{4}{8} \cdot \frac{4}{7} + \frac{4}{9} \cdot \frac{3}{8} \cdot \frac{5}{7} = \dots = \frac{5}{9}\end{aligned}$$

Capitolo 4

Random variables

4.1 Intro

4.1.1 Random variables linking probability spaces

Remark 76. A probability space $(\Omega, \mathcal{A}, \mathbb{P})$ is a particular measurable space.

Definition 4.1.1 (Measurable space). A pair (S, \mathcal{B}) , composed by a set S and a σ -field \mathcal{B} defined on it.

Definition 4.1.2 (Random variable X). A random variable is a *measurable* function $X : \Omega \rightarrow S$ which creates a mapping between a probability space $(\Omega, \mathcal{A}, \mathbb{P})$ and a measurable space (S, \mathcal{B}) by connecting the first two sets.

Definition 4.1.3 (Measurability). Being X *measurable* means that

$$\forall E \in \mathcal{B}, \exists X^{-1}(E) = \{\omega \in \Omega : X(\omega) \in E\} \in \mathcal{A}, \quad (4.1)$$

In words if I take any event of \mathcal{B} , there's a corresponding event in \mathcal{A} that does produce it through X . $X^{-1}(E)$ is called inverse image of the event E .

Remark 77. In practice in this course the measurable spaces (S, \mathcal{B}) of interest will be:

- $(\mathbb{R}, \beta(\mathbb{R}))$: X is called real or univariate random variable, and so is a function of type $X : \Omega \rightarrow \mathbb{R}$
- $(\mathbb{R}^n, \beta(\mathbb{R}^n))$: X is called n -variate random variable or n -dimensional random vector, a function of type $X : \Omega \rightarrow \mathbb{R}^n$

Remark 78 (Interpretation). The interpretation of rv is the following: one makes the experiment and see the resulting outcome $\omega \in \Omega$. Then after observing ω , $X(\omega)$ make a measurement on the outcome.

Remark 79. While the random variable is a *deterministic* mapping, the random part comes from the experiment.

Definition 4.1.4 (Rv support). It's the image $X(\Omega)$, the set of possible mappings, denoted by $R_X = \{x_1, x_2, \dots\}$

Example 4.1.1 (Two coin throws). Two coin throws can generate the following $\Omega = \{tt, th, ht, hh\}$. On this one we can define $X = \text{“sum of heads as follows”}$

$$X(tt) = 2; X(th) = 1; X(ht) = 1; X(hh) = 0;$$

Finally we have that the support is $R_X = \{0, 1, 2\}$.

Definition 4.1.5 (Probability distribution of X (and second probability space)). Given a probability space $(\Omega, \mathcal{A}, \mathbb{P})$, a measurable space (S, \mathcal{B}) , and a random variable $X : \Omega \rightarrow S$ connecting the twos, we can define a further probability space (S, \mathcal{B}, ν) , where the added probability function $\nu : \mathcal{B} \rightarrow [0, 1]$ is defined, using \mathbb{P} , in the following way:

$$\nu(E) = \mathbb{P}(X^{-1}(E)) = \mathbb{P}(\omega \in \Omega : X(\omega) \in E) = \mathbb{P}(X \in E), \quad \forall E \in \mathcal{B} \quad (4.2)$$

ν is called *probability distribution* of X .

Example 4.1.2. If the experiment is to draw one person from a class, $\Omega = \{\text{everyone}\}$, while the random variable X could be height, so if Luca is extracted ($\omega = \text{Luca}$), then $X(\text{Luca}) = 1.78$.

Distribution function ν of X is:

$$\nu(E) = \mathbb{P}(X \in E) = \mathbb{P}(\text{quelli di noi la cui altezza cade in } E)$$

Eg, if $E = (190, 195]$ and only Paolo and Francesca have an height such as that, then

$$\nu(B) = \mathbb{P}(\text{Paolo}) + \mathbb{P}(\text{Francesca})$$

Important remark 19 (Motivation for measurability request). A possible motivation for requiring measurability of X , as we did, is the need to define its distribution ν . Suppose we don't require X to be measurable; thus can be that:

$$\exists E \in \mathcal{B} : X^{-1}(E) \notin \mathcal{A}$$

there's an event of \mathcal{B} with no corresponding event in \mathcal{A} .

In that case $X^{-1}(E)$ does not belong to the domain of \mathbb{P} and thus we cannot define/write $\nu(E) = \mathbb{P}(X^{-1}(E)) = \mathbb{P}(X \in E)$.

Therefore the need to define ν forces us to require X to be measurable.

Important remark 20 (Notation). If we say:

- $X \sim \nu$ means that ν is the probability distribution of the rv X ; for instance considering a real random variable $X : \Omega \rightarrow \mathbb{R}$, if we say $X \sim N(0, 1)$ we are stating that probability distribution of X is standard normal;
- $X \sim Y$ means that X and Y have the same distribution (whatever it is).

4.1.2 Discrete and continuous rvs

Remark 80. Queste sotto sono definizioni utili per fissare i concetti (le definizioni Rigo style son sotto credo)

Definition 4.1.6 (Discrete rv). Rv which cardinality of support is finite or numerable (1-to-1 with \mathbb{N} .)

Example 4.1.3. Head count in two coin throwing is discrete since $\text{Card}(R_X) = |\{0, 1, 2\}| = 3$.

Definition 4.1.7 (Continuous rv). Rv which cardinality of support is not numerable (1-to-1 with \mathbb{R}).

Example 4.1.4. Numbers of minutes T of bulb lifetime is continue because $R_T = \{t \in \mathbb{R} : t > 0\}$

4.2 Distribution (and other) functions

Remark 81. In order to study random variables, an important concept is distribution function (which is the unifying one for continuous and discrete random variables); here we summarize/prove some results.

Important remark 21 (Jargon). When it's said distribution function we mean the cumulative distribution function.

Definition 4.2.1 (Distribution function). If X is a real valued rv, its distribution function $F : \mathbb{R} \rightarrow \mathbb{R}$ is defined as

$$F_X(x) = \mathbb{P}(X \leq x) = \mathbb{P}(X \in (-\infty, x]) = \nu((-\infty, x]), \quad \forall x \in \mathbb{R}$$

Remark 82. For any distribution function F , exists *one and only one* probability measure ν on $\beta(\mathbb{R})$ such that $F(x) = \nu((-\infty, x])$, $\forall x \in \mathbb{R}$ and viceversa.

Distribution functions are in bijection with probability measure on $\beta(\mathbb{R})$; thanks to this, in order to assign a ν on $\beta(\mathbb{R})$ it is enough to assign a distribution function F (once chosen F to it corresponds one and only one ν). And in practical terms, choosing a F is easier than assigning a ν

NB: considerazioni dalla triennale

Proposition 4.2.1 (Fundamental/characterizing properties). *The properties characterizing distribution functions are*

1. $\lim_{x \rightarrow -\infty} F(x) = 0, \quad \lim_{x \rightarrow +\infty} F(x) = 1,$
2. F is not decreasing: if $y > x$ then $F(y) \geq F(x)$;
3. F is right continuous $F(x) = \lim_{y \rightarrow x^+} F(y), \quad \forall x \in \mathbb{R}$

Important remark 22. Any function $F : \mathbb{R} \rightarrow \mathbb{R}$ which satisfies the three properties is a distribution function, that is, there exists a random variable X such that $F(x) = \mathbb{P}(X \leq x), \quad \forall x \in \mathbb{R}$.

Proposition 4.2.2. *Supposing we want to evaluate the probability of a certain point $\mathbb{P}(X = x) = \nu(\{x\}) = \nu((-\infty, x] \setminus (-\infty, x)) = \nu((-\infty, x]) - \nu((-\infty, x))$. The formula is*

$$\mathbb{P}(X = x) = F(x) - F(x^-) \quad (\text{jump of } F \text{ at } x) \quad (4.3)$$

where $F(x^-) = \lim_{y \rightarrow x^-} F(y)$ meaning limit with $y \rightarrow x$ from the left.

Dimostrazione. To prove this, recall (props 3.2.9 and 3.2.10) that for any probability measure \mathbb{P}

- if $A_1 \subseteq A_2 \subseteq \dots$ is a increasing sequence of events, $\mathbb{P}(\cup_n A_n) = \lim_n \mathbb{P}(A_n)$

- if $A_1 \supseteq A_2 \supseteq \dots$ is a decreasing sequence of events, $\mathbb{P}(\cap_n A_n) = \lim_n \mathbb{P}(A_n)$

Now suppose we want to evaluate

$$\mathbb{P}(X < x) = \mathbb{P}\left(\bigcup_{n=1}^{+\infty} \left\{X \leq x - \frac{1}{n}\right\}\right)$$

where we go nearer and nearer to x as n increases. These events are an increasing sequence of events, so

$$\begin{aligned} \mathbb{P}(X < x) &= \mathbb{P}\left(\bigcup_{n=1}^{+\infty} \left\{X \leq x - \frac{1}{n}\right\}\right) = \lim_{n \rightarrow +\infty} \mathbb{P}\left(X \leq x - \frac{1}{n}\right) = \lim_{n \rightarrow +\infty} F\left(x - \frac{1}{n}\right) \\ &= F(x^-) \end{aligned}$$

Finally in order to evaluate $\mathbb{P}(X = x)$ we have:

$$\mathbb{P}(X = x) = \mathbb{P}(X \leq x) - \mathbb{P}(X < x) = F(x) - F(x^-)$$

□

Important remark 23. As a consequences of 4.3, considering the set:

$$\{x \in \mathbb{R} : \mathbb{P}(X = x) > 0\} = \{x \in \mathbb{R} : \nu(\{x\}) > 0\} = \{x \in \mathbb{R} : F(x) > F(x^-)\}$$

- this set is *empty*, if the function is continuous: in other words the distribution function is *continuous* if and only if the jump is 0 at each point or in other words

$$F \text{ is continuous} \iff \mathbb{P}(X = x) = 0, \forall x \in \mathbb{R}$$

- its cardinality can *at most be countable* (for a calculus result): it can be countable (eg for Poisson, negative binomial etc) or can be finite as well. But can't be uncountable.

4.2.1 Types of RVs

Important remark 24 (RV types). Real random variables can be *discrete*, *singular continuous* (we can ignore it) or *absolutely continuous*. The following result is theoretically important.

Proposition 4.2.3. *If ν is any probability measure on $\beta(\mathbb{R})$, there exists a unique triplets (a, b, c) such that:*

- $a, b, c \geq 0$
- $a + b + c = 1$
- $\nu = a\nu_1 + b\nu_2 + c\nu_3$

where ν_1 is discrete probability measure, ν_2 is singular continuous probability measure, ν_3 is absolutely continuous probability measure.

Dimostrazione. We skip it.

□

Important remark 25. Thanks to the above thm

- if we are able to describe a discrete probability measure, a singular continuous probability measure and an absolute continuous probability measure, we are able to describe ANY probability measure on $\beta(\mathbb{R})$.
- any ν can be written as this mix of this three kind of rv. Clearly, eg

$$\begin{aligned} a = 1, b = c = 0 &\implies \nu = \nu_1 \text{ is discrete} \\ c = 1, a = b = 0 &\implies \nu = \nu_3 \text{ is absolutely continuous} \end{aligned}$$

This is the reason to focus on the three types, of which *only discrete and absolutely continuous are of interest for practical applications*.

Important remark 26. In this course we speak indifferently like:

$$X \text{ is discrete} \iff \nu \text{ is discrete} \iff F \text{ is discrete}$$

Similarly for singular and absolutely continuous rv

4.2.2 Discrete rvs

Definition 4.2.2 (Discrete rv). X is discrete if and only if $\exists B \subset \mathbb{R}$, with B finite or countable such that $\mathbb{P}(X \in B) = 1$.

Example 4.2.1 (Examples of discrete rvs). Some are:

- the degenerate rv, δ_a , where $B = \{a\}$ and thus $P(X \in \{a\}) = 1$; its distribution function is defined as

$$F(x) = \mathbb{P}(X \leq x) = \begin{cases} 1 & x \geq a \\ 0 & x < a \end{cases}$$

- binomial, then $B = \{0, 1, \dots, n\}$;
- Poisson, $B = \{0, 1, \dots\}$.

4.2.3 Singular continuous rvs

Remark 83. As we have said probability is a measure. In general

Definition 4.2.3. A measure m is a function that, considered a single set X and a *finite* or *numerable* set of incompatible events X_1, X_2, \dots

$$m(X) \geq 0, \quad \forall X \quad (4.4)$$

$$X_i \cap X_j = \emptyset, \forall i \neq j \implies m\left(\bigcup_i X_i\right) = \sum_i m(X_i) \quad (4.5)$$

Important remark 27. The *Lebesgue measure* in \mathbb{R} is the only measure on $\beta(\mathbb{R})$ that has this property, applied to an interval:

$$m(a, b] = b - a, \quad \forall a < b \quad (4.6)$$

where m is the Lebesgue measure of the interval. Regarding the measure a point, countable and uncountable sets (the real line) Lebesgue measure

$$\begin{aligned} m(\{x\}) &= 0, & \forall x \in \mathbb{R} \\ m(X) &= \sum_{x \in X} m(\{x\}) = \sum_{x \in X} 0 = 0 & \forall X \subset \mathbb{R} : X \text{ is countable} \\ m(\mathbb{R}) &= +\infty \end{aligned}$$

Definition 4.2.4 (Singular continuous rvs). X is a singular continuous random variable if both

1. the distribution function F is continuous
2. its first derivative is null ($F'(x) = 0$) *almost everywhere* with respect to the Lebesgue measure m (written concisely as “m.a.e.”):

$$m(\{x \in \mathbb{R} : F'(x) \neq 0\}) = 0$$

Important remark 28. Note that

- first derivative $\neq 0$ when it doesn't exist (eg left and right limit are different) or exists but is not 0;
- for this kind of rv, distribution may not be differentiable or with derivative 0 at every point
- however these $F'(x) \neq 0$ points are a finite or at most countable set of points.

Remark 84. For *discrete* RVs actually is the same: $F'(x) = 0$ mae (think step F functions) given that:

$$m(\{x \in \mathbb{R} : F'(x) \neq 0\}) = m(\{\text{jump points of } F\}) = 0$$

as the set $\{\text{jump points of } F\}$ is finite or countable.

However, if X is discrete, F is certainly discontinuous.

Remark 85. These variables

- seems to be a somewhat hybrid between discrete and absolutely continuous rv (since have characteristic from both the distribution), that is $F'(x) = 0$ mae from the discrete RV, and continuous F from absolutely continuous;
- are not usually used for describing real phenomena, and we will not consider them in what follows.

4.2.4 Absolutely continuous rvs

Example 4.2.2. eg exponential, beta, uniform, normal ...

Definition 4.2.5 (Absolutely continuous rv). X is absolutely continuous if and only if exists a function $f : \mathbb{R} \rightarrow \mathbb{R}$, called density, such that:

1. $f \geq 0$ (density is non negative)

2. f is integrable
3. the distribution function at point x can be written as (Lebesgue¹) integral of density function f

$$F(x) = \int_{-\infty}^x f(t) dt, \quad \forall x \in \mathbb{R}$$

Important remark 29 (Probability of an event). With absolutely continuous random variable the probability of an event $E \in \beta(\mathbb{R})$ is

$$\mathbb{P}(X \in E) = \int f(t) \mathbb{1}_E(t) dt = \int_E f(t) dt, \quad \forall E \in \beta(\mathbb{R})$$

where we denoted $\mathbb{1}_E(t)$ as the indicator function of the set E , that is

$$\mathbb{1}_E(t) = \begin{cases} 1, & t \in E \\ 0, & t \notin E \end{cases}$$

Important remark 30. Some properties for these RVs:

- $F' = f$ m.a.e:

$$m(\{x \in \mathbb{R} : f(x) \neq F'(x)\}) = 0$$

that is supposing we collect all the points where density doesn't equal the derivative of the distribution function, then they can differ at most in a countable set of $x \in \mathbb{R}$

- from the previous point, if f_1 and f_2 are both densities of the same RV X , can we say $f_1 = f_2$, that density is *unique*?
Since f_1 and f_2 are densities, $f_1 = F'$ m.a.e and $f_2 = F'$ m.a.e, so we have $f_1 = F' = f_2$ m.a.e that is

$$m(\{x \in \mathbb{R} : f_1(x) \neq f_2(x)\}) = 0$$

so the density f is *almost everywhere unique* (can be different but at most in a countable set of points).

Example 4.2.3. Regarding the last property, consider a standard normal $X \sim N(0, 1)$ which is absolutely continuous, having density

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

Now we define a new density which is different from the standard normal in a countable set \mathbb{Q} of points²:

$$g(x) = \begin{cases} f(x) & \text{if } x \notin \mathbb{Q} \\ 1 + \sin(\log|x| + 3), & \text{if } x \in \mathbb{Q} \end{cases}$$

¹In generale l'integrale di Lebesgue è una generalizzazione dell'integrale di Riemann e coincide con quest'ultimo sotto condizioni abbastanza generali

² \mathbb{Q} has two properties: it's a *countable* set and it's *dense*, that is $\forall a, b \in \mathbb{Q}, \exists q \in \mathbb{Q}$ such that $a < q < b$

We can say that:

$$m(\{f \neq g\}) \leq \underbrace{m(\mathbb{Q}) = 0}_{\text{being countable}}$$

Therefore the function f agrees with g m.a.e.

Thus f and g are *both* densities for X standard normal.

Remark 86. Another important property of absolutely continuous rvs is the following characterization

Theorem 4.2.4 (Absolutely continuous RV characterization). *X is absolutely continuous if and only if, for every set (event) with lebesgue measure 0, this set has probability 0*

$$X \text{ is absolutely continuous} \iff \begin{cases} \mathbb{P}(X \in A) = 0 \\ \forall A \in \beta(\mathbb{R}), \text{ such that } m(A) = 0 \end{cases}$$

Remark 87. Quindi non solo punti singoli hanno probabilità nulla ma anche un insieme finito o al più numerabile la ha.

4.3 OLD: Functions of random variables

4.3.1 Discrete rvs: PMF, CDF

Definition 4.3.1 (Probability mass function). Given a rv $X : \Omega \rightarrow \mathbb{R}$, PMF is a function $p : \mathbb{R} \rightarrow \mathbb{R}$ taking the outcome of the rv and giving its probability

$$p_X(x) = \mathbb{P}(X = x) = \begin{cases} \mathbb{P}(X(\omega) = x) & \text{se } x \in X(\Omega) \\ 0 & \text{se } x \in \mathbb{R} \setminus X(\Omega) \end{cases} \quad (4.7)$$

Proposition 4.3.1 (Valid PMF). *If X is a discrete rv with support $X(\Omega) = \{x_1, x_2, \dots\}$, a valid PMF p_X satisfies:*

$$p_X(x) \geq 0, \quad \forall x \in \mathbb{R} \quad (4.8)$$

$$\sum_{x \in \mathbb{R}} p_X(x) = 1 \quad (4.9)$$

Dimostrazione. Il primo criterio deve esser valido dato che la probabilità è non negativa. Il secondo deve essere valido dato che gli eventi $X = x_1, X = x_2, \dots$ sono disgiunti e X dovrà assumere pur qualche valore:

$$\begin{aligned} \sum_{x \in \mathbb{R}} p_X(x) &= \sum_{x \in X(\Omega)} p_X(x) = \sum_j \mathbb{P}(X = x_j) = \mathbb{P}\left(\bigcup_j \{X = x_j\}\right) \\ &= \mathbb{P}(X = x_1 \text{ or } X = x_2 \dots) = 1 \end{aligned}$$

□

Example 4.3.1. In two coins throwing 4.1.1

$$p_X(X = 0) = 1/4$$

$$p_X(X = 1) = 1/2$$

$$p_X(X = 2) = 1/4$$

and $p_X(x) = 0$ for $x \notin \{0, 1, 2\}$.

Definition 4.3.2 ((Cumulative) distribution function (CDF)). Given a discrete rv X its defined as:

$$F_X(x) = \mathbb{P}(X \leq x) = \sum_{x_j \in X(\Omega): x_j \leq x} p_X(x_j) \quad (4.10)$$

Remark 88 (Function shape). If X is discrete, $F_X(x)$ has stairway shape with finite or numerable steps on values of the support x_1, x_2, \dots : the step height is $p_X(x_1), p_X(x_2), \dots$

Proposition 4.3.2 (Valid CDF). If X is a discrete rv with support $X(\Omega) = \{x_1, x_2, \dots\}$, a valid CDF F_X must satisfy

$$x_1 \leq x_2 \implies F_X(x_1) \leq F_X(x_2) \quad (4.11)$$

$$\lim_{x \rightarrow x_j^+} F_X(x) = F_X(x_j) \quad (\text{right continuous}) \quad (4.12)$$

$$\lim_{x \rightarrow -\infty} F_X(x) = 0, \quad \lim_{x \rightarrow +\infty} F_X(x) = 1 \quad (4.13)$$

Dimostrazione. La prima è giustificata dal fatto che dato che, dato che l'evento $\{X \leq x_1\}$ si verifica sempre quando si verifica $\{X \leq x_2\}$ allora $\mathbb{P}(X \leq x_1) \leq \mathbb{P}(X \leq x_2)$.

La continuità da destra deriva dall'aver definito $F_X(x_0)$ come $\mathbb{P}(X \leq x_0)$ (coerentemente con la letteratura internazionale prevalente); altri autori definiscono $F_X(x_0) = \mathbb{P}(X < x_0)$, il che implica la continuità da sinistra.

Per la terza, dato che $F_X(x_{\min}) = 0$ con $x_{\min} = \min(x_1, x_2, \dots)$ e $-\infty < x_{\min}$ allora per la prima proprietà si ha che $F(-\infty) \leq 0$, ma non potendo una probabilità essere negativa, sarà nulla, dunque si conclude che $\lim_{x \rightarrow -\infty} F_X(x) = 0$. Altresì sfruttando sempre il fatto che $\{X = x_j\}$ sono eventi indipendenti

$$\lim_{x \rightarrow +\infty} F_X(x) = \sum_{x_j \in X(\Omega)} p_X(x_j) = 1$$

□

Example 4.3.2. Dato l'esperimento lancio di due dati, l'evento X somma degli esiti ha PMF e CMF riportate in figura 4.1. Ad esempio $\mathbb{P}(X = 2) = \mathbb{P}(\{1, 1\}) = (\frac{1}{6})^2 = 1/36 \approx 0.02778$. I "salti" nella CDF sono di entità pari alla PMF

4.3.2 Continuous rvs: PDF, CDF

Remark 89. PDF is the equivalent of PMF, CDF the same.

Definition 4.3.3 ((Probability) density function (PDF)). If X is a continuous rv density is a $f: \mathbb{R} \rightarrow \mathbb{R}$, $f_X(x)$ such as, considered $X \in B \subseteq \mathbb{R}$:

$$\mathbb{P}(X \in B) = \int_{x \in B} f_X(x) dx \quad (4.14)$$

Eg, if $a, b \in \mathbb{R}$, $a < b$:

$$\mathbb{P}(X \in [a, b]) = \int_a^b f_X(x) dx \quad (4.15)$$

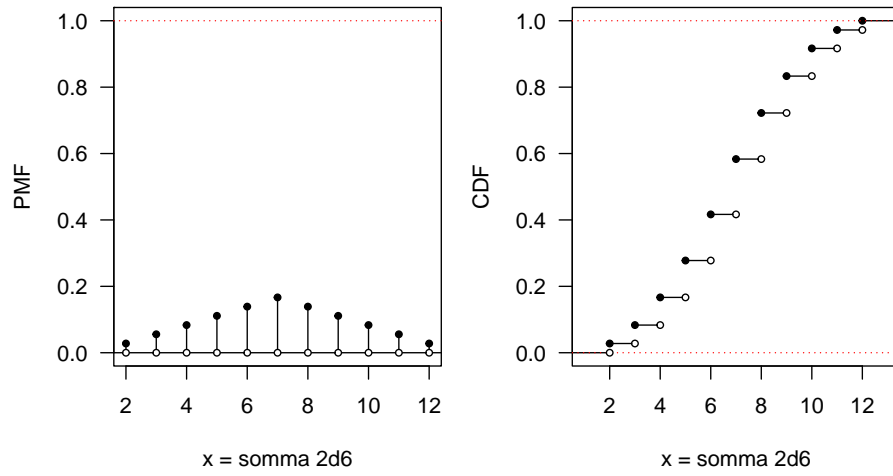


Figura 4.1: Somma del lancio di due d6

Proposition 4.3.3 (Valid PDF). *Must satisfy*

$$f_X(x) \geq 0 \quad (4.16)$$

$$\int_{-\infty}^{\infty} f_X(t) dt = 1 \quad (4.17)$$

Dimostrazione. Il primo criterio è necessario perché la probabilità è non negativa: se $f_X(x_0)$ fosse negativa, allora potremmo integrare su un piccolo intorno di x_0 e ottenere una probabilità negativa.

Il secondo criterio è necessario dato che la X , variabile quantitativa, deve avere un esito che sta in \mathbb{R} . \square

Remark 90. Differently from the discrete case (where PMF can't be more than 1) pdf can be more than 1, as long as integral sums on \mathbb{R} sums up to 1.

Definition 4.3.4 ((Cumulative) distribution function (CDF)). If X is a continuous rv, it's the function $F : \mathbb{R} \rightarrow \mathbb{R}$ defined as:

$$F_X(x) = \mathbb{P}(X \leq x) = \int_{-\infty}^x f_X(t) dt \quad (4.18)$$

Proposition 4.3.4 (Valid CDF). *It must satisfy*

$$x_1 \leq x_2 \implies F_X(x_1) \leq F_X(x_2) \quad (4.19)$$

$$\lim_{x \rightarrow x_0^+} F_X(x) = F_X(x_0) \quad (\text{continuità da destra}) \quad (4.20)$$

$$\lim_{x \rightarrow -\infty} F_X(x) = 0 \quad \lim_{x \rightarrow +\infty} F_X(x) = 1 \quad (4.21)$$

Example 4.3.3 (Esempio crash course). Let's check if

$$F(x) = \begin{cases} 0 & x < 0 \\ 1 - e^{-x} & x \geq 0 \end{cases}$$

is a distribution function. We have

1. $\lim_{x \rightarrow -\infty} F(x) = 0$, $\lim_{x \rightarrow +\infty} F(x) = \lim_{x \rightarrow +\infty} 1 - e^{-x} = 1$, so check for the first
2. for $y > x$ we must show that $F(y) \geq F(x)$ to ensure non decreasing nature. Let's check the sign of $F(y) - F(x)$ (since if $F(y) - F(x) \geq 0$ then $F(y) \geq F(x)$): we have

$$1 - e^{-y} - 1 + e^{-x} = e^{-x} - e^{-y} \stackrel{(1)}{\geq} 0$$

with (1) since $e^{-y} < e^{-x}$ given that $y < x$

3. because $F(x)$ is continuous, it is also right continuous

So yes, $F(x)$ is a CDF ($X \sim \text{Exp}(1)$).

Remark 91 (Probability calculation with CDF). If we know CDF we can evaluate probability of an interval $a \leq X \leq b$, $a, b \in \mathbb{R}$ as follows:

$$\mathbb{P}(a \leq X \leq b) = \mathbb{P}(X \leq b) - \mathbb{P}(X \leq a) = F_X(b) - F_X(a)$$

Remark 92 (Probability of a single value). A differenza delle variabili discrete, nel caso continuo si ha che:

$$\mathbb{P}(X = a) = \int_a^a f_X(x) dx = F_X(a) - F_X(a) = 0$$

Intuitively, if there are infinite outcomes probability of each of them is null.

Remark 93 (Irrelevance of extremes of integration). For the same reason $a, b \in \mathbb{R}$, $a < b$:

$$\mathbb{P}(X \in [a, b]) = \mathbb{P}(X \in (a, b]) = \mathbb{P}(X \in [a, b)) = \mathbb{P}(X \in (a, b)) = \int_a^b f_X(x) dx$$

Example 4.3.4 (Logistic rv). Logistic random variable, plotted in figure 4.2, is defined by:

$$F(x) = \frac{e^x}{1 + e^x}; \quad f(x) = \frac{e^x}{(1 + e^x)^2}$$

4.3.3 Other useful rv functions

4.3.3.1 Support indicator

Remark 94. Nel seguito servirà essere compatti/sicuri sul fatto che, al di fuori del supporto R_X della vc X , la probabilità/densità sia nulla. Per farlo si moltiplicherà la PMF/PDF per la funzione indicatrice applicata al supporto della variabile casuale.

Definition 4.3.5 (Funzione indicatrice del supporto di una vc). Definita come:

$$\mathbb{1}_{R_X}(x) = \begin{cases} 1 & \text{se } x \in R_X \\ 0 & \text{se } x \notin R_X \end{cases}$$

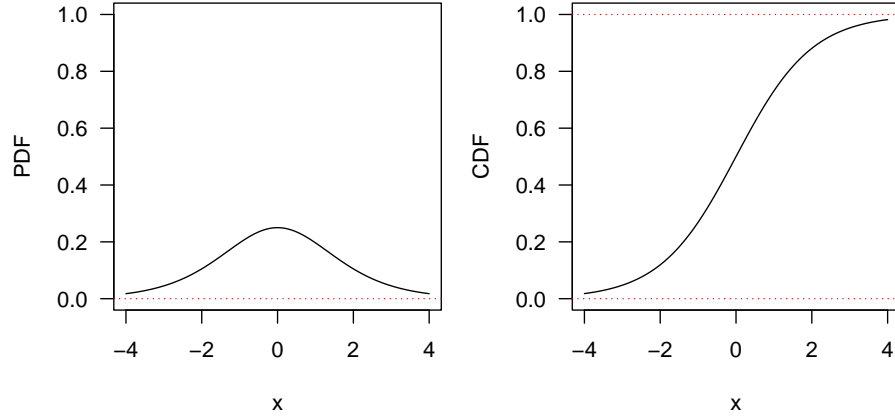


Figura 4.2: Logistic distribution

4.3.3.2 Survival and hazard function

Remark 95. If rv T has non negative support (eg lifetime), then two function are useful (survival for both discrete and continuous rvs, hazard for continuous)

Definition 4.3.6 (Survival function). Given a rv T such as $\mathbb{P}(T \geq 0) = 1$, it's defined as complement to 1 of cumulative distribution function

$$S(t) = \mathbb{P}(T > t) = 1 - \mathbb{P}(T \leq t) = 1 - F_T(t) \quad (4.22)$$

Definition 4.3.7 (Funzione di azzardo (o rischio)). Given a continuous rv T such as $\mathbb{P}(T \geq 0) = 1$, hazard function is defined as

$$H(t) = \frac{f_T(t)}{1 - F_T(t)} = -\frac{d}{dt} \log(1 - F_T(t)) = -\frac{d}{dt} \log(S(t)) \quad (4.23)$$

Remark 96. Hazard function can be interpreted as the probability that T stops at t given that it arrived to t

Remark 97. Relationship between Hazard, survival, density and distribution function can be retrieved by the equation. Eg integrating both members between $-\infty$ and x we have

$$\begin{aligned} H(t) &= -\frac{d}{dt} \log(S(t)) \\ \int_{-\infty}^x H(t) dt &= \int_{-\infty}^x -\frac{d}{dt} \log(S(t)) \\ \int_{-\infty}^x H(t) dt &= -\log(S(t)) \end{aligned}$$

Therefore:

$$\begin{aligned} \log(S(t)) &= -\int_{-\infty}^x H(t) dt \\ S(t) &= \exp\left(-\int_{-\infty}^x H(t) dt\right) \end{aligned} \quad (4.24)$$

While for what concerns $F_T(t)$ e $f_T(t)$ we have:

$$F_T(t) = 1 - \exp\left(-\int_{-\infty}^x H(t) dt\right) \quad (4.25)$$

$$f_T(t) = H(t) \cdot \exp\left(-\int_{-\infty}^x H(t) dt\right) \quad (4.26)$$

Btw, in the lower limit of integration we could have write 0 instead of $-\infty$.

4.4 Transformation

Definition 4.4.1 (Trasform of rv $g(X)$). Considered an experiment with sample space Ω , a random variable X on it and a function $g : \mathbb{R} \rightarrow \mathbb{R}$, then $Y = g(X)$ is the random variable mapping $\omega \rightarrow g(X(\omega))$, $\forall \omega \in \Omega$ and having support $R_{g(X)} = \{g(X(\omega_1)), g(X(\omega_2)), \dots\}$. We're interested in finding the distribution of Y , knowing the distribution of X

Remark 98. The logic behind is that, if X is a rv and g is a “well behaved” function (mainly *strictly increasing* or *strictly decreasing*), then $g(X)$ is also a rv. Our main aim is determine density function of $g(X)$.

Remark 99. More generally let X be a n -variate random vector and $Y = g(X)$ where $g : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is (borel) measurable. Given the distribution of X , we're interested in finding the distribution of Y .

This problem may be easy but also extremely difficult. Here we discuss a couple of simple cases where $m = n = 1$ (in blue in the continuous area).

4.4.1 Discrete rv transform

Remark 100. In the discrete case finding PMF of $g(X)$ is usually easy, the following are some example.

Remark 101. Given a discrete rv X with known PMF, how to get PMF of $Y = g(X)$? If:

- g è *injective*, $X(s_1) \neq X(s_2) \implies g(X(s_1)) \neq g(X(s_2))$, then PMF Y will be the same of X :

$$\mathbb{P}(Y = g(x)) = \mathbb{P}(g(X) = g(x)) = \mathbb{P}(X = x)$$

- otherwise there could be cases where $X(s_1) \neq X(s_2)$ but $\implies g(X(s_1)) = g(X(s_2))$: here we have to sum probability of different x that with g ends in the same y .

The following result is general and is ok for both cases

Proposition 4.4.1 (PMF of $g(X)$). Let X be a discrete rv and $g : \mathbb{R} \rightarrow \mathbb{R}$. Then support of $g(X)$ is the set of y such as that $g(x) = y$ for at least one $x \in R_X$ and PMF of $g(X)$ is

$$\mathbb{P}(g(X) = y) = \sum_{x: g(x)=y} \mathbb{P}(X = x), \quad \forall y \in R_{g(X)} \quad (4.27)$$

X	$\mathbb{P}(X=x)$	$Y=2X$	$\mathbb{P}(Y=y)$	$Z=X^2$	$\mathbb{P}(Z=z)$
-1	0.33	-2	0.33	1	0.66
0	0.33	0	0.33	0	0.33
1	0.33	2	0.33		

Tabella 4.1: PMF of discrete rv transform, an example

Example 4.4.1. In table 4.1 an example with X , $Y = 2X$ ($g(x) = 2 \cdot x$, injective) e $Z = X^2$ ($g(x) = x^2$ not injective).

Remark 102. It's a common error to apply g to the PMF (it could take probability over 1): g have to be applied to domain/support of PMF.

Example 4.4.2 (Transformation of a bernoulli). Let $X \sim \text{Bern}(p)$ and we're interested in $g(X) = e^X$. What is the dist of $g(X)$. We have that

$$X = \begin{cases} 1 & \text{with prob } p \\ 0 & \text{with prob } 1-p \end{cases}, \quad g(X) = \begin{cases} e^1 = e & \text{with prob } p \\ e^0 = 1 & \text{with prob } 1-p \end{cases}$$

Therefore

$$\mathbb{P}(g(X) = e) = \mathbb{P}(X = g^{-1}(e)) = \mathbb{P}(X = 1) = p$$

4.4.2 Continuous rvs transform (linear case)

Definition 4.4.2 (Scale-location transform for continuous rv). Let X be a continuous rv; $Y = \sigma X + \mu$ with $\sigma, \mu \in \mathbb{R}$ is a random variable obtained using a (linear) transform of both position and scale.

Remark 103. Here σ set the scale (if positive spread Y compared to X) while μ the location (if positive moves Y distribution toward right compared to X).

Remark 104. In order to go back to X we standardize Y , aka apply the transformation $X = \frac{Y - \mu}{\sigma}$.

Proposition 4.4.2. Y has the same family of distribution as X .

Dimostrazione. It has been obtained by a linear, injective transformation. \square

Remark 105. If this kind of transformation is applied to a discrete rv we have a distribution no more of the same family, considered that support changes (eg linear transform of a binomial does not give a binomial, defined on support $0, 1, \dots$).

4.4.3 Continuous rvs (monotonic) transform

Proposition 4.4.3 ($g = F$ con F funzione di ripartizione di X). *Let X be a real r.v. with distribution function F (this means that $F(x) = \mathbb{P}(X \leq x)$, $\forall x \in \mathbb{R}$). If F is continuous the $Y = F(X) \sim \text{Unif}(0, 1)$ that is*

$$\mathbb{P}(Y \leq y) = \begin{cases} 0 & y < 0 \\ y & 0 \leq y < 1 \\ 1 & y \geq 1 \end{cases}$$

So $Y = F(X) \sim \text{Unif}(0, 1)$

Dimostrazione. To prove it, for the sake of simplicity assume that F is not only continuous, but also strictly *increasing* (this is not actually needed for thm to hold).

In this case $\forall y \in (0, 1)$ one obtain

$$\mathbb{P}(Y \leq y) = \mathbb{P}(F(X) \leq y) = \mathbb{P}(X \leq F^{-1}(y)) = F(F^{-1}(y)) = y$$

ma questa è proprio la funzione di ripartizione di una uniforme in $(0, 1)$. \square

Proposition 4.4.4. *Let X be absolutely continuous and suppose that $P(X \in I) = 1$ where I is the interval where a function $g : I \rightarrow R$ is defined. Suppose also that g is everywhere differentiable $g' \neq 0$. Then $Y = g(X)$ is still absolutely continuous with density*

$$h_Y(y) = f(g^{-1}(y)) \left| g^{-1'}(y) \right| \cdot \mathbb{1}_{g(I)}(y), \quad \forall y \in g(I)$$

where f denotes the density of X

Example 4.4.3 $(-\log \text{Unif}(0, 1))$. Sia X uniforme in $(0, 1)$, voglio la legge di $Y = -\log(X)$. **NB:** dalla triennale

Basta porre $I = (0, 1)$, $g(x) = -\log(x)$, da cui $g'(x) = -\frac{1}{x} \neq 0 \forall x \in (0, 1)$ e $\mathbb{P}(X \in (0, 1)) = 1$. Quindi posso concludere che Y è assolutamente continua con densità

$$\begin{aligned} h(y) &= f(g^{-1}(y)) \cdot |(g^{-1}(y))'| \cdot \mathbb{1}_{g(I)}(y) \\ &= f(e^{-y}) \cdot |(e^{-y})'| \cdot \mathbb{1}_{(0, +\infty)}(y) \\ &= f(e^{-y}) \cdot e^{-y} \mathbb{1}_{(0, +\infty)}(y) \\ &= 1 \cdot e^{-y} \mathbb{1}_{(0, +\infty)}(y) \end{aligned}$$

dove:

- $y = -\log x \iff x = e^{-y}$ ovvero $g^{-1}(y) = e^{-y}$
- $g(I) = \{g(x) : x \in I\} = \{-\log x : x \in (0, 1)\} = (0, +\infty)$
- $(e^{-y})' = -e^{-y}$ ma poi prendendone il valore assoluto il meno va via
- essendo X uniforme si ha che

$$f(x) = \begin{cases} 1 & 0 < x < 1 \\ 0 & \text{altrimenti} \end{cases}$$

e quindi poiché $e^{-y} \in (0, 1), \forall y \in (0, +\infty)$ si ha $f(e^{-y}) = 1$

In definitiva, poiché $h(y) = \mathbb{1}_{(0, +\infty)}(y) e^{-y}$ è la densità di una esponenziale con $\lambda = 1$ abbiamo dimostrato che

$$-\log \text{Unif}(0, 1) \sim \text{Exp}(1)$$

Example 4.4.4. Per definizione Y ha legge lognormale se $Y > 0$ e $\log Y \sim N(\mu, \sigma^2)$. Al fine di avere una legge esplicita per Y possiamo considerare che deve essere $Y \sim e^X$ con $X \sim N(\mu, \sigma^2)$. Per ottenerla dunque basta applicare il teorema precedente con $X \sim N$ e $g(x) = \exp(x)$. Posso dunque concludere che $Y = g(X)$ ha legge assolutamente continua con densità

$$\begin{aligned} h(y) &= f(g^{-1}(y)) \cdot |(g^{-1}(y))'| \cdot \mathbb{1}_{g(I)}(y) \\ &= f(\log(y)) \cdot |(\log(y))'| \cdot \mathbb{1}_{(0,+\infty)}(y) \\ &= \frac{f(\log(y))}{y} \mathbb{1}_{(0,+\infty)}(y) \\ &= \frac{\exp\left[-\frac{1}{2}\left(\frac{\log y - \mu}{\sigma}\right)^2\right]}{y\sqrt{2\pi\sigma^2}} \mathbb{1}_{(0,+\infty)}(y) \end{aligned}$$

Proposition 4.4.5. If X is a continuous random variable, g a monotonic function (strictly increasing or decreasing), the density function of the random variable $g(X)$, $f_{g(X)}$, is obtained as:

$$f_{g(X)}(x) = f_X(g^{-1}(x)) \cdot \left| \frac{\partial g^{-1}(x)}{\partial x} \right| \quad (4.28)$$

Dimostrazione. For the continuous case we have that, in order to obtain $f_{g(X)}(x)$ we need to differentiate $F_{g(X)}(x)$

$$F_{g(X)}(x) = \mathbb{P}(g(X) \leq x)$$

Now

- if the function g is *decreasing* we have

$$\begin{aligned} F_{g(X)}(x) &= \mathbb{P}(g(X) \leq x) = \mathbb{P}(X \geq g^{-1}(x)) = 1 - \mathbb{P}(X < g^{-1}(x)) \\ &= 1 - F_X(g^{-1}(x)) \end{aligned}$$

- viceversa if g is *increasing*

$$F_{g(X)}(x) = \mathbb{P}(g(X) \leq x) = \mathbb{P}(X \leq g^{-1}(x)) = F_X(g^{-1}(x))$$

In any case after that we have that

$$\begin{aligned} f_{g(X)}(x) &= \frac{\partial}{\partial x} F_{g(X)}(x) = \begin{cases} \frac{\partial(1 - F_X(g^{-1}(x)))}{\partial x} & \text{if increasing} \\ \frac{\partial(F_X(g^{-1}(x)))}{\partial x} & \text{if decreasing} \end{cases} \\ &= \begin{cases} -f_X(g^{-1}(x)) \cdot \frac{\partial}{\partial x} g^{-1}(x) \\ f_X(g^{-1}(x)) \cdot \frac{\partial}{\partial x} g^{-1}(x) \end{cases} \end{aligned}$$

The two cases can be combined in the single formula (not clear how to me for the moment) which is the thorem \square

Example 4.4.5 (Esercizio Berk Tan). Let $X \sim \text{Unif}(0, 1)$ and be $g(x) = e^x$; then what is the pdf of $Y = g(X)$? We have that $g^{-1}(Y) = \log Y$, so

$$\frac{\partial}{\partial y}(g^{-1}(y)) = \frac{1}{y}$$

Applying the formula

$$f_Y(y) = \mathbb{1}_{[0,1]}(\log y) \frac{1}{y}$$

and expressing $\mathbb{1}_{[0,1]}(\log y)$ in terms of y we have

$$\begin{aligned} 0 &\leq \log y \leq 1 \\ 1 &\leq y \leq e \end{aligned}$$

so finally

$$f_Y(y) = \mathbb{1}_{[1,e]}(y) \frac{1}{y} = \begin{cases} \frac{1}{y} & \text{if } y \in [1, e] \\ 0 & \text{elsewhere} \end{cases}$$

Example 4.4.6 (Esame vecchio viroli). Let X have the probability density function given by

$$f_X(x) = \frac{x}{2}$$

with $X \in [0, 2]$. Find the density function of $Y = 6X - 3$.

Qua il dominio diventa palesemente $Y \in [-3, 9]$, per quanto riguarda la funzione si ha che

$$\begin{aligned} f_Y(y) &= \left| \frac{\partial}{\partial y} g^{-1}(y) \right| f_X(g^{-1}(y)) \\ g(X) &= 6X - 3 \quad g^{-1}(Y) = \frac{Y + 3}{6} \\ f_Y(y) &= \frac{1}{6} \left(\frac{Y + 3}{6 \cdot 2} \right) = \frac{1}{6} \left(\frac{Y + 3}{12} \right) \end{aligned}$$

the answer is $f_Y(y) = \frac{3+y}{12} \frac{1}{6}$.

Si può verificare che $\int_{-3}^9 f_Y(y) = 1$ mediante sympy. Qui non c'è il problema di resprimere le variabili indicatrici (perché non è una uniforme 0,1 e la densità non ne fa uso).

Example 4.4.7 (Assignment 1 Viroli, Exercise 2). Let $X \sim \text{Unif}(0, 1)$. Find the PDF of $X^{1/\alpha}$ with $\alpha > 0$.

Let $X \sim \text{Unif}(0, 1)$ and $Y = X^{\frac{1}{\alpha}}$, with $\alpha > 0$. Let's obtain $f_Y(y)$ by applying:

$$f_Y(y) = f_X(g^{-1}(y)) \cdot \left| \frac{\partial g^{-1}(y)}{\partial y} \right| \quad (4.29)$$

Being $X \sim \text{Unif}(0, 1)$ we have that $f_X(x) = \mathbb{1}_{[0,1]}(x)$. Given the transformation $y = x^{1/\alpha}$, its inverse is

$$y = x^{1/\alpha} \iff y^\alpha = x$$

so $g^{-1}(Y) = Y^\alpha$; doing the derivative with respect to y we obtain:

$$\frac{\partial}{\partial y} g^{-1}(y) = \alpha y^{\alpha-1}$$

so putting things together:

$$f_Y(y) = f_X(g^{-1}(y)) \cdot \left| \frac{\partial g^{-1}(y)}{\partial y} \right| = \mathbb{1}_{[0,1]}(y^\alpha) \cdot \alpha y^{\alpha-1}$$

Now we need to express the indicator $\mathbb{1}_{[0,1]}(y^\alpha)$ in terms of y , therefore:

$$\begin{aligned} 0 &\leq y^\alpha \leq 1 \\ 0 &\leq y \leq 1 \end{aligned}$$

Finally:

$$f_Y(y) = \mathbb{1}_{[0,1]}(y) \cdot \alpha y^{\alpha-1} = \begin{cases} \alpha y^{\alpha-1} & \text{if } y \in [0, 1] \\ 0 & \text{elsewhere} \end{cases}$$

If $\alpha = 1$, as expected

$$f_Y(y) = \begin{cases} 1 & \text{if } y \in [0, 1] \\ 0 & \text{elsewhere} \end{cases} = \mathbb{1}_{[0,1]}(y) \implies Y \sim \text{Unif}(0, 1)$$

Example 4.4.8 (Esercizio virol). If $X \sim \text{Unif}(0, 1)$ and $Y = -2 \log X$, show that $Y \sim \chi^2_2$. We apply 4.28 and compare with χ^2_n one.

We have the transformation $y = -2 \log x$ so to obtain the inverse

$$-\frac{1}{2}y = \log x \iff x = e^{-\frac{1}{2}y}$$

therefore $g^{-1}(Y) = \exp(-\frac{Y}{2})$. We have, being X a uniform on $0,1$, that $f_X(x) = 1 \cdot \mathbb{1}_{[0,1]}(x)$. Now

$$\frac{\partial}{\partial y} g^{-1}(y) = -\frac{1}{2} e^{-y/2}$$

So applying the formula we arrive at

$$f_Y(y) = \mathbb{1}_{[0,1]}(e^{-y/2}) \cdot \frac{1}{2} e^{-y/2}$$

Now we need to express $\mathbb{1}_{[0,1]}(e^{-y/2})$ in terms of y . The domain of y so

$$\begin{aligned} 0 &\leq e^{-y/2} \leq 1 \\ -\infty &< -y/2 \leq 0 \\ 0 &< y \leq +\infty \end{aligned}$$

Finally

$$f_Y(y) = \mathbb{1}_{[0,+\infty)}(y) \cdot \frac{1}{2} e^{-y/2} = \begin{cases} \frac{1}{2} e^{-y/2} & \text{if } y \in [0, +\infty) \\ 0 & \text{elsewhere} \end{cases}$$

which is a χ^2 with 2 degrees of freedom.

4.5 Independence

4.5.1 Independence

Remark 106 (Notation). We can write intersections of events as follows

$$\begin{aligned}\mathbb{P}(X \in A, Y \in B) &= \mathbb{P}(X \in A \cap Y \in B) \\ \mathbb{P}(X \leq x, Y \leq y) &= \mathbb{P}(X \leq x \cap Y \leq y)\end{aligned}$$

Remark 107. The concept of independence for random variables is similar to events independence.

Definition 4.5.1 (RVs independence (general case)). Given *any* collection (finite, countable, non countable) of random variables $\mathcal{V} = \{X_1, X_2, \dots\}$, the elements of \mathcal{V} are said to be independent if, for any *finite* subset of events $\mathcal{X} \subset \mathcal{V}$

$$\begin{aligned}\mathbb{P}(X_j \in B_j, \dots, X_k \in B_k) &= \mathbb{P}(X_j \in B_j) \cdot \dots \cdot \mathbb{P}(X_k \in B_k) \\ X_j, \dots, X_k \in \mathcal{X} \quad \forall B_j, \dots, B_k \in \mathcal{B}\end{aligned}$$

or equivalently

$$\begin{aligned}\mathbb{P}(X_j \leq x_j, \dots, X_k \leq x_k) &= \mathbb{P}(X_j \leq x_j) \cdot \dots \cdot \mathbb{P}(X_k \leq x_k) \\ X_j, \dots, X_k \in \mathcal{X}, \quad \forall x_j, \dots, x_k \in \mathbb{R}\end{aligned} \quad (4.30)$$

Example 4.5.1 (Independence of two RVs $X \perp\!\!\!\perp Y$). Two rvs X, Y are independent, and we write $X \perp\!\!\!\perp Y$, if

$$\mathbb{P}(X \leq x, Y \leq y) = \mathbb{P}(X \leq x) \cdot \mathbb{P}(Y \leq y), \quad \forall x, y \in \mathbb{R} \quad (4.31)$$

Remark 108. In the discrete case 4.31 is equivalent to

$$\mathbb{P}(X = x, Y = y) = \mathbb{P}(X = x) \cdot \mathbb{P}(Y = y), \quad \forall x, y \in \mathbb{R}$$

Example 4.5.2. Let be X the result of first dice thrown and Y the second; sum and difference of results random variables $X + Y$, $X - Y$ are not independent considered that:

$$\begin{aligned}\mathbb{P}(X + Y = 12, X - Y = 1) &= 0 \\ \mathbb{P}(X + Y = 12) \cdot \mathbb{P}(X - Y = 1) &= \frac{1}{6} \cdot \frac{5}{6}\end{aligned}$$

This does make sense: knowing that the sum is 12, tells that their difference must be 0 so the two rv gives information of each other

Proposition 4.5.1. If X_1, \dots, X_n are independent, then they are pairwise, 3-wise, \dots $(n-1)$ -wise independent. Viceversa implication does not hold.

Dimostrazione. If X_1, \dots, X_n are independent si ha (considerando a titolo di esempio la coppia X_1, X_2) che

$$\mathbb{P}(X_1 \leq x_1, X_2 \leq x_2) = \mathbb{P}(X_1 \leq x_1) \cdot \mathbb{P}(X_2 \leq x_2)$$

Per vedere perché sia così basta far tendere a $+\infty$ gli x_3, \dots, x_n in maniera tale che a sinistra dell'uguale, nella definizione 4.30, entro parentesi si abbiano eventi certi e a destra dell'uguale si moltiplichino per 1.

The example of why contrary implication does not hold can be done via counterexample. \square

NB: La definizione generale di sopra vale anche qui perché in un set finito chiediamo un subset sempre finito di 2 variabili poi lui credo usi $\mathcal{X} \subset \mathcal{V}$ per poter intendere anche $\mathcal{X} \subseteq \mathcal{V}$

Example 4.5.3. Example of three variables which are pairwise independent but not independent. Let X, Y be iid with

$$\mathbb{P}(X = 1) = \mathbb{P}(X = -1) = 1/2$$

and $Z = XY$ so that

$$\begin{aligned} \mathbb{P}(Z = 1) &= \mathbb{P}(X = Y) = \mathbb{P}(X = 1, Y = 1) + \mathbb{P}(X = -1, Y = -1) \\ &= \mathbb{P}(X = 1) \cdot \mathbb{P}(Y = 1) + \mathbb{P}(X = -1) \cdot \mathbb{P}(Y = -1) \\ &= \frac{1}{2} \cdot \frac{1}{2} + \frac{1}{2} \cdot \frac{1}{2} = \frac{1}{2} \end{aligned}$$

and thus $\mathbb{P}(Z = -1) = 1/2$.

The set $\{X, Y, Z\}$ is not independent since, for example

$$\mathbb{P}(X = 1, Y = -1, Z = 1) = \mathbb{P}(\emptyset) = 0 \neq \mathbb{P}(X = 1) \cdot \mathbb{P}(Y = -1) \cdot \mathbb{P}(Z = 1) = \frac{1}{8}$$

However the three random variables are pairwise independent since

$$\begin{aligned} \mathbb{P}(X = 1, Y = 1) &= \mathbb{P}(X = 1) \cdot \mathbb{P}(Y = 1) = \frac{1}{2} \frac{1}{2} = \frac{1}{4} \\ \mathbb{P}(X = 1, Y = -1) &= \mathbb{P}(X = 1) \cdot \mathbb{P}(Y = -1) = \frac{1}{2} \frac{1}{2} = \frac{1}{4} \\ \mathbb{P}(X = 1, Z = 1) &= \mathbb{P}(X = 1) \cdot \mathbb{P}(Z = 1) = \frac{1}{2} \frac{1}{2} = \frac{1}{4} \\ \mathbb{P}(X = 1, Z = -1) &= \mathbb{P}(X = 1) \cdot \mathbb{P}(Z = -1) = \frac{1}{2} \frac{1}{2} = \frac{1}{4} \\ &\dots \end{aligned}$$

and in general one obtains

$$\mathbb{P}(X = a, Z = b) = \mathbb{P}(X = a) \mathbb{P}(Z = b), \quad \forall a, b \in \{-1, 1\}$$

Proposition 4.5.2 (Transform of independent rv). *If X and Y are independent, then any transformation of X and Y are independent as well.*

Dimostrazione. Not shown. □

4.5.2 IID RVs

Remark 109. A very important case is IID variables; this assumption is involved in the *law of large number* and *central limit theorem*.

Definition 4.5.2 (i.i.d. rvs). Random variables in the set $\mathcal{V} = \{X_1, X_2, \dots\}$ are *independent* and *identically* distributed if

- the elements of \mathcal{V} are independent
- $X_i \sim X_j$, for all $X_i, X_j \in \mathcal{V}$ (have the same distribution function).

Important remark 31 (Notation). If the elements of $\mathcal{X} = \{X_1, X_2, \dots\}$ are iid, to communicate the common distribution of the X_i it suffices to write $X_i \sim \nu$

4.5.3 Conditional independence

Definition 4.5.3 (Conditional independence). X and Y are conditional independent given Z if $\forall x, y \in \mathbb{R}$ and $\forall z \in R_Z$ it is:

$$\mathbb{P}(X \leq x, Y \leq y | Z = z) = \mathbb{P}(X \leq x | Z = z) \cdot \mathbb{P}(Y \leq y | Z = z) \quad (4.32)$$

Remark 110. For discrete rvs, an equivalent definition based on the mass function is

$$\mathbb{P}(X = x, Y = y | Z = z) = \mathbb{P}(X = x | Z = z) \cdot \mathbb{P}(Y = y | Z = z) \quad (4.33)$$

Proposition 4.5.3. *Rvs independence does not imply conditional independence and viceversa.*

Dimostrazione. By counterexamples, see Blitzstein pag 121. □

4.6 Moments

Remark 111. Distribution functions are the unifying concepts for continuous and discrete rvs; furthermore knowing F_X is to know the entire probabilistic structure of the rv.

In order to compare different rvs, however, often synthetic indicators are needed and these are the moments. Different indicators are available for different features of the distribution.

Definition 4.6.1 (Moment of a rv). A statistic of this kind, if it exists

$$\begin{aligned} \sum_{x \in R_X}^{\infty} g(x) \cdot p_X(x) & \quad \text{if } X \text{ is discrete} \\ \int_{-\infty}^{+\infty} g(x) \cdot f_X(x) \, dx & \quad \text{if } X \text{ is abs. continuous} \end{aligned}$$

Different g functions defines different moments

Important remark 32 (Important moments). These are expected value, variance, asymmetry and kurtosis; all can be seen as a specialized (for g) version of the equations above.

4.6.1 Expected value

Remark 112 (Expected value existence check). Let X be a real r.v.; we aim to define its expectation. Before doing this however, it should be noted that such expectation (which involves series or integrals) may fail to exist (not finite).

To define the expectation of X (whether it is discrete or continuous), one should previously evaluate the expectation of $|X|$, that is $\mathbb{E}[|X|]$; this can be done through the formula

$$\mathbb{E}[|X|] = \int_0^{+\infty} \mathbb{P}(|X| > t) \, dt$$

Incidentally if X is absolutely continuous, the above integral can be written as

$$\mathbb{E}[|X|] = \int_0^{+\infty} \mathbb{P}(|X| > t) dt = \int_{-\infty}^{+\infty} |x| f(x) dx$$

where f is the density of X . Now there are two situations:

$$\begin{cases} \mathbb{E}[|X|] = +\infty & \implies \text{we stop: expectation of } X \text{ does not exist} \\ \mathbb{E}[|X|] < \infty & \implies \text{expectation of } X \text{ exists and may be evaluated with following formulas} \end{cases}$$

Definition 4.6.2 (Expected value). If $\mathbb{E}[|X|] < +\infty$ the expectation of X , denoted by $\mathbb{E}[X]$ or μ , gives a probability weighted mean of X and can be evaluated by

$$\mathbb{E}[X] = \begin{cases} \sum_{x \in R_X} x \cdot \mathbb{P}(X = x) & \text{if } X \text{ is discrete} \\ \int_{-\infty}^{+\infty} x \cdot f_X(x) dx & \text{if } X \text{ is abs. continuous} \\ \int_0^{+\infty} \mathbb{P}(X > t) dt & \text{if } X \geq 0 \end{cases}$$

Remark 113. The cases above don't cover all the possible cases (eg there are other formulas if X is not discrete, absolutely continuous or non negative) but are more than enough for us

Example 4.6.1 (Single dice). Let X be the result of a single fair dice with $p_X(1) = \dots = p_X(6) = 1/6$:

$$\mathbb{E}[X] = 1 \cdot \frac{1}{6} + 2 \cdot \frac{1}{6} + 3 \cdot \frac{1}{6} + 4 \cdot \frac{1}{6} + 5 \cdot \frac{1}{6} + 6 \cdot \frac{1}{6} = \frac{7}{2}$$

Example 4.6.2. For the Cauchy random variable, the expected value does not exist. If $X \sim \text{Ca}(\theta)$, with $\theta > 0$, X is absolutely continuous with support \mathbb{R} and density

$$f(x) = \frac{\theta}{\pi} \frac{1}{\theta^2 + x^2} \cdot \mathbb{1}_{R_X}(x)$$

In order to check it, we start evaluating the test for expected value existence (suppose $\theta = 1$):

$$\begin{aligned} \mathbb{E}[|X|] &= \int_{-\infty}^{+\infty} |x| \cdot \frac{1}{\pi} \frac{1}{1+x^2} \stackrel{(1)}{=} 2 \int_0^{+\infty} x \cdot \frac{1}{\pi} \frac{1}{1+x^2} \\ &= 2 \cdot \frac{1}{\pi} \int_0^{+\infty} \frac{x}{1+x^2} dx \stackrel{(2)}{=} +\infty \end{aligned}$$

where in:

- (1) because it's an even function (symmetry with respect to y axis) so we can double the integral on the positive part (taking x out of absolute value);
- (2) if we want to check very well, integrating by parts we have:

NB: rigo non ha fatto l'integrazione

$$\int \frac{x}{1+x^2} dx = \frac{1}{2} \int \frac{2x}{1+x^2} dx = \frac{1}{2} \log(1+x^2) + c$$

Therefore

$$\mathbb{E}[|X|] = \frac{2}{\pi} \left(\left[\frac{1}{2} \log(1+x^2) \right]_0^{+\infty} \right) = \frac{2}{\pi} (+\infty - 0) = +\infty$$

Therefore the expected value does not exist.

Remark 114. Generalizing a bit, expectation is the *first* moment of a random variable X .

Definition 4.6.3 (Moment of order r (r -th moment) of X). Adopting as g the r -power of X in the definition 4.6.1

$$\mu_r = \mathbb{E}[X^r] = \begin{cases} \sum_{x \in R_X} x^r \cdot \mathbb{P}(X = x) & \text{if } X \text{ is discrete} \\ \int_{-\infty}^{+\infty} x^r \cdot f_X(x) dx & \text{if } X \text{ is abs. continue} \end{cases} \quad (4.34)$$

Definition 4.6.4 (Moment of order r existence). In general moment of order r for X exists (or X has moment of order r) if $\mathbb{E}[|X|^r] < +\infty$.

Remark 115. A useful results is the following.

Theorem 4.6.1. If $\mathbb{E}[|X|^r] < +\infty$ for some $r > 0$, then all the moments of order $q \leq r$ exists/are finite as well:

$$\mathbb{E}[|X|^q] < +\infty, \quad \forall q \in (0, r]$$

Remark 116. From now on, all the involved rv are assumed to have the mean. The following properties are very useful since they hold for any rv (regardless the type). The only needed assumption is that the involved rv has the mean.

Proposition 4.6.2 (Main properties of the operator $\mathbb{E}[\cdot]$). *We have*

1. $\mathbb{E}[aX + bY] = a\mathbb{E}[X] + b\mathbb{E}[Y]$ (*linearity*)
2. if $c \in \mathbb{R}$, $\mathbb{E}[c] = c$ (*expval of constant/dirac*)
3. $X \geq 0$ a.s. $\mathbb{E}[X] \geq 0$ (*positivity, just \geq*)
4. if $X \geq 0$ and $\mathbb{P}(X > 0) > 0$ then $\mathbb{E}[X] > 0$ (*strict positivity*)

Proposition 4.6.3 (Expected value properties (old non Rigo version)).

$$\mathbb{E}[aX + b] = a\mathbb{E}[X] + b \quad (4.35)$$

$$\mathbb{E}[X + Y] = \mathbb{E}[X] + \mathbb{E}[Y] \quad (4.36)$$

$$X \geq 0 \implies \mathbb{E}[X] \geq 0 \quad (4.37)$$

$$X \geq 0, \mathbb{P}(X > 0) > 0 \implies \mathbb{E}[X] > 0 \quad (4.38)$$

$$\mathbb{E}[g(X)] = \sum_i g(x_i) \cdot p_X(x_i) \quad (4.39)$$

$$X \perp\!\!\!\perp Y \implies \mathbb{E}[XY] = \mathbb{E}[X]\mathbb{E}[Y] \quad (4.40)$$

$$\min(X) \leq \mathbb{E}[X] \leq \max(X) \quad (4.41)$$

$$\mathbb{E}[X - \mathbb{E}[X]] = 0 \quad (4.42)$$

$$\text{minimizes } \mathbb{E}[(X - \mathbb{E}[X])^2] \quad (4.43)$$

Remark 117. Congiuntamente alle 4.35 e 4.36 ci si riferisce come linearità del valore atteso, che torna spesso comodo per il calcolo soprattutto se si riesce a scrivere una vc come somma di due o più vc. La linearità è un mero fatto algebrico e di bello c'è che, ad esempio per 4.36, non è necessaria l'indipendenza tra X e Y affinché valga.

TODO: da chiarire sta
nota di colore

Important remark 33. If $f : \mathbb{R} \rightarrow \mathbb{R}$ is a measurable function, to evaluate the expectation of $f(X)$, that is $E(f(X))$, we can repeat the previous properties with $f(X)$ instead of X .

Dimostrazione. Mostriamo con riferimento alle variabili discrete. Per la 4.35

$$\begin{aligned}\mathbb{E}[aX + b] &= \sum_{x \in R_X} (ax + b) \cdot \mathbb{P}(aX + b = ax + b) = \sum_{x \in R_X} (ax + b) \cdot \mathbb{P}(X = x) \\ &= \sum_{x \in R_X} ax \cdot \mathbb{P}(X = x) + \sum_{x \in R_X} b \cdot \mathbb{P}(X = x) \\ &= a \sum_{x \in R_X} x \cdot \mathbb{P}(X = x) + b \underbrace{\sum_{x \in R_X} \mathbb{P}(X = x)}_1 \\ &= a \mathbb{E}[X] + b\end{aligned}$$

Viceversa nel caso continuo

$$\mathbb{E}[aX + b] = \int_{-\infty}^{+\infty} (ax+b)f(x) dx = a \int_{-\infty}^{+\infty} xf(x) dx + b \underbrace{\int_{-\infty}^{+\infty} f(x) dx}_{=1} = a \mathbb{E}[X] + b$$

Per 4.36 facendo un passo indietro, possiamo scrivere un generico valore atteso facendo riferimento all'evento $\omega \in \Omega$ e applicando la funzione X ad esso, al fine di ottenere x :

$$\mathbb{E}[X] = \sum_{x \in R_X} x \cdot \mathbb{P}(X = x) = \sum_{\omega \in \Omega} X(\omega) \cdot \mathbb{P}(\{\omega\})$$

Da questa possiamo generalizzare alla somma di due funzioni:

$$\begin{aligned}\mathbb{E}[X + Y] &= \sum_{\omega \in \Omega} (X + Y)(\omega) \cdot \mathbb{P}(\{\omega\}) = \sum_{\omega} (X(\omega) + Y(\omega)) \cdot \mathbb{P}(\{\omega\}) \\ &= \sum_{\omega \in \text{samplesp}} X(\omega) \cdot \mathbb{P}(\{\omega\}) + \sum_{\omega \in \Omega} Y(\omega) \cdot \mathbb{P}(\{\omega\}) \\ &= \mathbb{E}[X] + \mathbb{E}[Y]\end{aligned}$$

Per il valore atteso della trasformazione g , 4.39, sfruttiamo la stessa tecnica facendo un passo indietro (rispetto all'applicazione della funzione X agli eventi dello spazio campionario): sia $\omega \in \Omega$ un outcome dello spazio campionario e X la vc considerata. Come detto possiamo scrivere il valore atteso $\mathbb{E}[X]$ come prodotto del risultato di X per la probabilità che si verifichi quell'evento:

$$\mathbb{E}[X] = \sum_{\omega \in \Omega} X(\omega) \mathbb{P}(\{\omega\})$$

L'applicazione della trasformazione g porta il valore atteso $\mathbb{E}[g(X)]$:

$$\begin{aligned}\mathbb{E}[g(X)] &= \sum_{\omega \in \Omega} g(X(\omega)) \cdot \mathbb{P}(\{\omega\}) \\ &\stackrel{(1)}{=} \sum_{x \in R_X} \sum_{\omega: X(\omega)=x} g(X(\omega)) \mathbb{P}(\{\omega\}) \\ &= \sum_{x \in R_X} g(x) \sum_{\omega: X(\omega)=x} \mathbb{P}(\{\omega\}) \\ &= \sum_{x \in R_X} g(x) \cdot \mathbb{P}(X = x) \\ &= \sum_{x \in R_X} g(x) \cdot p_X(x)\end{aligned}$$

dove in (1) semplicemente raggruppiamo per i diversi ω che attraverso X forniscono lo stesso x .

Per 4.40 (mostrando il caso delle discrete) se $X \perp\!\!\!\perp Y$, allora $\mathbb{P}(X = x, Y = y) = \mathbb{P}(X = x) \cdot \mathbb{P}(Y = y)$, da questo

$$\begin{aligned}\mathbb{E}[XY] &= \sum_{x \in R_X} \sum_{y \in R_Y} x \cdot y \cdot \mathbb{P}(X = x, Y = y) = \sum_{x \in R_X} \sum_{y \in R_Y} x \cdot y \cdot \mathbb{P}(X = x) \mathbb{P}(Y = y) \\ &= \sum_{x \in R_X} x \cdot \mathbb{P}(X = x) \sum_{y \in R_Y} y \cdot \mathbb{P}(Y = y) = \mathbb{E}[X] \cdot \mathbb{E}[Y]\end{aligned}$$

La 4.41 è ovvia essendo $\mathbb{E}[X]$ una media pesata da probabilità dei valori assunti da X ; l'uguaglianza vale in caso di variabili degeneri.

La 4.42 è una applicazione della linearità

$$\mathbb{E}[X - \mathbb{E}[X]] = \mathbb{E}[X] - \mathbb{E}[\mathbb{E}[X]] = \mathbb{E}[X] - \mathbb{E}[X] = 0$$

□

Example 4.6.3 (Valore atteso di trasformazione). Supponiamo che X sia una vc che assuma i valori $-1, 0, 1$ con probabilità pari a $\mathbb{P}(x = -1) = 0.2$, $\mathbb{P}(x = 0) = 0.5$, $\mathbb{P}(x = 1) = 0.3$. Calcoliamo $\mathbb{E}[X^2]$ applicando prima la trasformazione e poi moltiplicando per la probabilità:

$$\mathbb{E}[X^2] = (-1)^2(0.2) + 0^2 \cdot (0.5) + 1^2(0.3) = 0.5$$

Proposition 4.6.4 (Valore atteso di funzioni non lineari di vc). *In generale non vale $\mathbb{E}[g(X)] = g(\mathbb{E}[X])$ per una qualsiasi funzione g .*

Example 4.6.4. Sia X il lancio di un dado: calcoliamo $\exp(\mathbb{E}[X])$ e $\mathbb{E}[\exp X]$; ricordando che $\mathbb{E}[X] = 7/2$ si ha

$$\begin{aligned}g(\mathbb{E}[X]) &= \exp(7/2) \approx 33.12 \\ \mathbb{E}[g(X)] &= e^1 \cdot \frac{1}{6} + \dots + e^6 \cdot \frac{1}{6} \approx 106.1\end{aligned}$$

Considerando invece una trasformazione lineare $g(x) = 2x + 1$ i due risultati coincidono, come in mostrato 4.35. Si ha:

$$\begin{aligned}g(\mathbb{E}[X]) &= 2 \cdot \frac{7}{2} + 1 = 8 \\ \mathbb{E}[g(X)] &= 3 \frac{1}{6} + 5 \frac{1}{6} + 7 \frac{1}{6} + 9 \frac{1}{6} + 11 \frac{1}{6} + 13 \frac{1}{6} = 8\end{aligned}$$

4.6.2 Variance

Definition 4.6.5 (Variance). If $\mathbb{E}[|X|^2] = \mathbb{E}[X^2] < +\infty$ we can define the variance of X as

$$\bar{\mu}_2 = \text{Var}[X] = \sigma^2 = \mathbb{E}[(X - \mathbb{E}[X])^2] \quad (4.44)$$

measure dispersion of the rv around its mean value.

Proposition 4.6.5 (Formula to use for evaluation).

$$\text{Var}[X] = \mathbb{E}[X^2] - (\mathbb{E}[X])^2 \quad (4.45)$$

Remark 118. In the computation formula 4.45, its easier to see that to have a variance it must be $\mathbb{E}[|X|^2] \mathbb{E}[X^2] < +\infty$

Dimostrazione. We expand $(X - \mathbb{E}[X])^2$ and used expected value linearity:

$$\begin{aligned} \text{Var}[X] &= \mathbb{E}[(X - \mathbb{E}[X])^2] = \mathbb{E}[X^2 - 2\mathbb{E}[X]X + (\mathbb{E}[X])^2] \\ &= \mathbb{E}[X^2] - 2\mathbb{E}[X]\mathbb{E}[X] + (\mathbb{E}[X])^2 \\ &= \mathbb{E}[X^2] - (\mathbb{E}[X])^2 \end{aligned}$$

□

Example 4.6.5 (Dice variance). If X is result of a dice throw, previously we computed $\mathbb{E}[X] = 7/2$; furthermore we have

$$\mathbb{E}[X^2] = 1^2\left(\frac{1}{6}\right) + 2^2\left(\frac{1}{6}\right) + 3^2\left(\frac{1}{6}\right) + 4^2\left(\frac{1}{6}\right) + 5^2\left(\frac{1}{6}\right) + 6^2\left(\frac{1}{6}\right) = \left(\frac{1}{6}\right)(91)$$

Therefore

$$\text{Var}[X] = \frac{91}{6} - \left(\frac{7}{2}\right)^2 = \frac{35}{12}$$

Remark 119 (Interpretation). We have that:

- X can be regarded as the outcome of a numerical experiment
- $\mathbb{E}[X]$ our best prediction of X (before making the experiment)
- $X - \mathbb{E}[X]$ can be seen as the error
- the variance is $\mathbb{E}[\text{error}^2]$ if we adopt the best prediction possible (which minimizes error squared) for outcome of our experiment (which is $\mathbb{E}[X]$). Infact, in general, if we predict X by a real number t the error becomes $X - t$. Defining the function

$$e(t) = \mathbb{E}[\text{error}^2] = \mathbb{E}[(X - t)^2]$$

we aim to minimize e . To this end, we note that

$$\begin{aligned} e(t) &= \mathbb{E}[(X - t)^2] = \mathbb{E}[(X - \mathbb{E}[X] + (\mathbb{E}[X] - t))^2] \\ &= \mathbb{E}[(X - \mathbb{E}[X])^2] + (\mathbb{E}[X] - t)^2 + 2(\mathbb{E}[X] - t)\underbrace{\mathbb{E}[X - \mathbb{E}[X]]}_{=0} \\ &= \text{Var}[X] + (t - \mathbb{E}[X])^2 \end{aligned}$$

Hence e attains its minimum at the point $t = \mathbb{E}[X]$ and $\text{Var}[X]$ is our estimate of error in prediction.

Proposition 4.6.6. *Si ha che*

$$X \text{ è degenere} \iff \text{Var}[X] = 0$$

Dimostrazione. La si può provare utilizzando una disuguaglianza (come faremo a tempo debito) o con le proprietà del valore atteso:

NB: dimostrazione dalla triennale che non fa uso di Jensen

- supposing $X = a$ almost surely ($\mathbb{P}(X = a) = 1$), then $\mathbb{E}[X] = a$ and also $\mathbb{E}[X^2] = a^2$, thus

$$\text{Var}[X] = \mathbb{E}[X^2] - (\mathbb{E}[X])^2 = a^2 - a^2 = 0$$

- otherwise suppose

$$0 = \text{Var}[X] = \mathbb{E}[(X - \mathbb{E}[X])^2]$$

Ma $(X - \mathbb{E}[X])^2 \geq 0$ e quindi per la proprietà 4 del valore atteso (strict positivity), se fosse $\mathbb{P}(X \neq \mathbb{E}[X]) > 0$ si avrebbe $\mathbb{P}((X - \mathbb{E}[X])^2 > 0) > 0$ e quindi $\text{Var}[X] = \mathbb{E}[(X - \mathbb{E}[X])^2] > 0$.

Quindi deve essere $\mathbb{P}(X \neq \mathbb{E}[X]) = 0$ ossia $\mathbb{P}(X = \mathbb{E}[X]) = 1$ ovvero X degenere con $a = \mathbb{E}[X]$

□

Remark 120. Generalizing a bit, variance is the *second* moment of a random variable with respect to its mean.

Definition 4.6.6 (r -th moments of X with respect to mean). In the definition 4.6.1 is obtained by adopting as g the r -power of difference between X and its expected value, $g = (x - \mathbb{E}[X])^r$:

$$\bar{\mu}_r = \mathbb{E}[(X - \mathbb{E}[X])^r] = \begin{cases} \sum (x_i - \mathbb{E}[X])^r \cdot p_X(x_i) & \text{se } X \text{ è discreta} \\ \int_{-\infty}^{+\infty} (x - \mathbb{E}[X])^r \cdot f_X(x) \, dx & \text{se } X \text{ è continua} \end{cases} \quad (4.46)$$

Remark 121. Since $\bar{\mu}_0 = 1, \bar{\mu}_1 = 0$, these moments become interesting starting from $r = 2$.

Proposition 4.6.7 (Properties of variance). *Given $a, b, c \in \mathbb{R}$:*

$$\text{Var}[X] \geq 0 \quad (4.47)$$

$$\text{Var}[X] = 0 \iff \mathbb{P}(X = c) = 1 \quad (4.48)$$

$$\text{Var}[aX + b] = a^2 \text{Var}[X] \quad (4.49)$$

$$X \perp\!\!\!\perp Y \implies \text{Var}[X + Y] = \text{Var}[X] + \text{Var}[Y] \quad (4.50)$$

Dimostrazione. Per la 4.47, la varianza è il valore atteso della vc nonnegativa $(X - \mathbb{E}[X])^2$, motivo per cui è non negativa date le proprietà del valore atteso. Per 4.48 se $\mathbb{P}(X = c) = 1$ per qualche costante c allora $\mathbb{E}[X] = c$ e $\mathbb{E}[X^2] = c^2$, pertanto $\text{Var}[X] = 0$; viceversa se $\text{Var}[X] = 0$ allora $\mathbb{E}[(X - \mathbb{E}[X])^2] = 0$ che mostra che $(X - \mathbb{E}[X])^2 = 0$ ha probabilità 1, che a sua volta mostra che X è uguale alla sua media con probabilità 1.

Per la 4.49 e per la linearità del valore atteso si ha:

$$\begin{aligned} \text{Var}[aX + b] &= \mathbb{E}[(aX + b - (a\mathbb{E}[X] + b))^2] \\ &= \mathbb{E}[(aX + b - a\mathbb{E}[X] - b)^2] \\ &= \mathbb{E}[(aX - a\mathbb{E}[X])^2] \\ &= \mathbb{E}[a^2(X - \mathbb{E}[X])^2] \\ &= a^2 \mathbb{E}[(X - \mathbb{E}[X])^2] \\ &= a^2 \text{Var}[X] \end{aligned}$$

La 4.50 verrà dimostrata/generalizzata in seguito, per ora verifichiamola:

$$\begin{aligned} \text{Var}[X + Y] &= \mathbb{E}[(X + Y)^2] - (\mathbb{E}[X + Y])^2 = \mathbb{E}[X^2 + 2XY + Y^2] - (\mathbb{E}[X] + \mathbb{E}[Y])^2 \\ &\stackrel{(1)}{=} \mathbb{E}[X^2] + 2\mathbb{E}[X]\mathbb{E}[Y] + \mathbb{E}[Y^2] - \mathbb{E}[X]^2 - 2\mathbb{E}[X]\mathbb{E}[Y] - \mathbb{E}[Y]^2 \\ &= \text{Var}[X] + \text{Var}[Y] \end{aligned}$$

where in (1) we used that if $X \perp\!\!\!\perp Y$ we have $\mathbb{E}[XY] = \mathbb{E}[X]\mathbb{E}[Y]$. \square

Remark 122 (Variance is nonlinear). Differently from expected value a is squared and b omitted, therefore variance of sum of different random variable could be different from sum of their variance.

Definition 4.6.7 (Standard deviation).

$$\sigma = \sigma_X = \sqrt{\text{Var}[X]} \quad (4.51)$$

4.6.3 Asymmetry/skewness and kurtosis

Definition 4.6.8 (Standardized rvs). Given any RV X such that $\mathbb{E}[X^2] < +\infty$ and variance $\text{Var}[X] \in (0, +\infty)$, standardized rv Z is defined as:

$$Z = \frac{X - \mathbb{E}[X]}{\sqrt{\text{Var}[X]}} = \frac{X - \mathbb{E}[X]}{\sigma} \quad (4.52)$$

Remark 123. Note that $\mathbb{E}[Z] = 0$ and $\text{Var}[Z] = 1$:

$$\begin{aligned} \mathbb{E}[Z] &= \frac{\mathbb{E}[X] - \mathbb{E}[X]}{\sqrt{\text{Var}[X]}} = 0 \\ \text{Var}[Z] &= \text{Var}\left[\frac{X - \mathbb{E}[X]}{\sqrt{\text{Var}[X]}}\right] = \frac{1}{\text{Var}[X]} \cdot \text{Var}[X - \mathbb{E}[X]] = \frac{\text{Var}[X]}{\text{Var}[X]} = 1 \end{aligned}$$

Thus transform make rv independent from measure unit.

Definition 4.6.9 (r -th standardized moments of X). We have them if $g = \left(\frac{x - \mathbb{E}[X]}{\sigma}\right)^r$:

$$\bar{\mu}_r = \mathbb{E} \left[\left(\frac{X - \mathbb{E}[X]}{\sigma} \right)^r \right] = \begin{cases} \sum \left(\frac{x_i - \mathbb{E}[X]}{\sigma} \right)^r \cdot p_X(x_i) & \text{se } X \text{ è discreta} \\ \int_{-\infty}^{+\infty} \left(\frac{x - \mathbb{E}[X]}{\sigma} \right)^r \cdot f_X(x) dx & \text{se } X \text{ è continua} \end{cases} \quad (4.53)$$

Remark 124. Since for any rv $\bar{\mu}_0 = 1$, $\bar{\mu}_1 = 0$, $\bar{\mu}_2 = 1$ moments of interest are where $r = 3$ and $r = 4$.

4.6.3.1 Asymmetry/Skewness

Definition 4.6.10 (Symmetric rv). X is symmetric (respect to $\mathbb{E}[X]$) if $X - \mathbb{E}[X]$ has the same distribution of $\mathbb{E}[X] - X$.

Remark 125 (Intuizione significato). $X - \mathbb{E}[X]$ sposta la densità/probabilità, così com'è, centrandola sullo 0. Intuitivamente $-X$ ha l'effetto di ottenere la densità probabilità simmetrica/specchiata rispetto a $x = 0$; infine $-X + \mathbb{E}[X]$ specchia la densità/probabilità rispetto a 0 e poi la ricentra su 0. Pertanto se $X - \mathbb{E}[X]$ e $-X + \mathbb{E}[X]$ coincidono, è perché la distribuzione di partenza X è simmetrica rispetto al centro.

Proposition 4.6.8 (Simmetria di una vc continua (PDF)). *Sia X una vc continua con PDF f . Allora è simmetrica su $\mathbb{E}[X]$ se e solo se $f(x) = f(2\mathbb{E}[X] - x)$.*

Remark 126. La definizione è meramente quella di una funzione simmetrica rispetto a $x = \mu$ (vedi calcolo).

Dimostrazione. Sia F la CDF di X ; dimostriamo la doppia implicazione. Se la simmetria vale ($X - \mathbb{E}[X] = \mathbb{E}[X] - X$) abbiamo:

$$\begin{aligned} F(x) &= \mathbb{P}(X - \mathbb{E}[X] \leq x - \mathbb{E}[X]) \stackrel{(1)}{=} \mathbb{P}(\mathbb{E}[X] - X \leq x - \mathbb{E}[X]) \stackrel{(2)}{=} \mathbb{P}(X \geq 2\mathbb{E}[X] - x) \\ &= 1 - F(2\mathbb{E}[X] - x) \end{aligned}$$

dove in (1) abbiamo sfruttato la simmetria ($X - \mathbb{E}[X] = \mathbb{E}[X] - X$) e in (2) abbiamo elaborato algebricamente. Facendo la derivata dei membri estremi dell'equazione si ottiene $f(x) = f(2\mathbb{E}[X] - x)$.

Viceversa supponendo che $f(x) = f(2\mathbb{E}[X] - x)$ valga *forall* x , vogliamo dimostrare che $\mathbb{P}(X - \mathbb{E}[X] \leq t) = \mathbb{P}(\mathbb{E}[X] - X \leq t)$, ossia vi è simmetria e le cumulate CDF coincidono. Si ha

$$\begin{aligned} \mathbb{P}(X - \mathbb{E}[X] \leq t) &= \mathbb{P}(X \leq \mathbb{E}[X] + t) = \int_{-\infty}^{\mathbb{E}[X] + t} f(x) dx \stackrel{(1)}{=} \int_{-\infty}^{\mathbb{E}[X] + t} f(2\mathbb{E}[X] - x) dx \\ &\stackrel{(2)}{=} \int_{\mathbb{E}[X] - t}^{\infty} f(w) dw = \mathbb{P}(\mathbb{E}[X] - X \leq t) \end{aligned}$$

dove in abbiamo sfruttato che $f(x) = f(2\mathbb{E}[X] - x)$, mentre in (2) deve avvenire qualche trick di integrazione (integra $f(-x)$ ad indici invertiti e moltiplicati direi). \square

Definition 4.6.11 (Skewness). It's the 3-rd standardized moment:

$$\text{Asym}(X) = \bar{\mu}_3 = \mathbb{E} \left[\left(\frac{X - \mathbb{E}[X]}{\sigma} \right)^3 \right] \quad (4.54)$$

Remark 127. A negative skewness means a left longer tail, while positive a right longer one.

4.6.3.2 Kurtosis

Definition 4.6.12 (Kurtosis). It's the 4-th standardized moment

$$\text{Kurt}(X) = \bar{\mu}_4 = \mathbb{E} \left[\left(\frac{X - \mathbb{E}[X]}{\sigma} \right)^4 \right] \quad (4.55)$$

Remark 128. Some defines kurtosis by centering on 3 (value assumed by the normal) as in:

$$\text{Kurt}(X) = \mathbb{E} \left[\left(\frac{X - \mathbb{E}[X]}{\sigma} \right)^4 \right] - 3 \quad (4.56)$$

In this way the normal will have 0 kurtosis and the remaining a value a negative or positive value, related to givin less or more weight to the tail of the distribution.

Remark 129. Una distribuzione con eccesso di curtosi (4.56) negativo (detta *platicurtica*) tende ad avere un profilo più piatto della normale e una minore importanza delle code. Produce outlier in misura minore o meno estremi rispetto alla normale. Un esempio è l'uniforme.

Viceversa una distribuzione con eccesso di curtosi positivo è detta *leptocurtica* (ad esempio distribuzione T di Student, logistica, Laplace): ha code che si avvicinano allo zero più lentamente rispetto una gaussiana, per cui produce più outlier della stessa.

In fig 4.3 alcune distribuzioni (con media 0 e varianza 1) e relativa curtosi.

4.6.4 Other indicators

4.6.4.1 Mediana

4.6.4.2 Moda

4.7 Random vectors

4.7.1 Random vectors and their distribution

Definition 4.7.1. A random vector X (or n -variate random variable) is a function $X : \Omega \rightarrow \mathbb{R}^n$ that maps the occurrence of the experiment to a real vector of n components. It's denoted as

$$X = \begin{bmatrix} X_1 \\ \dots \\ X_n \end{bmatrix}$$

where X_1, \dots, X_n are real random variables.

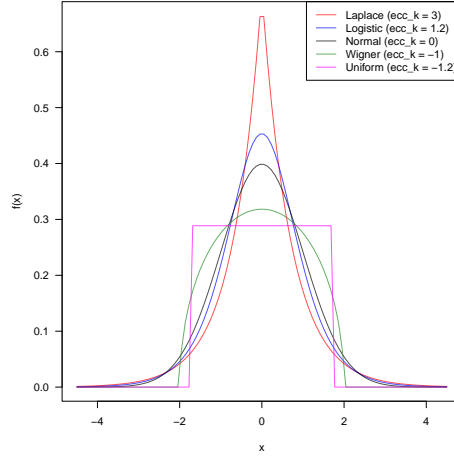


Figura 4.3: PDF for some rv (mean 0, variance 1) and their kurtosis

Example 4.7.1 (Two dice roll). With $\Omega = \{\{1, 1\}, \dots, \{6, 6\}\}$, we could construct following are *bivariate* random vectors:

- $X = (X_1, X_2)$ with X_1 outcome for the first dice, X_2 outcome of the second one;
- $X = (X_1, X_2)$ with X_1 sum of the two dice, X_2 difference

Important remark 34 (Probability of event E). It is defined as

$$\nu(E) = \mathbb{P}(X \in E), \quad \forall E \in \beta(\mathbb{R}^n)$$

Definition 4.7.2 (Distribution function). Distribution of random vector X is a function $F : \mathbb{R}^n \rightarrow \mathbb{R}$ defined as

$$F(x_1, \dots, x_n) = \mathbb{P}(X_1 \leq x_1, \dots, X_n \leq x_n) \quad (4.57)$$

Remark 130 (Remarks on distribution function). Again

- there's a 1-to-1 correspondance between F and ν expressed by

$$F(x_1, \dots, x_n) = \nu((-\infty, x_1] \times \dots \times (-\infty, x_n]), \quad \forall (x_1, \dots, x_n) \in \mathbb{R}^n$$

- F determines the probability distribution ν of X in the sense that

$$X \sim Y \iff X \text{ and } Y \text{ have the same distribution function}$$

4.7.2 Type of random vectors

Important remark 35 (Types of random vectors). Random vector X is

- **multivariate discrete** iff $\exists B \subset \mathbb{R}^n$, B finite or countable such that $\mathbb{P}(X \in B) = 1$
- **multivariate absolutely continuous** iff exists a density (called *joint*) $f : \mathbb{R}^n \rightarrow \mathbb{R}$ such that

1. $f \geq 0$
2. f is integrable
3. distribution is defined as integral of density

$$F(x_1, \dots, x_n) = \int_{-\infty}^{x_1} \cdots \int_{-\infty}^{x_n} f(t_1, \dots, t_n) dt_1 \dots dt_n, \quad \forall (x_1, \dots, x_n) \in \mathbb{R}^n$$

and for which

$$\int_{-\infty}^{+\infty} \cdots \int_{-\infty}^{+\infty} f(t_1, \dots, t_n) dt_1 \dots dt_n = 1$$

- **singular continuous** (ignored) is not easily to handle (splits in several cases)

Remark 131. Essentially the same remarks done for random variable holds. Next we extend to random vector the characterization theorem (already discussed in the $n = 1$ case) useful for proving that X is absolutely continuous. For this we need the concept of Lebesgue measure in $2+$ dimension.

Definition 4.7.3 (Lebesgue measure on \mathbb{R}^n). It's the only measure on $\beta(\mathbb{R}^n)$ such that the measure of the cartesian product of interval is equal to the product of the length of the intervals:

$$m(I_1 \times \dots \times I_n) = \text{len}(I_1) \cdot \dots \cdot \text{len}(I_n), \quad \forall I_i$$

where $\text{len}(I_i)$ is the length of the interval I_i (eg if $I_i = [a, b]$ that is $b - a$).

Example 4.7.2. Intuitively, if $A \in \beta(\mathbb{R}^2)$ then $m(A)$ is the area of A ; in $\beta(\mathbb{R}^3)$ is a volume and so on.

Theorem 4.7.1 (Absolutely continuous random vector characterization). *A random vector X is absolutely continuous if and only if any set (event) with null lebesgue measure has null probability as well*

$$X \text{ is absolutely continuous} \iff \begin{cases} \mathbb{P}(X \in E) = 0 \\ \forall E \in \beta(\mathbb{R}^n), \text{ such that } m(E) = 0 \end{cases}$$

Example 4.7.3. In a distribution in 2d (X_1, X_2) , if any point x_i, y_i (which has a 2d lebesgue measure of 0) has zero probability then that is an absolutely continuous random variable.

Theorem 4.7.2. *If X_1, \dots, X_n are absolutely continuous, this does not imply that the vector X is absolutely continuous.*

Example 4.7.4. As example of X not being absolutely continuous even if X_1, \dots, X_n are follows.

With $n = 2$, consider $X_1 \sim N(0, 1)$ (absolutely continuous because it's a standard normal), $X_2 = X_1$ (equal, so absolutely continuous). Is $\mathbf{X} = (X_1, X_2)$ absolutely continuous?

To check that X is not absolutely continuous, we apply the theorem, letting the event to be we extracted a point on the diagonal $y = x$

$$E = \{(x, y) \in \mathbb{R}^2 : x = y\}$$

We have that:

- $\mathbb{P}(X \in E) = 1$: infact once we extracted $X_1 = x_1$ we have $X_2 = x_1$ as well so the vector will be on the diagonal $y = x$;
- however $m(E) = 0$ (that is the area of the line $y = x$ compared to 2 space dimension \mathbb{R}^2 is 0).

So X is not absolutely continuous

4.7.3 Marginals

Definition 4.7.4 (Marginal of random vector $X = (X_1, \dots, X_n)^\top$). It is any subvector $(X_{j_1}, \dots, X_{j_k})$ where $\{j_1, \dots, j_k\}$ is a subset of $1, \dots, n$.

Remark 132. It is just a vector with less random variables; there are

- n marginals of only 1 variable (these are $X_1 \dots X_n$);
- $\binom{n}{2}$ marginals of 2 random variables $(X_i, X_j)^\top$;
- $\binom{n}{3}$ marginals of 3 random variables $(X_i, X_j, X_k)^\top$;

Theorem 4.7.3 (Density of marginal of absolute continuous random vectors). If $X = (X_1, \dots, X_n)$ is absolutely continuous with multivariate density f

- all marginal of X are still absolutely continuous (converse is not true in general but special case, see thm 4.7.5);
- the density g of the marginal $(X_1, \dots, X_k)^\top$ is obtained by making $n - k$ integral of f , that is integrating on the remaining $n - k$ variables one wants to eliminate:

$$g(x_1, \dots, x_k) = \underbrace{\int_{-\infty}^{+\infty} \dots \int_{-\infty}^{+\infty}}_{n-k \text{ integrals}} f(t_1, \dots, t_n) dt_{k+1} \dots dt_n \quad (4.58)$$

Example 4.7.5. If $n = 3$, $\mathbf{X} = (X, Y, Z)^\top$

- the density of $(Y, Z)^\top$ is

$$g(y, z) = \int_{-\infty}^{+\infty} f(x, y, z) dx$$

- density of Y is

$$g(y) = \int_{-\infty}^{+\infty} \int_{-\infty}^{+\infty} f(x, y, z) dx dz$$

4.7.4 Independence

Theorem 4.7.4 (Independence). Let $X = (X_1, \dots, X_n)^\top$ be any random vector with distribution function F then X_1, \dots, X_n are independent if and only if the

joint distribution function F is the product of the marginal distribution functions F_i :

$$X_1, \dots, X_n \text{ are independent} \iff F(x_1, \dots, x_n) = F_1(x_1) \cdot \dots \cdot F_n(x_n), \quad \forall \begin{bmatrix} X_1 \\ \dots \\ X_n \end{bmatrix} \in \mathbb{R}^n$$

Similarly if X is absolutely continuous we can replace distribution with densities, that is:

$$X_1, \dots, X_n \text{ are independent} \iff f(x_1, \dots, x_n) = f_1(x_1) \cdot \dots \cdot f_n(x_n), \quad \forall \begin{bmatrix} X_1 \\ \dots \\ X_n \end{bmatrix} \in \mathbb{R}^n$$

Theorem 4.7.5 (Independence and absolutely continuous random vectors). In $X = (X_1, \dots, X_n)$, if X_1, \dots, X_n are independent then

$$X \text{ is absolutely continuous} \iff X_1, \dots, X_n \text{ are absolutely continuous}$$

NB: non mostrato da rigo ma usato, trovato su <https://math.stackexchange.com/questions/606205>

Theorem 4.7.6 (Independence and expected value). If $X \perp\!\!\!\perp Y$ then $\mathbb{E}[XY] = \mathbb{E}[X] \mathbb{E}[Y]$

Dimostrazione. In case of absolutely continuous random variables, by definition they are independent if $f_{XY}(x, y) = f_X(x)f_Y(y)$. Then we have

$$\begin{aligned} \mathbb{E}[XY] &= \int_{-\infty}^{+\infty} x \cdot y \cdot f_{XY}(x, y) \, dx \, dy \\ &= \int_{-\infty}^{+\infty} x \cdot y \cdot f_X(x) f_Y(y) \, dx \, dy \\ &= \int_{-\infty}^{+\infty} x \cdot f_X(x) \, dx \int_{-\infty}^{+\infty} y \cdot f_Y(y) \, dy \\ &= \mathbb{E}[X] \mathbb{E}[Y] \end{aligned}$$

while the proof in the discrete case is analogous □

4.8 Relationship between RVs

4.8.1 Covariance

Definition 4.8.1 (Covariance). Considering two random variables X and Y , if $\mathbb{E}[|X|] < +\infty$, $\mathbb{E}[|Y|] < +\infty$ and $\mathbb{E}[|XY|] < +\infty$ we can define the covariance as

$$\text{Cov}(X, Y) = \mathbb{E}[(X - \mathbb{E}[X])(Y - \mathbb{E}[Y])] \quad (4.59)$$

Remark 133. La covarianza misura la forza del legame lineare tra X e Y

Important remark 36. A sufficient condition for the existence of $\text{Cov}(X, Y)$ is that $\mathbb{E}[X^2] < +\infty$ and $\mathbb{E}[Y^2] < +\infty$

Proposition 4.8.1 (Another useful formula).

$$\text{Cov}(X, Y) = \mathbb{E}[XY] - \mathbb{E}[X] \mathbb{E}[Y] \quad (4.60)$$

Dimostrazione. We have

$$\begin{aligned}\text{Cov}(X, Y) &= \mathbb{E}[(X - \mathbb{E}[X])(Y - \mathbb{E}[Y])] \\ &= \mathbb{E}[XY - \mathbb{E}[X] \cdot Y - \mathbb{E}[Y] \cdot X + \mathbb{E}[X] \mathbb{E}[Y]] \\ &= \mathbb{E}[XY] + \mathbb{E}[X] \mathbb{E}[Y] - 2 \mathbb{E}[X] \mathbb{E}[Y] \\ &= \mathbb{E}[XY] - \mathbb{E}[X] \mathbb{E}[Y]\end{aligned}$$

□

Example 4.8.1. In particular if $Y = X$

$$\text{Cov}(X, X) = \mathbb{E}[X^2] - (\mathbb{E}[X])^2 = \text{Var}[X]$$

Proposition 4.8.2 (Covariance and independence). *Assuming the covariance exists, $X \perp\!\!\!\perp Y \implies \text{Cov}(X, Y) = 0$. The converse implication is false.*

Dimostrazione. If $X \perp\!\!\!\perp Y$, then $\mathbb{E}[XY] = \mathbb{E}[X] \mathbb{E}[Y]$ so the covariance is 0. A counterexample for counterimplication follows. □

Example 4.8.2 (Counterexample where $\text{Cov}(X, Y) = 0$ but X, Y are not independent). Let $X \sim N(0, 1)$ and $Y = X^2$. Let's prove:

- $\text{Cov}(X, Y) = 0$. We have that

$$\text{Cov}(X, Y) = \mathbb{E}[XY] - \underbrace{\mathbb{E}[X] \mathbb{E}[Y]}_{=0} = \mathbb{E}[XY] = \mathbb{E}[X^3]$$

Since X is absolutely continuous (normal) the expectation of X to the power 3 can be written as

$$\mathbb{E}[X^3] = \int_{-\infty}^{+\infty} x^3 \cdot \frac{e^{-\frac{x^2}{2}}}{\sqrt{2\pi}}$$

and since the integrand it's an odd function evaluated on a symmetric interval, the integral is 0

- $X \not\perp\!\!\!\perp Y$. It's intuitive these are not independent, however let's prove it formally. To prove that we consider the probability $\mathbb{P}(-1 \leq X \leq 1, Y > 1)$ which under independence should be equal to $\mathbb{P}(-1 \leq X \leq 1) \cdot \mathbb{P}(Y > 1)$. Now actually:

$$\begin{aligned}\mathbb{P}(|X| \leq 1, Y > 1) &\stackrel{(1)}{=} \mathbb{P}(|X| \leq 1, |X| > 1) = \mathbb{P}(\emptyset) = 0 \\ \mathbb{P}(-1 \leq X \leq 1) \cdot \mathbb{P}(Y > 1) &= \underbrace{\mathbb{P}(|X| \leq 1)}_{>0} \cdot \underbrace{\mathbb{P}(|X| > 1)}_{>0} > 0\end{aligned}$$

where in (1) since $Y = X^2$.

Since the first is null and the second positive, they can't be equal and so random variables are not independent

Important remark 37 (Special case). There is an important special case in which independence amounts to null covariance:

$$\text{If } \begin{bmatrix} X \\ Y \end{bmatrix} \sim N, \text{ then } X \perp\!\!\!\perp Y \iff \text{Cov}(X, Y) = 0$$

Proposition 4.8.3 (Variance of sum of RVs). *Assuming the covariances exists, the variance of the sum of random variables is*

$$\begin{aligned} \text{Var} \left[\sum_{i=1}^n a_i X_i \right] &= \sum_{i=1}^n a_i^2 \text{Var} [X_i] + \sum_{i \neq j} a_i a_j \text{Cov} (X_i, X_j) \\ &\stackrel{(1)}{=} \sum_{i=1}^n a_i^2 \text{Var} [X_i] + 2 \sum_{1 \leq i < j \leq n} a_i a_j \text{Cov} (X_i, X_j) \end{aligned} \quad (4.61)$$

where (1) because $\text{Cov} (X_i, X_j) = \text{Cov} (X_j, X_i)$

Important remark 38. If X_1, \dots, X_n are independent, this formula reduces to

$$\text{Var} \left[\sum_{i=1}^n a_i X_i \right] = \sum_{i=1}^n a_i^2 \text{Var} [X_i]$$

Example 4.8.3. With $n = 2$, by just letting $a_1 = 1, a_2 = \pm 1$

$$\begin{aligned} \text{Var} [X + Y] &= \text{Var} [X] + \text{Var} [Y] + 2 \text{Cov} (X, Y) \\ \text{Var} [X - Y] &= \text{Var} [X] + \text{Var} [Y] - 2 \text{Cov} (X, Y) \end{aligned}$$

In case X, Y are independent covariance is null and 1) variance of sum is sum of variance 2) $\text{Var} [X + Y] = \text{Var} [X - Y]$

Example 4.8.4. Sia $X \perp\!\!\!\perp Y$,

$$\text{Var} [X - Y + 2Z] = \text{Var} [X] + \text{Var} [Y] + 4 \text{Var} [Z] + 4 \text{Cov} (X, Z) - 4 \text{Cov} (Y, Z)$$

Proposition 4.8.4 (Proprietà covarianza (wikipedia, non fatte da rigo)). *If X, Y, W, V are real-valued random variables and $a, b, c, d \in \mathbb{R}$, then the following facts are a consequence of the definition of covariance:*

$$\text{Cov} (X, a) = 0 \quad (4.62)$$

$$\text{Cov} (X, X) = \text{Var} [X] \quad (4.63)$$

$$\text{Cov} (X, Y) = \text{Cov} (Y, X) \quad (4.64)$$

$$\text{Cov} (aX, bY) = ab \text{Cov} (X, Y) \quad (4.65)$$

$$\text{Cov} (X + a, Y + b) = \text{Cov} (X, Y) \quad (4.66)$$

$$\text{Cov} (aX + bY, cW + dV) = ac \text{Cov} (X, W) + ad \text{Cov} (X, V) + bc \text{Cov} (Y, W) + bd \text{Cov} (Y, V) \quad (4.67)$$

$$X \perp\!\!\!\perp Y \implies \text{Cov} (X, Y) = 0 \quad (4.68)$$

Example 4.8.5 (Esame vecchio viroli). Let X_1 and X_2 be two random variables with distribution $X_1 \sim N(0, 2)$ and $X_2 \sim N(-2, 1)$ (parameters are mean and variance) and covariance -1 . Compute $\text{Cov} (X_1 + X_2, X_1 - X_2)$. We have that

$$\begin{aligned} \text{Cov} (X_1 + X_2, X_1 - X_2) &= \text{Cov} (X_1, X_1) - \text{Cov} (X_1, X_2) + \text{Cov} (X_2, X_1) - \text{Cov} (X_2, X_2) \\ &= \text{Var} [X_1] - \text{Var} [X_2] = 2 - 1 = 1 \end{aligned}$$

Example 4.8.6 (Esame vecchio viroli). Let X_1, X_2 be two standard gaussian variables with covariance -1. Compute $\text{Cov}(X_1 + X_2, X_1 - X_2)$.
With the same developmet as above we have:

$$\text{Cov}(X_1 + X_2, X_1 - X_2) = \text{Var}[X_1] - \text{Var}[X_2] = 1 - 1 = 0$$

Example 4.8.7 (Esame vecchio viroli). Let X and Y be two independent bernoulli random variables with same parameter p . Compute $\text{Cov}(Y - X, 2X + 2Y)$.

$$\begin{aligned} \text{Cov}(Y - X, 2X + 2Y) &= 2\text{Cov}(X, Y) + 2\text{Cov}(Y, Y) - 2\text{Cov}(X, X) - 2\text{Cov}(X, Y) \\ &= 2\text{Var}[X] - 2\text{Var}[Y] = 0 \end{aligned}$$

taluni suggeriscono -1 ma mi pare na gran cacata

```
## [1] -8.398084e-07
```

Example 4.8.8 (Esame vecchio viroli). Let $X = (X_1, X_2)$ be a bivariate gaussian vector with $\mu = [0, 0]$ and

$$\Sigma = \begin{bmatrix} 1 & 0.5 \\ 0.5 & 1 \end{bmatrix}$$

What is the distribution of $Y = 3X_1 - 2X_2$?

Si ha che

$$\begin{aligned} \mathbb{E}[Y] &= \mathbb{E}[3X_1 - 2X_2] = 3\mathbb{E}[X_1] - 2\mathbb{E}[X_2] = 0 \\ \text{Var}[Y] &= \text{Var}[3X_1 - 2X_2] \\ &= 3^2 \text{Var}[X_1] + (-2)^2 \text{Var}[X_2] + 2(3 \cdot (-2)) \text{Cov}(X_1, X_2) \\ &= 9 \cdot 1 + 4 \cdot 1 + 2 \cdot (-6) \cdot 0.5 = 7 \end{aligned}$$

quindi è $Y \sim N(0, 7)$ come confermato da taluni

4.8.2 Correlation coefficient

Definition 4.8.2 (Correlation coefficient). Considered two rv X, Y , if $\mathbb{E}[X^2] < +\infty$, $\mathbb{E}[Y^2] < +\infty$, $\text{Var}[X] > 0$, $\text{Var}[Y] > 0$, we can define the correlation coefficient as

$$\rho(X, Y) = \text{Corr}(X, Y) = \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}[X]}\sqrt{\text{Var}[Y]}} \quad (4.69)$$

Proposition 4.8.5. *Some properties:*

- it can be easily seen that correlation coefficient can be written as covariance between the two standardized variables (def 4.6.8)

$$\text{Corr}(X, Y) = \text{Cov}\left(\frac{X - \mathbb{E}[X]}{\sqrt{\text{Var}[X]}}, \frac{Y - \mathbb{E}[Y]}{\sqrt{\text{Var}[Y]}}\right)$$

quindi $\text{Cov}(X, Y)$ e $\text{Corr}(X, Y)$ danno essenzialmente la stessa informazione, entrambi misurano l'intensità del legame lineare tra X ed Y .
 $\text{Corr}(\cdot, \cdot)$ ha il vantaggio di essere normalizzato

- it ranges in

$$-1 \leq \text{Corr}(X, Y) \leq 1$$

Note that this amounts to

$$\text{Cov}(X, Y)^2 \leq \text{Var}[X] \text{Var}[Y]$$

Quindi se il legame lineare tra X ed Y è molto alto $\text{Corr}(X, Y)$ è vicino a 1 in termini assoluti, ma non riesco ad esprimere questo fatto in termini di covarianza perché quest'ultima dipende dalle unità di misura di X ed Y

- the limit cases are the following:

$$\text{Corr}(X, Y) = 1 \iff Y = a + bX, b > 0$$

$$\text{Corr}(X, Y) = -1 \iff Y = a + bX, b < 0$$

Example 4.8.9 (Esame vecchio viroli). Let X and Y be two gaussian variables with zero mean $\text{Var}[X] = 1$, $\text{Var}[Y] = 9$, covariance -1 , compute $\rho(X + Y, X)$.

$$\begin{aligned} \text{Corr}(X + Y, X) &= \frac{\text{Cov}(X + Y, X)}{\sqrt{\text{Var}[X + Y]}\sqrt{\text{Var}[X]}} = \frac{\text{Cov}(X, X) + \text{Cov}(Y, X)}{\sqrt{\text{Var}[X] + \text{Var}[Y] + 2\text{Cov}(X, Y)}\sqrt{\text{Var}[X]}} \\ &= \frac{\text{Var}[X] + \text{Cov}(X, Y)}{\sqrt{\text{Var}[X] + \text{Var}[Y] + 2\text{Cov}(X, Y)}\sqrt{\text{Var}[X]}} = \frac{1 + (-1)}{\sqrt{1 + 9 + 2(-1)}\sqrt{1}} \\ &= 0 \end{aligned}$$

Il risultato è confermato dal Bigo.

Example 4.8.10 (Esame vecchio viroli). Let X and Y be two gaussian variables with zero mean $\text{Var}[X] = 1$, $\text{Var}[Y] = 9$, covariance -1 , compute $\rho(1 - 2X + 2, 3 + Y)$.

We have:

$$\begin{aligned} \text{Corr}(1 - 2X + 2, 3 + Y) &= \text{Corr}(3 - 2X, 3 + Y) = \frac{\text{Cov}(3 - 2X, 3 + Y)}{\sqrt{\text{Var}[3 - 2X]}\sqrt{\text{Var}[3 + Y]}} \\ &= \frac{\text{Cov}(3, 3) + \text{Cov}(3, Y) + \text{Cov}(-2X, 3) + \text{Cov}(-2X, Y)}{\sqrt{4\text{Var}[X]}\sqrt{\text{Var}[Y]}} \\ &= \frac{0 + 0 + 0 + 2}{2 \cdot 1 \cdot 3} = \frac{1}{3} \end{aligned}$$

come confermato da taluni

4.9 Exercises

4.9.1 Functions and moments

Example 4.9.1 (crashcourse, day 1 es 4 pag 7). Let $F(x) = \frac{c}{2}(1 - \frac{1}{x^2})$ for $x \in [1, \infty)$:

1. obtain $f(x)$

2. obtain c
3. $\mathbb{E}[X]$
4. $\text{Var}[X]$

we have

1.

$$f(x) = \frac{\partial}{\partial x} F(x) = \frac{\partial}{\partial x} \frac{c}{2} \left(1 - \frac{1}{x^2} \right) = \frac{c}{x^3}$$

2.

$$c \int_1^\infty \frac{1}{x^3} dx = c \left[-\frac{1}{2} \frac{1}{x^2} \right]_1^\infty = \frac{c}{2} = 1 \implies c = 2$$

3.

$$2 \int_1^\infty x \frac{1}{x^3} dx = 2 \int_1^\infty \frac{1}{x^2} dx = 2 \left[-\frac{1}{x} \right]_1^\infty = 2$$

4. first we find

$$\mathbb{E}[X^2] = 2 \int_1^\infty x^2 \frac{1}{x^3} dx = 2 \int_1^\infty \frac{1}{x} dx = 2 [\log x]_1^\infty = +\infty$$

4.9.2 Random vectors

Example 4.9.2 (Esame vecchio viroli). Let $\mathbf{X} = (X, Y)^\top$ be a random vector with joint density

$$f(x, y) = ky$$

where $0 < x < y < 1$. Compute k .

In order to compute k it must be:

$$\begin{aligned} 1 &= \int_0^1 \int_0^y ky \, dx \, dy = k \int_0^1 y \int_0^y 1 \, dx \, dy \\ &= k \int_0^1 y [x]_0^y \, dy = k \int_0^1 y^2 \, dy = k \left[\frac{y^3}{3} \right]_0^1 = \frac{k}{3} \end{aligned}$$

da cui $k = 3$

Example 4.9.3. Consider the function

$$f(x|y) = \begin{cases} \frac{y^x e^{-y}}{x!} & \text{for } x = 0, 1, 2, \dots \text{ and } y \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

1. if the marginal pdf of Y is $\text{Exp}(1)$, what is the joint probability function of (X, Y)
2. derive the marginal probability function of X

We have:

1. for the joint probability

$$f_{X,Y}(x, y) = f(y) \cdot f(x|y) = e^{-y} \frac{y^x e^{-y}}{x!} = \frac{y^x e^{-2y}}{x!}$$

2. for the marginal probability of X

$$\begin{aligned} f_X(x) &= \int_0^{+\infty} \frac{y^x e^{-2y}}{x!} dy = \frac{1}{x!} \underbrace{\int_0^{+\infty} y^x e^{-2y} dy}_{(1)} \\ &= \frac{1}{x!} \frac{\Gamma(x+1)}{2^{x+1}} = \frac{1}{2^{x+1}} \end{aligned}$$

where (1) is the kernel of a Gamma ($\alpha = x + 1, \beta = 2$)

Capitolo 5

Discrete random variables

5.1 Recap notazione

Remark 134. Note Soffritti:

- Azzalini: occhio alla notazione, non sempre coerente (stesso simbolo associato a cose diverse)
- per esercizi inferenza fare meno riferimento ad azzalini e guardare il resto del materiale didattico

Remark 135. Al fine della carrellata di variabili casuali univariate (discrete e continue) fissiamo alcuni punti e aggiungiamo i concetti di famiglia di variabili casuali e spazio parametrico

Remark 136 (Ripasso Soffritti). We have that

- a random variable X is a particular function which maps outcome of an experiment, belonging to the sample space, to another set $X : \Omega \rightarrow S$ (basically always numerical $S = \mathbb{R}$)
- *support/range space* R_X (supporto) is the subset of \mathbb{R} with all possible outcome of the random variable: $R_X = X(\Omega)$.
If R_X is finite or countable the random variable is discrete otherwise if it's uncountable (absolutely) continuous;
- X (maiuscolo) is considered a random variable, x (minuscolo) one of its realization
- all random variable have a *distribution function* $F_X(x) = \mathbb{P}(X \leq x), \forall x \in \mathbb{R}$ with following characterizing properties

$$\begin{aligned}\lim_{x \rightarrow -\infty} F_X(x) &= 0 = F_X(-\infty) \\ \lim_{x \rightarrow +\infty} F_X(x) &= 1 = F_X(+\infty) \\ x_1 \leq x_2 &\implies F_X(x_1) \leq F_X(x_2) \\ \lim_{h \rightarrow 0^+} F_X(x+h) &= F_X(x), \forall x \in \mathbb{R}\end{aligned}$$

Conosco	Voglio ottenere	Operazione svolta
$p_X(x)$	$F_X(t) = \sum_{x \leq t} p_X(x)$	ottengo la f. di distribuzione mediante somma
$f_X(x)$	$F_X(t) = \sum_{-\infty}^t f_X(x) dx$	ottengo la f. di distribuzione mediante integrazione
$F_X(x)$	$f_X(x) = F'_X(x) = \frac{\partial F_X(x)}{\partial x}$	ottengo densità mediante derivazione

Tabella 5.1: Procedimenti

- if X is discrete it has a *probability mass function* $p_X(x) = \mathbb{P}(X = x), \forall x \in \mathbb{R}$ which is a function $p_X : \mathbb{R} \rightarrow [0, 1]$ with following properties:

$$0 \leq p_X(x) \leq 1$$

$$\sum_{x \in R_X} p_X(x) = 1$$

- if X is continuous it has a *density function* $f_X(x)$ of type $f : \mathbb{R} \rightarrow [0, \infty)$ with properties

$$f_X(x) \geq 0, \forall x \in \mathbb{R}$$

$$\int_{-\infty}^{+\infty} f_X(x) dx = 1$$

$$\mathbb{P}(X = x) = 0, \forall x \in \mathbb{R}$$

$$\int_a^b f_X(x) dx = \mathbb{P}(a \leq X \leq b) = F_X(b) - F_X(a)$$

In merito alla terza, $\mathbb{P}(X = x) = 0$, quest'ultima se pensiamo in termini di probabilità come eventi favorevoli su possibili, la probabilità è 0 perché abbiamo un unico caso favorevole ($X=x$) e infiniti casi possibili (essendo il supporto infinito non numerabile).

Si chiama densità perché non ci fornisce direttamente una probabilità (a differenza della pmf per le discrete) ma ci permette di calcolarla via integrazione

- procedimenti per passare da una funzione all'altra sono riportati in tabella 5.1
- una funzione molto utilizzata nel prosieguo è la funzione indicatrice di insieme. Se A è un insieme ed x un oggetto, la funzione indicatrice dell'insieme A è così definita:

$$\mathbb{1}_A(x) = \begin{cases} 1 & \text{if } x \in A \\ 0 & \text{if } x \notin A \end{cases}$$

Remark 137. Le variabili casuali possono essere classificate in famiglie, all'interno delle quali le diverse variabili casuali differiscono per il valore assunto da 1+ parametri.

Definition 5.1.1 (Family of random variables). Set of distribution function $F(x; \Theta)$ having the same functional form but different for one or more parameters.

Important remark 39. Per individuare i parametri di una distribuzione guardare alla formula della density o pmf: tutte le lettere, ad eccezione di x , sono considerati come parametri della distribuzione e ne determinano la forma.

Definition 5.1.2 (Parameters space). Θ , it's the set of possible value for the parameters of a distribution function.

Remark 138. In questo capitolo guardiamo alle famiglie di funzioni discrete

5.2 Dirac

5.2.1 Defintion

Remark 139. Sono la traduzione in termini di variabile casuale di una costante.

Example 5.2.1. Un esempio di Dirac è X = “numero di occhi in individuo sano”

Definition 5.2.1 (Dirac rv (degenere)). $X \sim \delta_c$ if $\mathbb{P}(X = c) = 1$.

5.2.2 Functions

Remark 140 (Support and parametric space). If $X \sim \delta_c$

$$\begin{aligned} R_X &= \{c\} \\ \Theta &= \{c \in \mathbb{R}\} \end{aligned}$$

Proposition 5.2.1 (Dirac PMF). *Defined as*

$$p_X(x) = \begin{cases} 1, & \text{if } x = c \\ 0, & \text{if } x \neq c \end{cases} = \mathbb{1}_{R_X}(x)$$

Proposition 5.2.2 (Dirac distribution function).

$$F_X(x) = \mathbb{P}(X \leq x) = \begin{cases} 0, & \text{if } x < c \\ 1, & \text{if } x \geq c \end{cases} = \mathbb{1}_{[c, +\infty)}(x) \quad (5.1)$$

5.2.3 Moments

Proposition 5.2.3 (Moments).

$$\begin{aligned} \mathbb{E}[X] &= c \\ \text{Var}[X] &= 0 \end{aligned}$$

Dimostrazione.

$$\begin{aligned} \mathbb{E}[X] &= c \cdot 1 = c \\ \text{Var}[X] &= \mathbb{E}[X^2] - (\mathbb{E}[X])^2 = c^2 \cdot 1 - c^2 = 0 \end{aligned}$$

□

5.2.4 Shape

5.2.5 Extras

TODO: plot dirac PMF distribution function

Remark 141. Dirac is the only random variable having null variance.

5.3 Bernoulli

5.3.1 Definition

Remark 142. Viene utilizzata quando si ha a che fare con un esperimento il cui esito possibile è dicotomico (es $X = 1$ successo, $X = 0$ insuccesso).

Definition 5.3.1 (vc di Bernoulli). X is distributed as Bernoulli with parameter p ($0 \leq p \leq 1$), written $X \sim \text{Bern}(p)$, if $\mathbb{P}(X = 1) = p$ and $\mathbb{P}(X = 0) = 1 - p$.

Remark 143. If $p = 0 \vee p = 1$ we obtain a Dirac.
In inferenza p (probabilità di successo) è il parametro di interesse da trovare.

5.3.2 Functions

Remark 144 (Support and parametric space).

$$\begin{aligned} R_X &= \{0, 1\} \\ \Theta &= \{p \in \mathbb{R} : 0 \leq p \leq 1\} \end{aligned}$$

Definition 5.3.2 (PMF).

$$p_X(x) = \mathbb{P}(X = x) = p^x \cdot (1 - p)^{1-x} \cdot \mathbb{1}_{R_X}(x) \quad (5.2)$$

Remark 145. Mediante l'utilizzo della indicatrice posso

- evitare di scrivere la pmf più compattamente del caso seguente (equivalente)

$$p_X(x) = \mathbb{P}(X = x) = \begin{cases} p^x \cdot (1 - p)^{1-x} & \text{se } x \in R_X \\ 0 & \text{altrimenti} \end{cases}$$

- posso potenzialmente valutare $p_X(x)$ per qualunque valore di \mathbb{R} (non solo quelli appartenenti al supporto) settando ciò che non è rilevante a 0

Definition 5.3.3 (PDF).

$$F_X(x) = \mathbb{P}(X \leq x) = \begin{cases} 0 & \text{se } x < 0 \\ 1 - p & \text{se } 0 \leq x < 1 \\ 1 & \text{se } x \geq 1 \end{cases} \quad (5.3)$$

5.3.3 Moments

Proposition 5.3.1 (Momenti caratteristici).

$$\mathbb{E}[X] = p \quad (5.4)$$

$$\text{Var}[X] = p(1-p) \quad (5.5)$$

$$\text{Asym}(X) = \frac{1-2p}{\sqrt{p(1-p)}} \quad (5.6)$$

$$\text{Kurt}(X) = \frac{3p^2 - 3p + 1}{p(1-p)} \quad (5.7)$$

Dimostrazione. Per il valore atteso

$$\mathbb{E}[X] = 1 \cdot p + 0 \cdot (1-p) = p$$

Per la varianza, dato che $X^2 = X$ e dunque $\mathbb{E}[X^2] = \mathbb{E}[X]$ si ha:

$$\text{Var}[X] = \mathbb{E}[X^2] - \mathbb{E}[X]^2 = p - p^2 = p(1-p)$$

□

Remark 146. In particolare il valore atteso coincide con la probabilità di successo e la varianza è sempre compresa nell'intervallo $[0; 0.25]$, raggiungendo il massimo per $p = 1/2$.

5.3.4 Shape

5.3.5 Extras

5.3.5.1 Indicator RV

Important remark 40. Any event $A \subseteq \Omega$ is associated to a Bernoulli indicator random variable. This rv provides a link between probability of an event and expected value

Definition 5.3.4 (Indicator rv of event A). Let $\Omega = \{\omega_1, \omega_2, \dots\}$ be the sample space of the experiment considered and $A \subseteq \Omega$ a possible event; suppose that ω is the outcome that currently happens as a result of the experiment. Then:

$$I_A = \mathbb{1}_A(\omega) = \begin{cases} 1 & \text{if } A \text{ verifies: } \omega \in A \\ 0 & \text{if } A \text{ does not: } \omega \notin A \end{cases}$$

therefore if $\mathbb{P}(A) = p$, then $I_A \sim \text{Bern}(p)$

Proposition 5.3.2. Probability of an event A to occur and the expected value of its indicator rv I_A are 1-1 linked as follows:

$$\mathbb{P}(A) = \mathbb{E}[I_A] \quad (5.8)$$

Dimostrazione. For any event A we can define a rv I_A , and viceversa for each I_A there's one event A defined as $A = \{\omega \in \Omega : I_A(\omega) = 1\}$. If $p = \mathbb{P}(A)$ and thus $I_A \sim \text{Bern}(p)$ we have

$$\mathbb{E}[I_A] = \mathbb{E}[\text{Bern}(p)] = 0 \cdot \mathbb{P}(I(A) = 0) + 1 \cdot \mathbb{P}(I(A) = 1) \mathbb{P}(A) = p$$

□

TODO: fare grafico PMF e PDF della bernoulli. 0 è uno stick alto $1-p$, 1 è alto p

Proposition 5.3.3 (Indicator RV properties).

$$(I_A)^n = I_A, \quad \forall n \in \mathbb{N} : n > 0 \quad (5.9)$$

$$I_{\bar{A}} = 1 - I_A \quad (5.10)$$

$$I_{A \cap B} = I_A \cdot I_B \quad (5.11)$$

$$I_{A \cup B} = I_A + I_B - I_A \cdot I_B \quad (5.12)$$

Dimostrazione. La 5.9 vale dato che $0^n = 0$ e $1^n = 1$ per qualsiasi intero positivo n . La 5.10 vale dato che $1 - I_A$ è 1 se A non accade e 0 se accade. Per la 5.11, $I_A \cdot I_B$ è 1 solo se sia I_A che I_B sono 1 e 0 altrimenti. Per la 5.12,

$$\begin{aligned} I_{A \cup B} &\stackrel{(1)}{=} 1 - I_{\bar{A} \cap \bar{B}} = 1 - I_{\bar{A}} \cdot I_{\bar{B}} = 1 - (1 - I_A)(1 - I_B) \\ &= I_A + I_B - I_A I_B \end{aligned}$$

dove in (1) abbiamo sfruttato De Morgan. \square

Remark 147 (Usefulness). Previous result enable to express any probability as expected value; some examples come in the following section.

Furthermore indicator rvs are useful in exercises on expected value: often we can define a complex rv of unknown/complex distribution function as sum of indicator function (simpler). The so-called fundamental bridge enable then, applying expected value properties, to find expected value of unknown complex distribution function

5.3.5.2 Applications: probability

Proposition 5.3.4 (Boole inequality). *If E_1, \dots, E_n are events we have:*

$$\mathbb{P}(E_1 \cup \dots \cup E_n) \leq \mathbb{P}(E_1) + \dots + \mathbb{P}(E_n) \quad (5.13)$$

Dimostrazione. Let E_1, \dots, E_n be the events considered; we note that

$$I_{E_1 \cup \dots \cup E_n} \leq I_{E_1} + \dots + I_{E_n}$$

since left branch is 1 if all the events occur while right one is 1 even if only one does. Taking expected value:

$$\mathbb{E}[I_{E_1 \cup \dots \cup E_n}] \leq \mathbb{E}[I_{E_1} + \dots + I_{E_n}] \quad \text{by linearity of expectation} \dots$$

$$\mathbb{E}[I_{E_1 \cup \dots \cup E_n}] \leq \mathbb{E}[I_{E_1}] + \dots + \mathbb{E}[I_{E_n}] \quad \text{applying 5.8} \dots$$

$$\mathbb{P}(E_1 \cup \dots \cup E_n) \leq \mathbb{P}(E_1) + \dots + \mathbb{P}(E_n)$$

\square

Proposition 5.3.5 (Bonferroni inequality). *If E_1, \dots, E_n are events:*

$$\mathbb{P}(E_1 \cap \dots \cap E_n) \geq 1 - \sum_{i=1}^n \mathbb{P}(\bar{E}_i) \quad (5.14)$$

Dimostrazione. Similarly to the Boole inequality, applying DeMorgan

$$I_{E_1 \cap \dots \cap E_n} = 1 - I_{\bar{E}_1 \cup \dots \cup \bar{E}_n}$$

Taking expected value:

$$\begin{aligned}\mathbb{E}[I_{E_1 \cap \dots \cap E_n}] &= \mathbb{E}[1 - I_{\overline{E_1} \cup \dots \cup \overline{E_n}}] && \text{per linearità} \dots \\ \mathbb{E}[I_{E_1 \cap \dots \cap E_n}] &= 1 - \mathbb{E}[I_{\overline{E_1} \cup \dots \cup \overline{E_n}}] && \text{passando alle probabilità} \dots \\ \mathbb{P}(E_1 \cap \dots \cap E_n) &= 1 - \mathbb{P}(\overline{E_1} \cup \dots \cup \overline{E_n})\end{aligned}$$

Finally applying 5.13

$$\mathbb{P}(E_1 \cap \dots \cap E_n) = 1 - \mathbb{P}(\overline{E_1} \cup \dots \cup \overline{E_n}) \geq 1 - \mathbb{P}(\overline{E_1}) - \dots - \mathbb{P}(\overline{E_n})$$

□

Proposition 5.3.6 (Inclusion/exclusion principle). *In case of two events*

$$\mathbb{P}(A \cup B) = \mathbb{P}(A) + \mathbb{P}(B) - \mathbb{P}(A \cap B) \quad (5.15)$$

In general:

$$\mathbb{P}\left(\bigcup_{i=1}^n E_i\right) = \sum_{r=1}^n (-1)^{r+1} \sum_{i_1 < \dots < i_r} \mathbb{P}(E_{i_1} \cap E_{i_2} \cap \dots \cap E_{i_r}) \quad (5.16)$$

$$= \sum_i \mathbb{P}(E_i) - \sum_{i < j} \mathbb{P}(E_i \cap E_j) + \sum_{i < j < k} \mathbb{P}(E_i \cap E_j \cap E_k) - \dots + (-1)^{n+1} \mathbb{P}(E_1 \cap \dots \cap E_n) \quad (5.17)$$

Dimostrazione. Given 5.15 we take expected value of both branch of 5.12. Considering 5.16, we can apply indicator rv properties

$$\begin{aligned}1 - I_{E_1 \cup \dots \cup E_n} &= I_{\overline{E_1} \cap \dots \cap \overline{E_n}} \\ &= I_{\overline{E_1}} \cdot \dots \cdot I_{\overline{E_n}} \\ &= (1 - I_{E_1}) \cdot \dots \cdot (1 - I_{E_n}) \\ &\stackrel{(1)}{=} 1 - \sum_i I_{E_i} + \sum_{i < j} I_{E_i} I_{E_j} - \dots + (-1)^n I_{E_1} \cdot \dots \cdot I_{E_n}\end{aligned}$$

where in (1):

- il 1 significa selezionare tutti gli 1 negli n fattori;
- il $\sum_i I_{E_i}$ si ottiene selezionando tutti gli 1 a meno di un fattore a turno che ha sempre il segno $-$ davanti;
- $\sum_{i < j} I_{E_i} I_{E_j}$ si ottiene selezionando tutti gli 1 ad eccezione di due fattori.

Prendendo i valori attesi di ambo i membri si ha

$$\begin{aligned}\mathbb{E}[1 - I_{E_1 \cup \dots \cup E_n}] &= \mathbb{E}\left[1 - \sum_i I_{E_i} + \sum_{i < j} I_{E_i} I_{E_j} - \dots + (-1)^n I_{E_1} \cdot \dots \cdot I_{E_n}\right] \\ 1 - \mathbb{E}[I_{E_1 \cup \dots \cup E_n}] &\stackrel{(1)}{=} 1 - \mathbb{E}\left[\sum_i I_{E_i} - \sum_{i < j} I_{E_i} I_{E_j} + \dots + (-1)^{n+1} I_{E_1} \cdot \dots \cdot I_{E_n}\right] \\ \mathbb{E}[I_{E_1 \cup \dots \cup E_n}] &= \mathbb{E}\left[\sum_i I_{E_i}\right] - \mathbb{E}\left[\sum_{i < j} I_{E_i} I_{E_j}\right] + \dots + \mathbb{E}[(-1)^{n+1} I_{E_1} \cdot \dots \cdot I_{E_n}] \\ \mathbb{P}\left(\bigcup_{i=1}^n E_i\right) &= \sum_i \mathbb{P}(E_i) - \sum_{i < j} \mathbb{P}(E_i \cap E_j) + \dots + (-1)^{n+1} \mathbb{P}(E_1 \cap \dots \cap E_n)\end{aligned}$$

dove in (1) abbiamo raccolto un meno al secondo membro entro parentesi. \square

5.3.5.3 Applications: expected value evaluation

Example 5.3.1 (Matching carte). Abbiamo un mazzo di n carte numerate da 1 a n ben mischiato. Una carta è un match se la sua posizione nell'ordine del mazzo matcha con il suo numero. Sia X il numero totale di match nel mazzo: qual è il valore atteso di X ?

Se scriviamo $X = I_1 + \dots + I_n$ con

$$I_i = \begin{cases} 1 & \text{se l}'i\text{-esima carta matcha col proprio numero} \\ 0 & \text{altrimenti} \end{cases}$$

Si ha che, non condizionando a nulla e pensando ad un singolo shuffle/match

$$\mathbb{E}[I_i] = \frac{1}{n}$$

pertanto per linearità

$$\mathbb{E}[X] = \mathbb{E}[I_1] + \dots + \mathbb{E}[I_n] = n \cdot \frac{1}{n} = 1$$

Quindi il numero di match medi è 1, indipendentemente da n . Anche se I_i sono dipendenti in maniera complicata, la linearità del valore atteso vale sempre.

Example 5.3.2 (Valore atteso di Ipergeometrica Negativa). Un'urna contiene w palline bianche e b palline nere che sono estratte senza reinserimento. Il numero di palline nere estratte prima di pescare la prima bianca ha una distribuzione Ipergeometrica negativa (in tab 5.2 una sintesi dei casi). Trovare il valore atteso.

Trovarlo dalla definizione della variabile è complicato, ma possiamo esprimere la variabile come somma di indicatrici. Etichettiamo le palline nere con $1, 2, \dots, b$ e sia I_i l'indicatrice che la pallina nera i è stata estratta prima di qualsiasi bianca. Si ha che

$$\mathbb{P}(I_i = 1) = \frac{1}{w+1}$$

Fisso ...	con reinserimento	senza reinserimento
n trial	binomiale	ipergeometrica
n successi	binomiale negativa	ipergeometrica negativa

Tabella 5.2

dato nel listare l'ordine in cui la pallina nera i e le altre bianche son pescate (ignorando le altre) tutti gli ordine sono equiprobabili. Pertanto per linearità

$$\mathbb{E} \left[\sum_{i=1}^b I_i \right] = \sum_{i=1}^b \mathbb{E} [I_i] = \frac{b}{w+1}$$

La risposta ha n senso dato che aumenta con b , diminuisce con w ed è corretta nei casi estremi $b = 0$ (nessuna pallina nera sarà estratta) e $w = 0$ (tutte le palline nere saranno esaurite prima di pescare una non esistente bianca).

5.4 Binomial

5.4.1 Definition

Remark 148. Used to know the probability of having x success among $n \geq x$ independent Bernoulli trial with common probability success p .

Definition 5.4.1 (vc binomiale). Si utilizza quando eseguiamo n prove bernoulliane mantenendo le condizioni sperimentali costanti: ossia consideriamo le bernoulliane indipendenti, aventi comune probabilità di successo p (es estrazioni di palline bianche/nere da una urna con reimmissione).

Il numero di successi X ottenuti in n prove bernoulliane indipendenti si distribuisce come una vc binomiale di parametri n e p , e si scrive $X \sim \text{Bin}(n, p)$.

La cosa aleatoria è il numero di successi, mentre il numero di prove n è prefissato

Remark 149. Se $n = 1$ la distribuzione Binomiale coincide con quella di Bernoulli, ossia $\text{Bin}(1, p) = \text{Bern}(p)$

5.4.2 Functions

Remark 150 (Supporto e spazio parametrico).

$$\begin{aligned} R_X &= \{0, 1, \dots, n\} \\ \Theta &= \{n \in \mathbb{N} \setminus \{0\}, p \in \mathbb{R} : 0 \leq p \leq 1\} \end{aligned}$$

Definition 5.4.2 (Funzione di massa di probabilità).

$$p_X(x) = \mathbb{P}(X = x) = \binom{n}{x} \cdot p^x (1-p)^{n-x} \cdot \mathbb{1}_{R_X}(x) \quad (5.18)$$

con: x è il numero di successi, n è il numero di esperimenti, p probabilità di successo in ogni esperimento.

Remark 151. Nella 5.18 la prima parte (il coefficiente binomiale) serve per quantificare il numero di casi in cui si verificano il numero di successi desiderati; questa viene moltiplicata per la seconda che costituisce la probabilità di un tale esito (determinato come probabilità di eventi indipendenti di successo/insuccesso).

Definition 5.4.3 (Funzione di ripartizione).

$$F_X(x) = \mathbb{P}(X \leq x) = \sum_{k=0}^x \binom{n}{k} \cdot p^k (1-p)^{n-k}$$

Validità PMF. Si ha che

$$\sum_{x=0}^n p_X(x) = \sum_{x=0}^n \binom{n}{x} p^x (1-p)^{n-x} \stackrel{(1)}{=} (p + (1-p))^n = 1$$

dove in (1) si è sfruttata la proprietà del coefficiente binomiale:

$$(a+b)^n = \sum_{x=0}^n \binom{n}{x} a^x b^{n-x}$$

□

5.4.3 Moments

Proposition 5.4.1 (Momenti caratteristici).

$$\mathbb{E}[X] = np \quad (5.19)$$

$$\text{Var}[X] = np(1-p) \quad (5.20)$$

$$\text{Asym}(X) = \frac{1-2p}{\sqrt{np(1-p)}} \quad (5.21)$$

$$\text{Kurt}(X) = 3 + \frac{1-6p+6p^2}{np(1-p)} \quad (5.22)$$

Dimostrazione. Per il valore atteso, sfruttando il fatto che $X \sim \text{Bin}(n, p)$ sia descrivibile come la somma di n vc $X_i \sim \text{Bern}(p)$, sfruttando la linearità del valore atteso, il risultato è la somma di n valori attesi uguali:

$$\mathbb{E}[X] = \mathbb{E}\left[\sum_{i=1}^n X_i\right] = \sum_{i=1}^n \mathbb{E}[X_i] = n \mathbb{E}[X_i] = np$$

Alternativamente potevamo sviluppare l'algebra:

$$\begin{aligned} \mathbb{E}[X] &= \sum_{x=0}^n x \cdot \binom{n}{x} p^x (1-p)^{(n-x)} = \sum_{x=0}^n x \cdot \frac{n!}{x!(n-x)!} p^x (1-p)^{(n-x)} \\ &= \sum_{x=0}^n x \cdot \frac{n(n-1)!}{x(x-1)![(n-1)-(x-1)]!} p p^{x-1} (1-p)^{[(n-1)-(x-1)]} \end{aligned}$$

Ora dato che per $x=0$ il termine entro sommatoria è nullo possiamo portare avanti di uno l'indice inferiore della stessa:

$$\mathbb{E}[X] = \sum_{x=1}^n x \cdot \frac{n(n-1)!}{x(x-1)![(n-1)-(x-1)]!} p p^{x-1} (1-p)^{[(n-1)-(x-1)]}$$

ponendo $y = x - 1$ si giunge

$$\begin{aligned}\mathbb{E}[X] &= np \sum_{y=0}^{n-1} \underbrace{\frac{(n-1)!}{y![(n-1)-y]!} p^y (1-p)^{[(n-1)-y]}}_{\text{Bin}(n-1, p)} \\ &\stackrel{(1)}{=} np\end{aligned}$$

con (1) dato che la sommatoria è 1. \square

Dimostrazione. Sfruttando sempre il fatto che $X \sim \text{Bin}(n, p)$ sia descrivibile come la somma di n vc iid $X_i \sim \text{Bern}(p)$, con varianza comune $p(1-p)$, e applicando le proprietà della varianza:

$$\text{Var}[X] = \text{Var}\left[\sum_{i=1}^n X_i\right] \stackrel{(1)}{=} \sum_{i=1}^n \text{Var}[X_i] = n \text{Var}[X_i] = n \cdot p(1-p) \quad (5.23)$$

where in (1) there's no covariance since they are independent. \square

5.4.4 Shape

Proposition 5.4.2 (Shape). *La distribuzione (figura 5.1) è*

- *simmetrica se $p = 0.5$, ossia nel caso $p_X(x) = p_X(n-x)$*
- *asimmetrica positiva (grosso della distribuzione nella parte bassa e coda a destra) se $p < 0.5$,*
- *asimmetrica negativa (gross nella parte alta e coda a sinistra) se $p > 0.5$.*

Dimostrazione. Per $p = 0.5$ è simmetrica in quanto $p = 1 - p = \frac{1}{2}$ e

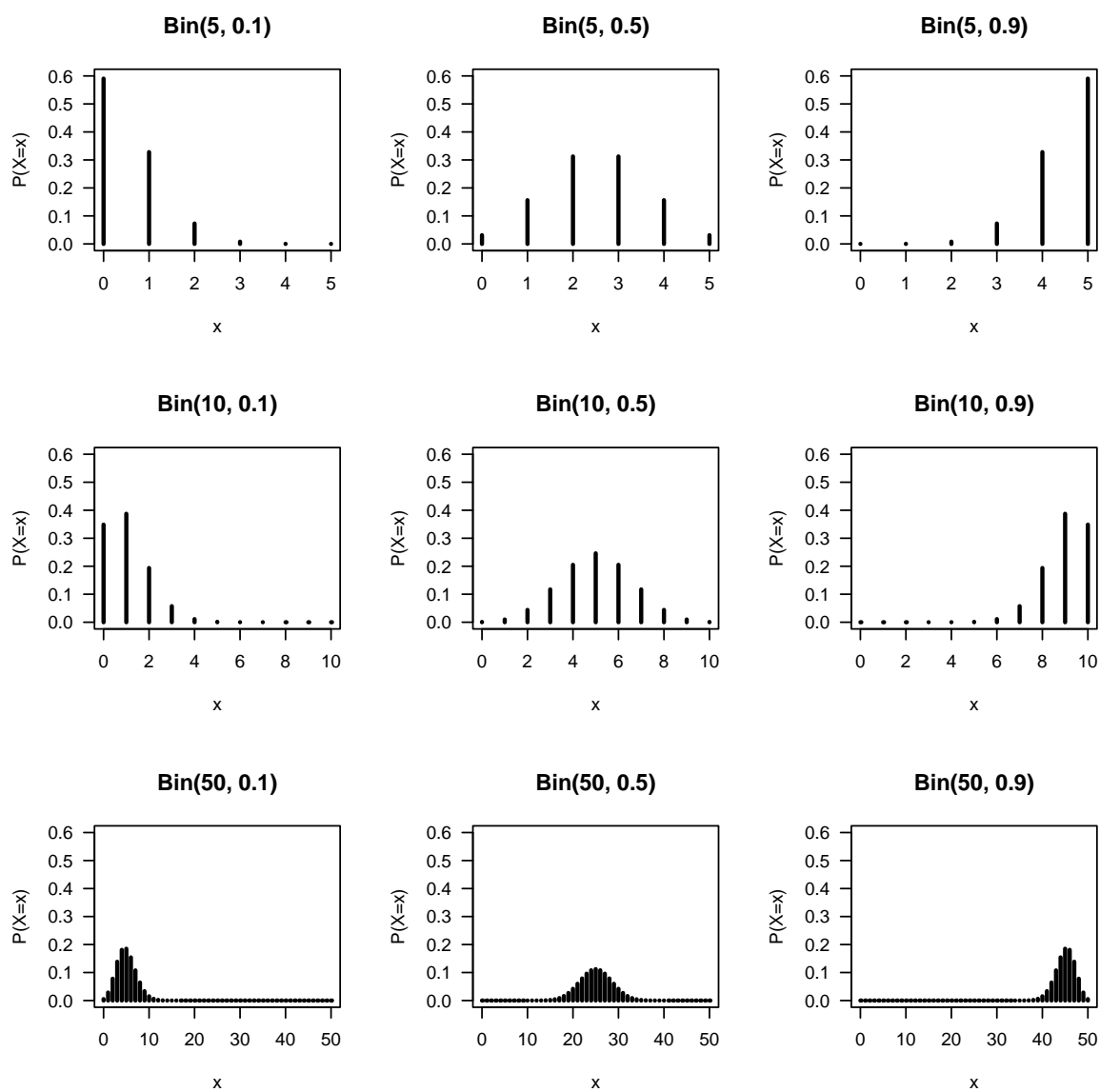
$$p_X(x) = \binom{n}{x} \left(\frac{1}{2}\right)^x \left(\frac{1}{2}\right)^{n-x} = p_X(n-x) = \binom{n}{n-x} \left(\frac{1}{2}\right)^{n-x} \left(\frac{1}{2}\right)^x \quad (5.24)$$

per le proprietà del coefficiente binomiale. E dato che $p_X(x) = p_X(n-x)$, $\forall x \in R_X$, allora la distribuzione è simmetrica attorno al centro del supporto. \square

Proposition 5.4.3. *In una binomiale di parametri n, p , la funzione di densità (per x che varia da 0 a n) è inizialmente strettamente crescente e successivamente strettamente decrescente. Si raggiunge il massimo in corrispondenza del più grande intero $x \leq (n+1)p$*

Dimostrazione. Consideriamo il rapporto $\mathbb{P}(X=x)/\mathbb{P}(X=x-1)$ e determiniamo per quali valori di x esso risulti maggiore (funzione crescente) o minore (decrescente) di 1:

$$\frac{\mathbb{P}(X=x)}{\mathbb{P}(X=x-1)} = \frac{\frac{n!}{(n-x)!x!} p^x (1-p)^{n-x}}{\frac{n!}{(n-x+1)!(x-1)!} p^{x-1} (1-p)^{n-x+1}} = \frac{(n-x+1)p}{x(1-p)}$$

Figura 5.1: Forma distribuzione $\text{Bin}(n, p)$

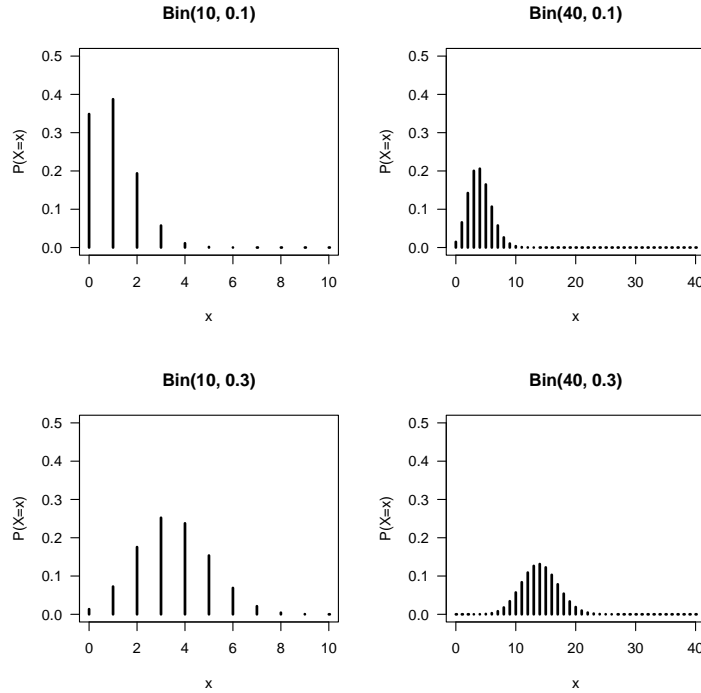


Figura 5.2: Convergenza alla normale della binomiale

Quindi tale rapporto ≥ 1 se e solo se:

$$\begin{aligned}(n-x+1)p &\geq x(1-p) \\ np - xp + p &\geq x - xp\end{aligned}$$

ossia $x \leq (n+1)p$

□

Remark 152 (Convergenza alla normale). La distribuzione converge verso la Normale (diviene simmetrica e la curtosi tende a 3) al crescere di $n \rightarrow \infty$; la convergenza è tanto più veloce per quanto più p è prossimo a 0.5. (figura 5.2)

5.4.5 Extras

5.4.5.1 Generazione mediante somma di bernoulliane

Proposition 5.4.4. *La binomiale può essere generata sommando bernoulliane iid; se X_i , $i = 1, \dots, n$ sono vc bernoulliane iid $X_i \sim \text{Bern}(p)$ allora la loro somma $X = \sum_{i=1}^n X_i \sim \text{Bin}(n, p)$*

Dimostrazione. Sia $X_i = 1$ se l' i -esimo trial ha successo o 0 in caso contrario. Se pensiamo di avere una persona per ciascun trial, chiediamo di alzare la mano se si ha successo e contiamo le mani alzate (che equivale a sommare X_i) otteniamo il numero totale di successi in n trial che è X . □

5.4.5.2 Somma di binomiali

Remark 153. Un fatto importante della binomiale è che la somma di binomiali indipendenti aventi la stessa probabilità di successo è un'altra binomiale

Proposition 5.4.5 (Somma di binomiali). *Se $X \sim \text{Bin}(n, p)$, $Y \sim \text{Bin}(m, p)$, $X \perp\!\!\!\perp Y$, allora $X + Y \sim \text{Bin}(n + m, p)$*

Dimostrazione. Un modo semplice è rappresentare X e Y come le somma di $X = X_1 + \dots + X_n$ e $Y = Y_1 + \dots + Y_m$ con $X_i, Y_i \sim \text{Bern}(p)$ iid. Allora $X + Y$ è la somma di $n + m$ $\text{Bern}(p)$ iid, pertanto la distribuzione è $\text{Bin}(n + m, p)$ per teorema 5.4.4.

Alternativamente, mediante la legge delle probabilità totali, possiamo trovare la PMF di $X + Y$ condizionando su X (oppure ugualmente su Y) e sommando:

$$\begin{aligned}
 \mathbb{P}(X + Y = k) &= \sum_{j=0}^k \mathbb{P}(X + Y = k | X = j) \cdot \mathbb{P}(X = j) \\
 &= \sum_{j=0}^k \mathbb{P}(Y = k - j | X = j) \cdot \mathbb{P}(X = j) \\
 &\stackrel{(1)}{=} \sum_{j=0}^k \mathbb{P}(Y = k - j) \cdot \mathbb{P}(X = j) \\
 &= \sum_{j=0}^k \binom{m}{k-j} p^{k-j} (1-p)^{m-k+j} \cdot \binom{n}{j} p^j (1-p)^{n-j} \\
 &= p^k (1-p)^{n+m-k} \sum_{j=0}^k \binom{m}{k-j} \binom{n}{j} \\
 &\stackrel{(2)}{=} \binom{n+m}{k} p^k (1-p)^{n+m-k} = \text{Bin}(n+m, p)
 \end{aligned}$$

dove in (1) abbiamo sfruttato l'indipendenza tra X e Y e in (2) l'identità di Vandermonde (eq 2.15). \square

5.4.5.3 Variabili derivate

Definition 5.4.4 (Vc frequenza relativa). A partire dalla distribuzione binomiale si può definire la vc frequenza relativa come il rapporto tra la vc X e il numero n complessivo delle sottoprove $Y = X/n$. Essa possiede valori medi e varianza pari a

$$\begin{aligned}
 \mathbb{E}[Y] &= \mathbb{E}[X]/n = p \\
 \text{Var}[Y] &= \frac{\text{Var}[X]}{n^2} = \frac{p(1-p)}{n}
 \end{aligned}$$

Proposition 5.4.6 (Vc numero di insuccessi). *Sia $X \sim \text{Bin}(n, p)$. Allora $n - X \sim \text{Bin}(n, 1 - p)$.*

Dimostrazione. Ad intuito basta invertire i ruoli di successo e insuccesso (si inverte anche la probabilità). Volendo tuttavia verificare, sia $Y = n - X$, la

PMF è:

$$\begin{aligned}\mathbb{P}(Y = x) &\stackrel{(1)}{=} \mathbb{P}(X = n - x) = \binom{n}{n-x} p^{n-x} (1-p)^x \\ &\stackrel{(1)}{=} \binom{n}{x} (1-p)^x p^{n-x} = \text{Bin}(n, 1-p)\end{aligned}$$

dove in (1) diciamo che in n estrazioni la probabilità di avere x fallimenti è uguale alla probabilità di avere $n - x$ successi, mentre in (2) abbiamo sfruttato la proprietà del coefficiente binomiale. \square

5.5 Hypergeometric

5.5.1 Definition

Remark 154. La variabile ipergeometrica descrive l'estrazione *senza reinserimento* di palline dicotomiche da un'urna. A differenza della binomiale dove la probabilità di successo p non cambiava da una sottoprova Bernoulliana all'altra, qui il *non reinserimento* fa sì che la probabilità di successo vari ad ogni prova.

Definition 5.5.1 (Distribuzione ipergeometrica). Supponiamo di dover estrarre un campione di n palline senza reinserimento da un'urna che contiene w palline bianche (successo) e b nere. Il numero X di palline bianche (successi) tra le estratte si distribuisce come una ipergeometrica con parametri w , b ed n e si scrive $X \sim \text{HGeom}(w, b, n)$.

5.5.2 Functions

Remark 155 (Supporto e spazio parametrico).

$$\begin{aligned}R_X &= \{0, 1, \dots, n\} \\ \Theta &= \{w, b \in \mathbb{N} : w + b \geq 1; n \in \{0, \dots, w + b\}\}\end{aligned}$$

Remark 156. Come si nota il supporto è lo stesso della binomiale (dalla quale cambia solo il modo di estrazione e dunque le probabilità generate).

Definition 5.5.2 (Funzione di massa di probabilità).

$$p_X(x) = \mathbb{P}(X = x) = \frac{\binom{w}{x} \binom{b}{n-x}}{\binom{w+b}{n}} \cdot \mathbb{1}_{R_X}(x) \quad (5.25)$$

Remark 157 (Interpretazione). Al denominatore sono quantificati il numero di modi con cui posso estrarre n palline qualsiasi dall'urna. Di queste estrazioni, al numeratore sono quantificati il numero di modi in cui nelle n palline estratte ci sono x bianche (successi); ossia devo averne x bianche scelte tra w ($\binom{w}{x}$ modi per farlo), e $n - x$ nere scelte tra b ($\binom{b}{n-x}$ modi).

Remark 158 (Binomiale e ipergeometrica). Come è intuitivo, se il numero di palline estratte n è molto vicino al numero di palline totali nell'urna $w + b$ la binomiale e l'ipergeometrica forniscono risultati molto diversi.

Tuttavia se n è molto più piccolo di $w + b$ binomiale e ipergeometrica danno risultati simili: nella pratica le estrazioni sono quasi sempre senza reimmissione ma ciò nonostante si usa la binomiale al posto dell'ipergeometrica perché l'errore che si commette è trascurabile.

Validità PMF. Facendo la somma del numeratore si ha:

$$\sum_{x=0}^n \binom{w}{x} \binom{b}{n-x} \stackrel{(1)}{=} \binom{w+b}{n}$$

con (1) per l'identità di Vandermonde (eq 2.15), per cui la PMF somma a 1. \square

Remark 159. In R per la PMF si usa `dhyper(x, m, n, k)` dove \mathbf{x} è il supporto (ossia il numero di palline bianche estratte), \mathbf{m} il numero di palline bianche nell'urna, \mathbf{n} il numero di palline nere e \mathbf{k} il numero di estrazioni.

Remark 160. L'argomento che abbiamo usato per l'ipergeometrica si adatta immediatamente al seguente caso più generale (equivalente estensione dalla binomiale alla multinomiale per la ipergeometrica)

Definition 5.5.3 (Definizione di Soffritti). Cambia la parametrizzazione, con n numero di estrazioni effettuate, N numero di palline totali nell'urna e M numero di "palline successo".

Una variabile casuale X con supporto $R_X = \{0, 1, \dots, n\}$ (stesso supporto della binomiale) ha distribuzione ipergeometrica, e scriviamo $X \sim \text{HGeom}(n, M, N)$ con n, M, N interi positivi con $n \leq N$ e $M \leq N$ se la sua funzione di massa di probabilità è

$$p_X(x) = \frac{\binom{M}{x} \binom{N-M}{n-x}}{\binom{N}{n}} \mathbb{1}_{R_X}(x)$$

Definition 5.5.4 (Multivariate Hypergeometric distribution). Supponiamo di avere una urna contenente palline di k colori, in numerosità n_1, \dots, n_k e di effettuare $n \leq \sum_i n_i$ estrazioni senza reimmissione. In tal caso la probabilità di estrarre j_1 palline del primo colore, \dots j_k palline del k -esimo è

$$\mathbb{P}(j_1, \dots, j_k) = \frac{\binom{n_1}{j_1} \dots \binom{n_k}{j_k}}{\binom{\sum_i n_i}{\sum_i j_i}}$$

5.5.3 Moments

Proposition 5.5.1 (Momenti caratteristici).

$$\mathbb{E}[X] = n \frac{w}{w+b} \tag{5.26}$$

$$\text{Var}[X] = np(1-p) \left(\frac{w+b-n}{w+b-1} \right), \quad \text{con } p = \frac{w}{w+b} \tag{5.27}$$

Dimostrazione. Per il valore atteso, come nel caso binomiale possiamo scrivere X come somma di Bernoulliane $I_i \sim \text{Bern}(p)$ con $p = w/(w+b)$.

$$X = I_1 + \dots + I_n$$

A differenza della binomiale le I_i non sono indipendenti, tuttavia la linearità del valore atteso non lo richiede, quindi

$$\mathbb{E}[X] = \mathbb{E}[I_1 + \dots + I_n] = \mathbb{E}[I_1] + \dots + \mathbb{E}[I_n] = np = n \frac{w}{w+b}$$

□

Dimostrazione. Per la varianza invece essendo variabili non indipendenti non possiamo sommare le varianze direttamente. Vedremo in seguito la dimostrazione della formula riportata. □

5.5.4 Extras

5.5.4.1 Esperimenti assimilabili

Remark 161. L'idea dell'Ipergeometrica è classificare una popolazione utilizzando due set di tag consecutivi (entrambi dicotomici successo/insuccesso) e ottenere il numero degli elementi caratterizzati dal successo in entrambi i tag. Nell'esempio delle palline il primo tag è il colore della pallina (bianco = successo), mentre il secondo è estrazione (estratta = successo).

Problemi aventi la stessa struttura presenteranno medesima distribuzione.

Example 5.5.1. Il numero A di assi estratti (sono 4 in un mazzo di 52 carte) in una mano di poker (5 carte estratte) si distribuirà come $A \sim \text{HGeom}(4, 48, 5)$.

Remark 162. La struttura essenziale ci permette di dimostrare facilmente l'uguaglianza di due ipergeometriche dove l'ordine dei set di tag viene invertito

Proposition 5.5.2. $\text{HGeom}(w, b, n)$ e $\text{HGeom}(n, w+b-n, w)$ sono identiche.

Dimostrazione. Sia $X \sim \text{HGeom}(w, b, n)$ è il numero di palline bianche tra le estratte campione; sia $Y \sim \text{HGeom}(n, w+b-n, w)$ il numero di palline estratte tra le bianche (pensando ad estratto/non estratto come il primo tag e al colore come secondo). Entrambe X, Y contano il numero di bianche estratte pertanto avranno la stessa distribuzione.

Alternativamente possiamo controllare algebricamente che

$$\begin{aligned} \mathbb{P}(X=x) &= \frac{\binom{w}{x} \binom{b}{n-x}}{\binom{w+b}{n}} = \frac{\frac{w!}{x!(w-x)!} \frac{b!}{(n-x)!(b-n+x)!}}{\frac{(w+b)!}{n!(w+b-n)!}} = \frac{w!b!n!(w+b-n)!}{k!(w-k)!(n-k)!(b-n+k)!} \\ \mathbb{P}(Y=y) &= \frac{\binom{n}{y} \binom{w+b-n}{w-y}}{\binom{w+b}{w}} = \frac{\frac{n!}{y!(n-y)!} \frac{(w+b-n)!}{(w-y)!(b-n+y)!}}{\frac{(w+b)!}{w!b!}} = \frac{w!b!n!(w+b-n)!}{k!(w-k)!(n-k)!(b-n+k)!} \end{aligned}$$

e dunque $\mathbb{P}(X=x) = \mathbb{P}(Y=y)$. □

5.5.4.2 Connessioni con la binomiale

Remark 163. Binomiale ed ipergeometrica sono connesse: possiamo ottenere la binomiale calcolando un limite sull'ipergeometrica, oppure ottenere una ipergeometrica condizionando una binomiale.

Dall'ipergeometrica alla binomiale

Proposition 5.5.3. *Se $X \sim \text{HGeom}(w, b, n)$ e $w + b \rightarrow \infty$ ma $p = w/(w + b)$ rimane fisso, allora la PMF di X converge a $\text{Bin}(n, p)$.*

Dimostrazione. Sviluppiamo algebricamente per essere comodi prima di applicare il limite:

$$\mathbb{P}(X = x) = \frac{\binom{w}{x} \binom{b}{n-x}}{\binom{w+b}{n}} \stackrel{(1)}{=} \binom{n}{x} \frac{\binom{w+b-n}{w-x}}{\binom{w+b}{w}}$$

dove in (1) abbiamo sfruttato che $\text{HGeom}(w, b, n) = \text{HGeom}(n, w + b - n, w)$ come nella dimostrazione di 5.5.2. Ora sviluppiamo il rapporto al secondo fattore ricordando che $\binom{n}{d} = \frac{n!}{d!(n-d)!}$; si ha:

$$\begin{aligned} \frac{\binom{w+b-n}{w-x}}{\binom{w+b}{w}} &= \frac{(w+b-n)!}{(w-x)!(w+b-n-w+x)!} : \frac{(w+b)!}{w!(w+b-w)!} \\ &= \frac{(w+b-n)!}{(w-x)!(b-n+x)!} \cdot \frac{w!b!}{(w+b)!} \\ &= \frac{w!}{(w-x)!} \frac{b!}{(b-n+x)!} \frac{(w+b-n)!}{(w+b)!} \\ &= \frac{w \cdot \dots \cdot (w-x+1)(w-x)!}{(w-x)!} \frac{b \cdot \dots \cdot (b-n+x+1)(b-n+x)!}{(b-n+x)!} \frac{(w+b-n)!}{(w+b) \cdot \dots \cdot (w+b-n+1)} \\ &= \frac{w \cdot \dots \cdot (w-x+1)}{1} \frac{b \cdot \dots \cdot (b-n+x+1)}{1} \frac{1}{(w+b) \cdot \dots \cdot (w+b-n+1)} \end{aligned}$$

ora al numeratore del primo rapporto abbiamo $w - (w - x + 1) + 1 = x$ fattori, al numeratore del secondo ne abbiamo $b - (b - n + x + 1) + 1 = n - x$ elementi. Pertanto complessivamente al numeratore abbiamo n fattori. Al denominatore invece abbiamo $(w + b) - (w + b - n + 1) + 1 = n$ fattori anche qui. Pertanto possiamo dividere per $(w + b)$, applicandolo n volte sia al numeratore che al denominatore, ottenendo

$$\frac{\binom{w+b-n}{w-x}}{\binom{w+b}{w}} = \frac{\frac{w}{w+b} \cdot \dots \cdot \left(\frac{w}{w+b} - \frac{x-1}{w+b}\right) \cdot \left(\frac{b}{w+b}\right) \cdot \dots \cdot \left(\frac{b}{w+b} - \frac{n-x-1}{w+b}\right)}{1 \cdot \dots \cdot \left(1 - \frac{n-1}{w+b}\right)}$$

ora sostituendo $p = \frac{w}{w+b}$, $1 - p = \frac{b}{w+b}$ e al denominatore $w + b = N$ dove utile si ha:

$$\frac{\binom{w+b-n}{w-x}}{\binom{w+b}{w}} = \frac{p \cdot \dots \cdot \left(p - \frac{x-1}{N}\right) \cdot (1-p) \cdot \dots \cdot \left(1 - p - \frac{n-x-1}{N}\right)}{\left(1 - \frac{1}{N}\right) \cdot \dots \cdot \left(1 - \frac{n-1}{N}\right)}$$

Ora tornando da dove siamo partiti abbiamo:

$$\mathbb{P}(X = x) = \binom{n}{x} \frac{p \cdot \dots \cdot \left(p - \frac{x-1}{N}\right) \cdot (1-p) \cdot \dots \cdot \left(1 - p - \frac{n-x-1}{N}\right)}{\left(1 - \frac{1}{N}\right) \cdot \dots \cdot \left(1 - \frac{n-1}{N}\right)}$$

Infine per $N \rightarrow +\infty$ il denominatore va a 1 mentre il numeratore va a $p^x(1-p)^{n-x}$ pertanto

$$\mathbb{P}(X = x) \rightarrow \binom{n}{x} p^x (1-p)^{n-x}$$

che è la $\text{Bin}(n, p)$.

Intuitivamente data un'urna con w palline bianche e b nere, la binomiale sorge dall'estrarre n palline con replacement, mentre l'ipergeometrica senza. Se il numero di palline nell'urna sale notevolmente rispetto al numero di palline estratte, il campionamento con ripetizione e senza diventano essenzialmente equivalenti. (l'estrazione di una pallina non cambia la probabilità delle prossime estrazioni perché data la grande numerosità nell'urna non modifica praticamente la probabilità di successo) \square

Remark 164. In termini pratici il teorema ci dice che se $N = w + b$ è grande rispetto a n possiamo approssimare la PMF di $\text{HGeom}(w, b, n)$ con $\text{Bin}(n, w/(w+b))$.

Dalla binomiale all'ipergeometrica

Proposition 5.5.4. *Se $X \sim \text{Bin}(n, p)$, $Y \sim \text{Bin}(m, p)$ e X è indipendente da Y , allora la distribuzione condizionata di X dato che $X + Y = r$ è $\text{HGeom}(n, m, r)$*

Remark 165. Dimostriamo attraverso un esempio (distribuzione del test esatto di Fisher).

Dimostrazione. Un ricercatore vuole studiare se la prevalenza di una data malattia sia uguale o meno tra maschi e femmine. Raccoglie un campione di n donne ed m uomini e testa la malattia. Sia $X \sim \text{Bin}(n, p_1)$ il numero di donne con la malattia nel campione e $Y \sim \text{Bin}(m, p_2)$ il numero di uomini. Qui p_1 e p_2 sono sconosciuti.

Supponiamo che siano osservate $X + Y = r$ persone malate. Siamo interessati a testare se $p_1 = p_2 = p$ (la cd ipotesi nulla); il test di Fisher si fonda sul condizionare sui totali di riga e colonna (quindi n, m, r sono considerati fissi) e verificare se il valore osservato X (numero di donne malate) sia estremo (dato che il tot malati è r) sotto ipotesi nulla. Assumendo l'ipotesi nulla vera troviamo la PMF condizionale di X dato che $X + Y = r$.

La tabella 2×2 di riferimento è la 5.3. Costruiamo PMF condizionata attraverso la regola di Bayes:

$$\begin{aligned} \mathbb{P}(X = x | X + Y = r) &= \frac{\mathbb{P}(X + Y = r | X = x) \mathbb{P}(X = x)}{\mathbb{P}(X + Y = r)} = \frac{\mathbb{P}(Y = r - x | X = x) \mathbb{P}(X = x)}{\mathbb{P}(X + Y = r)} \\ &\stackrel{(1)}{=} \frac{\mathbb{P}(Y = r - x) \mathbb{P}(X = x)}{\mathbb{P}(X + Y = r)} \end{aligned}$$

dove in (1) abbiamo sfruttato l'indipendenza di X e Y . Assumendo per buona l'ipotesi nulla e impostando $p_1 = p_2 = p$ si hanno le vc indipendenti $X \sim \text{Bin}(n, p)$ e $Y \sim \text{Bin}(m, p)$, per cui $X + Y \sim \text{Bin}(n + m, p)$ (per il risultato

	Donne	Uomini	Tot
Malato	x	$r - x$	r
Sano	$n - x$	$m - r + x$	$n + m - r$
Tot	n	m	$n + m$

Tabella 5.3

5.4.5). Pertanto sostituendo le formule per esteso si ha

$$\begin{aligned}\mathbb{P}(X = x | X + Y = r) &= \frac{\binom{m}{r-x} p^{r-x} (1-p)^{m-r+x} \cdot \binom{n}{x} p^x (1-p)^{n-x}}{\binom{n+m}{r} p^r (1-p)^{n+m-r}} \\ &= \frac{\binom{n}{x} \binom{m}{r-x}}{\binom{n+m}{r}} = \text{HGeom}(n, m, r)\end{aligned}$$

Intuitivamente questo avviene perché condizionatamente ad avere $X + Y = r$ malati (primo tag), X è il numero di donne (secondo tag) tra quelli. \square

5.6 Geometric (n. of trials)

Important remark 41. Supponiamo di ripetere in maniera indipendente diverse prove bernoulliane, ciascuna avente p probabilità di successo, sino a che si verifica il primo successo. La Geometrica può essere definita come

- X il numero di prove *fallimentari* necessari per ottenere il primo successo;
- X il numero di *trial complessivi*, incluso il primo successo che nel seguito è chiamata “First success distribution” (Soffritti considera questa)

In entrambi i casi si scrive che X si distribuisce come una variabile geometrica con parametro p e si scrive $X \sim \text{Geom}(p)$ ma è necessario capire a quale conteggio ci si stia riferendo, in quanto ha effetto su

- supporto: il numero di prove fallimentari include lo 0 (si ha subito successo) mentre il numero di trial complessivi parte da 1
- PMF (dove cambia l’esponente di $(1-p)$)

Remark 166. Situazioni reali per la geometrica: vogliamo studiare un evento associato alla ripetizione di prove bernoulliane indipendenti (similmente alla binomiale). Differentemente dalla binomiale, invece che fissare il numero di repliche n e chiederci il numero di successi x , ci chiediamo il numero di repliche (o insuccessi) necessari per avere il primo successo.

Se

- ho la sequenza [successo] (al primo colpo) il numero di interesse è rispettivamente 0 (numero di prove fallimentari prima del primo successo) o 1 (numero di trial complessivi sino al primo successo)
- ho la sequenza [insuccesso, successo] è rispettivamente 1 (n fallimenti) o 2 (n prove)

Tutti questi vettori hanno in comune che nell’ultima posizione si ha successo; si differenziano per la lunghezza e questo è il risultato della variabile casuale

5.6.1 Definition

Remark 167. Se definiamo X come il numero di *prove* necessarie per ottenere il primo successo (incluso quest'ultimo). Qui la chiamiamo first success distribution e la indichiamo con $X \sim \text{FS}(p)$

5.6.2 Functions

Remark 168 (Supporto e spazio parametrico).

$$\begin{aligned} R_X &= \{1, 2, 3, \dots\} \\ \Theta &= \{p \in (0, 1)\} \end{aligned}$$

Definition 5.6.1 (Funzione di massa di probabilità).

$$p_X(x) = \mathbb{P}(X = x) = (1 - p)^{x-1} p \cdot \mathbb{1}_{R_X}(x) \quad (5.28)$$

Remark 169 (Interpretazione). La probabilità di avere il primo successo all' n -esima estrazione è data dalla probabilità di $n - 1$ fallimenti per la probabilità di un successo.

Definition 5.6.2 (Funzione di ripartizione).

$$\begin{aligned} F_X(x) &= \mathbb{P}(X \leq x) = \sum_{k=1}^x \mathbb{P}(X = k) = \sum_{k=1}^x (1 - p)^{k-1} p \\ &= 1 - (1 - p)^x \end{aligned} \quad (5.29)$$

5.6.3 Moments

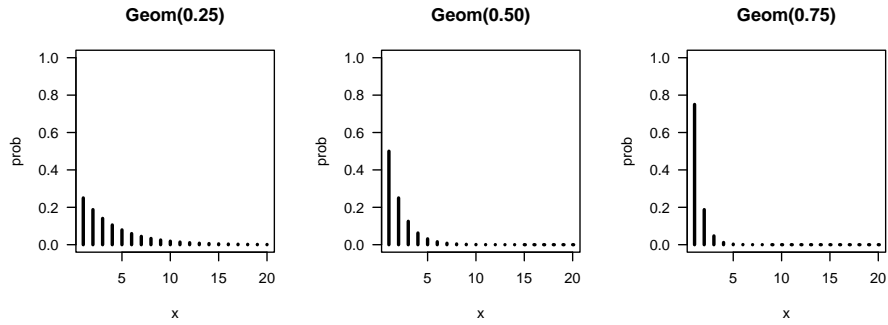
Proposition 5.6.1 (Momenti caratteristici).

$$\begin{aligned} \mathbb{E}[X] &= \frac{1}{p} \\ \text{Var}[X] &= \frac{1 - p}{p^2} \\ \text{Asym}(X) &= \frac{2 - p}{\sqrt{1 - p}} \\ \text{Kurt}(X) &= 9 + \frac{p^2}{1 - p} \end{aligned}$$

Dimostrazione. Sia $Y = X + 1 \sim \text{FS}(p)$ con $X \sim \text{Geom}(p)$. Allora sfruttando le conoscenze sulla geometrica e le proprietà di valore atteso e varianza

$$\begin{aligned} \mathbb{E}[Y] &= \mathbb{E}[X + 1] = \mathbb{E}[X] + 1 = \frac{1 - p}{p} + 1 = \frac{1}{p} \\ \text{Var}[Y] &= \text{Var}[X + 1] = \text{Var}[X] = \frac{1 - p}{p^2} \end{aligned}$$

□

Figura 5.3: Forma distribuzione Geom(p)

5.6.4 Shape

Remark 170 (Shape). Le geometriche (fig 5.3):

- sono decrescenti, con probabilità più alte associate ai valori più piccoli di x (asimmetria positiva).
- l'asimmetria positiva aumenta al crescere di p (più p è alto più facile stoppare l'esperimento presto e la PMF discende verso 0 velocemente).

5.6.5 Extras

5.6.5.1 Assenza di memoria

Proposition 5.6.2 (Assenza di memoria). *Si ha che la geometrica non ha memoria ossia $\mathbb{P}(X > t + s | X > t) = \mathbb{P}(X > s)$.*

Dimostrazione. Si ha:

$$\begin{aligned}
 \mathbb{P}(X > t + s | X > t) &= \frac{\mathbb{P}(X > t + s)}{\mathbb{P}(X > t)} = \frac{1 - F_X(t + s)}{1 - F_X(t)} \\
 &= \frac{(1 - p)^{t+s}}{(1 - p)^t} = (1 - p)^s \\
 &= \mathbb{P}(X > s)
 \end{aligned}$$

ovvero il ritardo accertato di un evento in t sottoprove indipendenti non modifica la probabilità che esso si verifichi entro ulteriori s sottoprove. \square

5.7 Geometric (n. of failures)

Remark 171. Qui definiamo X come il numero di *fallimenti* necessarie per ottenere il primo successo.

5.7.1 Functions

Remark 172 (Supporto e spazio parametrico).

$$\begin{aligned} R_X &= \{0, 1, 2, \dots\} \\ \Theta &= \{p \in (0, 1)\} \end{aligned}$$

Definition 5.7.1 (Funzione di massa di probabilità).

$$p_X(x) = \mathbb{P}(X = x) = (1-p)^x p \cdot \mathbb{1}_{R_X}(x) \quad (5.30)$$

Example 5.7.1. Il numero di croci ottenute sino alla prima testa verificatasi si distribuisce come $\text{Geom}(1/2)$.

Validità PMF. Si ha che

$$\sum_{x=0}^{\infty} (1-p)^x p = p \sum_{x=0}^{\infty} (1-p)^x \stackrel{(1)}{=} p \cdot \frac{1}{p} = 1$$

con l'uguaglianza (1) dovuta alla serie geometrica. \square

Remark 173. Come il teorema binomiale mostra che la PMF binomiale sia valida, la serie geometrica mostra che la PMF Geometrica sia valida.

Remark 174 (Interpretazione). La probabilità di avere x fallimenti consecutivi seguiti da un successo è data dalla probabilità di x fallimenti per la probabilità di un successo.

Definition 5.7.2 (Funzione di ripartizione). Si ha

$$F_X(x) = \mathbb{P}(X \leq x) = 1 - (1-p)^{x+1} \quad (5.31)$$

Derivazione della CDF. Si ha

$$F_X(x) = \mathbb{P}(X \leq x) = 1 - \mathbb{P}(X > x) = 1 - \sum_{k=x+1}^{\infty} (1-p)^k p$$

Espandendo la sommatoria:

$$\begin{aligned} \sum_{k=x+1}^{\infty} (1-p)^k p &= (1-p)^{x+1} \cdot p + (1-p)^{x+2} \cdot p + \dots + (1-p)^{\infty} \cdot p \\ &= p(1-p)^x [(1-p) + (1-p)^2 + \dots + (1-p)^{\infty}] \\ &= p(1-p)^x \left[\sum_{i=1}^{\infty} (1-p)^i \right] \\ &= p(1-p)^x \left[\sum_{i=0}^{\infty} (1-p)^i - 1 \right] \\ &= p(1-p)^x \left(\frac{1}{p} - 1 \right) = p(1-p)^x \frac{1-p}{p} \\ &= (1-p)^{x+1} \end{aligned}$$

Pertanto:

$$F_X(x) = 1 - (1-p)^{x+1}$$

\square

5.7.2 Moments

Proposition 5.7.1 (Momenti caratteristici).

$$\begin{aligned}\mathbb{E}[X] &= \frac{1-p}{p} \\ \text{Var}[X] &= \frac{1-p}{p^2}\end{aligned}$$

Dimostrazione. Per il valore atteso abbiamo

$$\mathbb{E}[X] = \sum_{x=0}^{\infty} x \cdot (1-p)^x p$$

Non può essere ricondotta a serie geometrica direttamente per la presenza entro sommatoria di x come primo fattore. Ma notiamo che il termine entro sommatoria assomiglia a $x(1-p)^{x-1}$ ossia la derivata di $(1-p)^x$ rispetto a $1-p$, quindi partiamo da lì:

$$\sum_{x=0}^{\infty} (1-p)^x = \frac{1}{p}$$

Questa serie converge dato che $0 < p < 1$. Derivando entrambi i membri rispetto a p .

$$\begin{aligned}\sum_{x=0}^{\infty} x(1-p)^{x-1} \cdot (-1) &= -\frac{1}{p^2} \\ \sum_{x=0}^{\infty} x(1-p)^{x-1} &= \frac{1}{p^2}\end{aligned}$$

e se moltiplichiamo entrambi i lati per $p(1-p)$ otteniamo la somma dalla quale siamo partiti

$$\begin{aligned}p(1-p) \sum_{x=0}^{\infty} x(1-p)^{x-1} &= \frac{1}{p^2} p(1-p) \\ \sum_{x=0}^{\infty} xp(1-p)^x &= \frac{1-p}{p}\end{aligned}$$

□

Dimostrazione. Per la varianza dobbiamo calcolare $\mathbb{E}[X^2]$:

$$\mathbb{E}[X^2] = \sum_{x=0}^{\infty} x^2 \cdot \mathbb{P}(X=x) = \sum_{x=0}^{\infty} x^2 \cdot (1-p)^x \cdot p \stackrel{(1)}{=} \sum_{x=1}^{\infty} x^2 \cdot (1-p)^x \cdot p$$

con (1) dato dal fatto che se $x=0$ il termine entro sommatoria è nullo e si può portare avanti l'indice della stessa. Anche qui cerchiamo di sfruttare la serie geometrica per arrivare ad una espressione compatta equivalente all'ultimo termine di sopra. La serie è

$$\sum_{x=0}^{\infty} (1-p)^x = \frac{1}{p}$$

Derivando rispetto a p entrambi i membri, come visto in precedenza si ha:

$$\sum_{x=0}^{\infty} x \cdot (1-p)^{x-1} = \frac{1}{p^2}$$

Possiamo portare avanti di 1 l'indice di sommatoria dato che se $x = 0$ è nullo il termine dentro

$$\sum_{x=1}^{\infty} x \cdot (1-p)^{x-1} = \frac{1}{p^2}$$

Ora, derivando ancora si andrebbe a $x(x-1)$ entro sommatoria, invece di x^2 desiderato, pertanto moltiplichiamo per $(1-p)$ entrambi i membri giungendo a:

$$\sum_{x=1}^{\infty} x \cdot (1-p)^x = \frac{1-p}{p^2}$$

Derivando ambo i membri nuovamente rispetto a p si va a

$$\begin{aligned} \sum_{x=1}^{\infty} x^2 \cdot (1-p)^{x-1} \cdot (-1) &= \frac{(-1) \cdot p^2 - 2p \cdot (1-p)}{p^4} \\ \sum_{x=1}^{\infty} x^2 \cdot (1-p)^{x-1} &= (-1) \frac{p^2 - 2p}{p^4} \\ \sum_{x=1}^{\infty} x^2 \cdot (1-p)^{x-1} &= \frac{2-p}{p^3} \end{aligned}$$

Moltiplicando entrambi i membri per $(1-p) \cdot p$ si arriva al punto dove eravamo rimasti con $\mathbb{E}[X^2]$

$$\sum_{x=1}^{\infty} x^2 \cdot (1-p)^x \cdot p = \frac{2-p}{p^3} \cdot (1-p) \cdot p = \frac{(2-p)(1-p)}{p^2}$$

Per cui

$$\mathbb{E}[X] = \sum_{x=1}^{\infty} x \cdot (1-p)^x \cdot p = \frac{(2-p)(1-p)}{p^2}$$

e dunque:

$$\begin{aligned} \text{Var}[X] &= \mathbb{E}[X^2] - (\mathbb{E}[X])^2 = \frac{(2-p)(1-p)}{p^2} - \frac{(1-p)^2}{p^2} \\ &= \frac{(1-p)(2-p-1+p)}{p^2} = \frac{1-p}{p^2} \end{aligned}$$

□

5.7.3 Extras

5.7.3.1 Conversione tra definizioni

Remark 175. Se $Y \sim \text{FS}(p)$ allora $Y-1 \sim \text{Geom}(p)$ e possiamo convertire tra le PMF di Y e $Y-1$ scrivendo

$$\mathbb{P}(Y = k) = \mathbb{P}(Y-1 = k-1)$$

Viceversa se $X \sim \text{Geom}(p)$ allora $X+1 \sim \text{FS}(p)$

5.7.3.2 Assenza di memoria

Remark 176. Una proprietà peculiare della geometrica è di esser l'unica vc discreta senza memoria (a parte la sua riformulazione).

Proposition 5.7.2 (Assenza di memoria).

$$\mathbb{P}(X > t + s | X > t) = \mathbb{P}(X > s) \quad (5.32)$$

Dimostrazione. Si ha:

$$\begin{aligned} \mathbb{P}(X > t + s | X > t) &= \frac{\mathbb{P}(X > t + s)}{\mathbb{P}(X > t)} = \frac{1 - F_X(t + s)}{1 - F_X(t)} = \frac{1 - 1 + (1 - p)^{t+s+1}}{1 - 1 + (1 - p)^{t+1}} \\ &= (1 - p)^s = 1 - F_X(s) = \mathbb{P}(X > s) \end{aligned}$$

□

5.8 Negative binomial (n. of trials)

Remark 177. Generalizza la distribuzione Geometrica.

Similmente alla binomiale si ripetono esperimenti bernoulliani in condizioni costanti (es estrazioni con reinserimento). A differenza della binomiale dove si fissa il numero n delle prove ed è aleatorio/di interesse il numero di successi qui faccio il contrario. Fisso il numero k di successi (non più solo il primo, ossia $k = 1$, diversamente dalla geometrica) dopodiché alternativamente conto:

- il numero di *insuccessi* che occorre fare per ottenere quel certo numero di successi (verriamo in questa sezione)
- il numero delle *prove* che occorre fare per ottenere quel certo numero di successi (sezione 5.8.1)

Anche qui fare attenzione a quale versione (Soffritti preferisce la seconda) perché ha influsso su

- supporto se contiamo il numero di fallimenti possono essere anche 0 mentre col numero di trial debbono essere almeno k
- pmf

Se ad esempio fisso $k = 3$:

- la sequenza [successo, successo, successo] porta a 0 nel primo caso, a 3 con la seconda definizione
- la sequenza [successo, insuccesso, successo, successo] porta 1 nel primo caso, 4 nel secondo

Tutte le sequenze di sopra sono caratterizzate dall'avere tre successi, ma si differenziano per lunghezza (e cosa contano)

5.8.1 Definition

Definition 5.8.1 (Definizione con numero di prove). Il numero di *prove* indipendenti X (ciascuna con probabilità p di essere successo) necessarie per avere $k \geq 1$ successi si distribuisce come una binomiale negativa di parametri k e p , ossia $X \sim \text{Nb}(k, p)$.

5.8.2 Functions

Remark 178 (Supporto e spazio parametrico).

$$\begin{aligned} R_X &= \{k, k+1, \dots\} \\ \Theta &= \{k \in \{1, 2, 3, \dots\}, p \in [0, 1]\} \end{aligned}$$

Definition 5.8.2 (Funzione di massa di probabilità).

$$p_X(x) = \mathbb{P}(X = x) = \binom{x-1}{k-1} p^k (1-p)^{x-k} \cdot \mathbb{1}_{R_X}(x) \quad (5.33)$$

Remark 179 (Interpretazione). La formula deriva dalla considerazione che per ottenere il k -esimo successo nella n -esima prova, ci dovranno essere $k-1$ successi nelle prime $n-1$ prove, la cui probabilità

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

è moltiplicata per la probabilità di un successo nella n -esima, ossia p .

Example 5.8.1. La geometrica (prima versione) è una specializzazione della binomiale negativa (prima versione) ossia

$$\text{Geom}(p) = \text{Nb}(1, p) \quad \text{ovvero } k = 1$$

Ad esempio se $X \sim \text{Geom}(p)$,

$$\mathbb{P}(X = j) = \binom{j-1}{k-1} p^k (1-p)^{j-k} \stackrel{(1)}{=} \binom{j-1}{0} p (1-p)^{j-1} = p(1-p)^{j-1}, \quad j = 1, 2, \dots$$

dove in (1) si è sostituito $k = 1$

5.8.3 Moments

Proposition 5.8.1 (Momenti caratteristici).

$$\begin{aligned} \mathbb{E}[X] &= \frac{k}{p} \\ \text{Var}[X] &= \frac{k(1-p)}{p^2} \\ \text{Asym}(X) &= \frac{2-p}{\sqrt{k(1-p)}} \\ \text{Kurt}(X) &= 3 + \frac{6}{k} + \frac{p^2}{k(1-p)} \end{aligned}$$

5.8.4 Shape

La binomiale negativa (fig 5.4) è asimmetrica positiva:

- se aumentiamo p a parità di k , l'asimmetria positiva è + grande (cresce la probabilità di valori più piccoli)

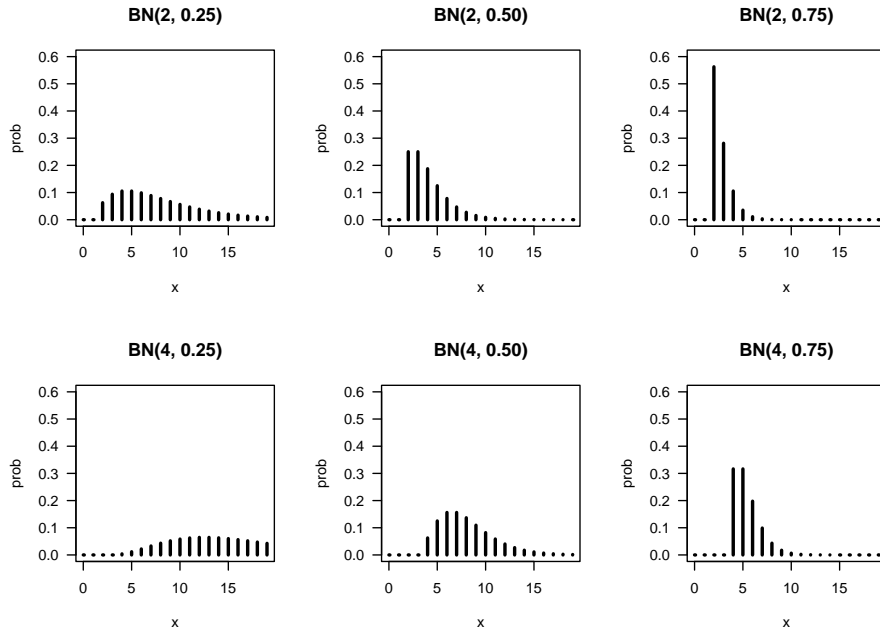


Figura 5.4: Distribuzione binomiale negativa(2)

- se fisso p e aumento k l'asimmetria diminuisce (voglio vedere più successi e la distribuzione di prove si sposta verso l'alto)

Come si nota in questa definizione (a differenza della precedente), non vi è massa di probabilità per valori del dominio inferiori a k .

5.9 Negative binomial (n. of failures)

5.9.1 Definition

Definition 5.9.1 (Definizione con numero di fallimenti). In una sequenza di prove Bernoulliane indipendenti con probabilità di successo p , se X è il numero di *fallimenti* prima del k -esimo successo, allora X ha una distribuzione binomiale negativa con parametri k e p e si scrive $X \sim \text{Nb}(k, p)$

Remark 180. Anche a livello di notazione, nei parametri, si nota subito la differenza con la binomiale: questa fissa il numero di trial mentre la binomiale negativa fissa il numero di successi.

5.9.2 Functions

Remark 181 (Supporto e spazio parametrico).

$$R_X = \{0, 1, 2, \dots\}$$

$$\Theta = \{k \in \mathbb{N} : k \geq 1, p \in \mathbb{R} : 0 \leq p \leq 1\}$$

Definition 5.9.2 (Funzione di massa di probabilità).

$$p_X(x) = \mathbb{P}(X = x) = \binom{x+k-1}{k-1} p^k (1-p)^x \cdot \mathbb{1}_{R_X}(x) \quad (5.34)$$

Remark 182 (Interpretazione). Ci sono $\binom{x+k-1}{k-1}$ sequenze possibili di x fallimenti e $k-1$ successi. Ciascuna di esse ha probabilità $p^{k-1}(1-p)^x$. Si termina con un success, quindi moltiplicando per p .

5.9.3 Moments

Proposition 5.9.1 (Momenti caratteristici).

$$\mathbb{E}[X] = k \frac{1-p}{p} \quad (5.35)$$

$$\text{Var}[X] = k \frac{1-p}{p^2} \quad (5.36)$$

Dimostrazione. Per il valore atteso sfruttiamo che X è scrivibile come somma di k vc Geometriche X_i . Il valore atteso è la somma dei valori attesi delle geometriche:

$$\mathbb{E}[X] = \mathbb{E}[X_1 + \dots + X_k] = \mathbb{E}[X_1] + \dots + \mathbb{E}[X_k] = k \frac{1-p}{p}$$

Per la varianza avviene lo stesso, dato che le variabili sono indipendenti:

$$\text{Var}[X] = \text{Var}[X_1 + \dots + X_k] = \text{Var}[X_1] + \dots + \text{Var}[X_k] = k \frac{1-p}{p^2}$$

□

5.9.4 Shape

Remark 183 (Shape). Si nota che così al crescere di k , la distribuzione diviene più simmetrica e la curtosi tende a 3 indicando convergenza alla normalità. All'aumentare di p assume asimmetria positiva. (figura 5.5)

5.9.5 Extras

5.9.5.1 Generazione mediante somma di geometriche

Remark 184. Come una binomiale può essere rappresentata da una somma di Bernoulli iid, una binomiale negativa può essere rappresentata come somma di Geometriche iid, come mostrato dal seguente teorema.

Proposition 5.9.2. *Sia $X \sim \text{Nb}(k, p)$ il numero di fallimenti prima del k -esimo successo in una sequenza di probe bernoulliane indipendenti con probabilità di successo p . Allora possiamo scrivere $X = X_1 + \dots + X_k$ dove gli X_i sono iid e $X_i \sim \text{Geom}(p)$.*

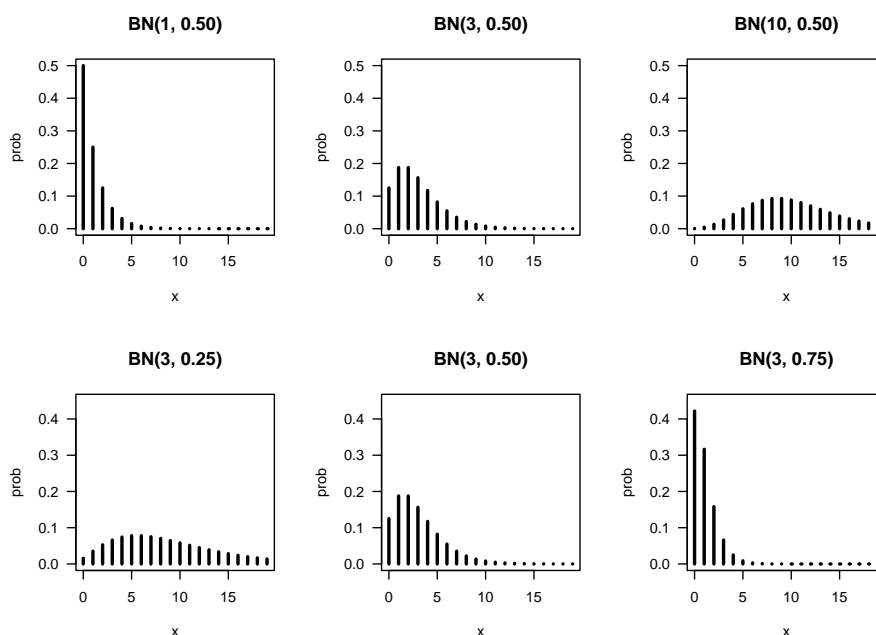


Figura 5.5: Distribuzione binomiale negativa

Dimostrazione. Sia X_1 il numero di fallimenti prima del primo successo, X_2 il numero di fallimenti tra il primo successo e il secondo e, in generale, X_i il numero di fallimenti tra $(i - 1)$ -esimo successo e l' i -esimo.

Allora $X_1 \sim \text{Geom}(p)$ per la definizione della geometrica, $X_2 \sim \text{Geom}(p)$ e così via. Inoltre le X_i sono indipendenti dato che le prove bernoulliane sono indipendenti l'un l'altra. Sommando gli X_i si ottiene il totale di fallimenti prima del k -esimo successo, che è X . \square

5.10 Poisson

5.10.1 Definition

Remark 185. L'esperimento che si descrive mediante la Poisson è il conteggio del numero di eventi di interesse (che si verificano in maniera indipendente tra loro e nelle stesse condizioni sperimentali) occorsi in un finito/prefissato intervallo di tempo o di spazio. Ad esempio:

- visite ad un sito web in un dato intervallo di tempo (es un'ora)
- numero di incidenti in un certo tratto di una autostrada
- numero di errori di battitura di un dattilografo in una pagina

Tutte le volte che si hanno “arrivi” aleatori, e tali arrivi soddisfano certe ipotesi, la probabilità di avere x arrivi in un dato lasso (temporale/spaziale) segue la distribuzione di Poisson $X \sim \text{Pois}(\lambda)$, con parametro $\lambda \in (0, +\infty)$ che indica il numero medio di eventi (centro/punto di maggiore probabilità)

Remark 186. Un risultato che ci servirà per questa distribuzione è il seguente

Proposition 5.10.1 (Sviluppo di Maclaurin della funzione esponenziale).

$$e^x = \sum_{k=0}^{\infty} \frac{x^k}{k!} \quad (5.37)$$

Dimostrazione. Si ha:

$$e^x = e^0 + \frac{e^0}{1!}(x-0) + \frac{e^0}{2!}(x-0)^2 + \dots + \frac{e^0}{m!}(x-0)^m + \dots = \sum_{k=0}^{\infty} \frac{x^k}{k!}$$

□

5.10.2 Functions

Remark 187 (Supporto e spazio parametrico).

$$R_X = \{0, 1, 2, \dots\}$$

$$\Theta = \{\lambda \in \mathbb{R} : \lambda > 0\}$$

Definition 5.10.1 (Funzione di massa di probabilità).

$$p_X(x) = \mathbb{P}(X = x) = \frac{e^{(-\lambda)} \cdot \lambda^x}{x!} \cdot \mathbb{1}_{R_X}(x) \quad (5.38)$$

Validità PMF. Si ha:

$$\sum_{x=0}^{\infty} p_X(x) = \sum_{x=0}^{\infty} \frac{e^{-\lambda} \lambda^x}{x!} = e^{-\lambda} \sum_{x=0}^{\infty} \frac{\lambda^x}{x!} \stackrel{(1)}{=} e^{-\lambda} e^{\lambda} = 1$$

dove in (1) abbiamo sfruttato la 5.37 con le dovute sostituzioni di lettere. □

5.10.3 Moments

Proposition 5.10.2 (Momenti caratteristici).

$$\mathbb{E}[X] = \lambda \quad (5.39)$$

$$\text{Var}[X] = \lambda \quad (5.40)$$

$$\text{Asym}(X) = \frac{1}{\sqrt{\lambda}} \quad (5.41)$$

$$\text{Kurt}(X) = 3 + \frac{1}{\lambda} \quad (5.42)$$

Dimostrazione. Per il valore atteso

$$\begin{aligned} \mathbb{E}[X] &= \sum_{x=0}^{\infty} x \cdot \frac{e^{-\lambda} \lambda^x}{x!} \stackrel{(1)}{=} e^{-\lambda} \sum_{x=1}^{\infty} x \frac{\lambda^x}{x!} = \lambda e^{-\lambda} \sum_{x=1}^{\infty} \frac{\lambda^{x-1}}{(x-1)!} \\ &\stackrel{(2)}{=} \lambda e^{-\lambda} \sum_{y=0}^{\infty} \frac{\lambda^y}{y!} = \lambda e^{-\lambda} e^{\lambda} = \lambda \end{aligned}$$

dove in (1) abbiamo anche portato avanti di 1 la sommatoria dato che il primo termine è nullo e in (2) abbiamo sostituito $y = x - 1$ e sfruttato 5.37. □

Dimostrazione. Per la varianza troviamo innanzitutto $\mathbb{E}[X^2]$:

$$\mathbb{E}[X^2] = \sum_{x=0}^{\infty} x^2 \cdot \mathbb{P}(X=x) = \sum_{x=0}^{\infty} x^2 \frac{e^{-\lambda} \lambda^x}{x!} = e^{-\lambda} \sum_{x=0}^{\infty} x^2 \frac{\lambda^x}{x!}$$

Ora prendiamo la serie dell'esponenziale e la deriviamo rispetto a λ ad entrambi i membri (x costante)

$$e^\lambda = \sum_{x=0}^{\infty} \frac{\lambda^x}{x!} \stackrel{(1)}{=} \sum_{x=0}^{\infty} x \frac{\lambda^{x-1}}{x!} \stackrel{(2)}{=} \sum_{x=1}^{\infty} x \frac{\lambda^{x-1}}{x!}$$

dove in (1) abbiamo effettuato la derivazione (il primo membro rimane invariato), in (2) abbiamo portato avanti l'indice di sommatoria perché il primo termine è nullo. Ora moltiplicando per λ entrambi i lati si ottiene

$$\lambda e^\lambda = \sum_{x=1}^{\infty} x \frac{\lambda^x}{x!}$$

Effettuando gli stessi passaggi, nell'ordine derivare entrambi i membri rispetto a λ e moltiplicandoli per λ si prosegue come

$$\begin{aligned} \sum_{x=1}^{\infty} x^2 \frac{\lambda^{x-1}}{x!} &= e^\lambda + \lambda e^\lambda = e^\lambda(1 + \lambda) \\ \sum_{x=1}^{\infty} x^2 \frac{\lambda^x}{x!} &= e^\lambda \lambda(1 + \lambda) \end{aligned}$$

E infine riprendendo da dove eravamo arrivati con la main quest

$$\mathbb{E}[X^2] = e^{-\lambda} \sum_{x=0}^{\infty} x^2 \frac{\lambda^x}{x!} = e^{-\lambda} e^\lambda \lambda(1 + \lambda) = \lambda(1 + \lambda)$$

per cui

$$\text{Var}[X] = \mathbb{E}[X^2] - (\mathbb{E}[X])^2 = \lambda(1 + \lambda) - \lambda^2 = \lambda$$

□

Dimostrazione. Dimostrazione alternativa per la varianza:

$$\begin{aligned}
 \text{Var}[X] &= \mathbb{E}[X^2] - [\mathbb{E}[X]]^2 \\
 &= \left(\sum_{x=0}^{\infty} x^2 \cdot \frac{e^{-\lambda} \lambda^x}{x!} \right) - \lambda^2 \\
 &= \left(\sum_{x=0}^{\infty} (x^2 + x - x) \cdot \frac{e^{-\lambda} \lambda^x}{x!} \right) - \lambda^2 \\
 &= \left(\sum_{x=0}^{\infty} (x(x-1) + x) \cdot \frac{e^{-\lambda} \lambda^x}{x!} \right) - \lambda^2 \\
 &= \left(\sum_{x=0}^{\infty} (x(x-1)) \frac{e^{-\lambda} \lambda^x}{x!} + \sum_{x=0}^{\infty} x \frac{e^{-\lambda} \lambda^x}{x!} \right) - \lambda^2 \\
 &= \left(\sum_{x=0}^{\infty} x(x-1) \frac{e^{-\lambda} \lambda^2 \lambda^{x-2}}{x(x-1)(x-2)!} + \sum_{x=0}^{\infty} x \frac{e^{-\lambda} \lambda \lambda^{x-1}}{x(x-1)!} \right) - \lambda^2 \\
 &= \left(\sum_{x=0}^{\infty} \frac{\lambda^{x-2}}{(x-2)!} e^{-\lambda} \lambda^2 + \sum_{x=0}^{\infty} \frac{\lambda^{x-1}}{(x-1)!} e^{-\lambda} \lambda \right) - \lambda^2 \\
 &\stackrel{(1)}{=} \left(\sum_{z=0}^{\infty} \frac{\lambda^z}{z!} e^{-\lambda} \lambda^2 + \sum_{y=0}^{\infty} \frac{\lambda^y}{y!} e^{-\lambda} \lambda \right) - \lambda^2 \\
 &= (e^{\lambda} e^{-\lambda} \lambda^2 + e^{\lambda} e^{-\lambda} \lambda) - \lambda^2 \\
 &= (\lambda^2 + \lambda) - \lambda^2 \\
 &= \lambda
 \end{aligned}$$

dove in (1) abbiamo posto $y = x-1$, $z = x-2$ per sfruttare 5.37 nel seguito. \square

5.10.4 Shape

Remark 188 (Shape). Per la Poisson (figura 5.6):

- valore medio e varianza della vc di Poisson coincidono con il parametro λ ; la distribuzione ha picco intorno a λ
- per valori bassi di λ la distribuzione è asimmetrica positiva: se $\lambda < 1$ la distribuzione ha un andamento decrescente, mentre se $\lambda > 1$ è prima crescente e poi decrescente;
- al crescere λ , la distribuzione diventa più simmetrica e la curtosi tende a 3 (convergenndo ad una Normale).

5.10.5 Extras

5.10.5.1 Origine e approssimazione

Remark 189. La Poisson si può utilizzare per contare arrivi aleatori in un fissato intervallo (es di tempo), se sono soddisfatte alcune ipotesi. Se:

NB: Rigo approach

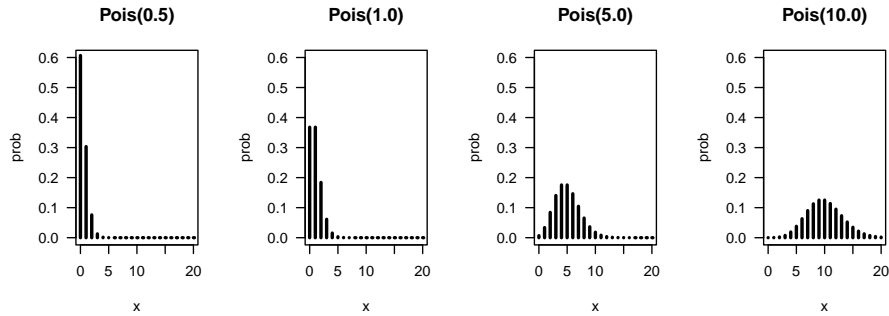


Figura 5.6: Distribuzione Poisson

1. ad ogni intervallo $I \subset [0, 1]$ possiamo associare una vc

$$N(I) = \text{numero arrivi nell'intervallo } I$$

(l'osservazione totale è $[0, 1]$ suddivisibile in intervalli)

2. possiamo ipotizzare che gli arrivi corrispondenti ad intervalli disgiunti diano luogo a vc indipendenti ($N(I_1), \dots, N(I_k)$ sono v.c. indipendenti se I_1, \dots, I_k sono intervalli a due a due disgiunti);
3. si ha che la probabilità di avere un evento in un dato intervallo

$$\mathbb{P}(N(I) = 1) = \lambda m(I) + f[m(i)]$$

dipende da

- un parametro λ
- $m(I)$ è la lunghezza dell'intervallo I
- f , una funzione tale che $\lim_{x \rightarrow 0} \frac{f(x)}{x} = 0$

4. la probabilità di avere più di un evento

$$\mathbb{P}(N(I) > 1) = g[m(i)]$$

dipende da una funzione g caratterizzata da $\lim_{x \rightarrow 0} \frac{g(x)}{x} = 0$ (ossia direi è infinitamente piccola se l'ampiezza dell'intervallo diminuisce)

Se valgono le condizioni di sopra allora

$$\mathbb{P}(N([0, 1]) = x) = \frac{e^{-\lambda} \lambda^x}{x!}, \quad \forall x = 0, 1, 2, \dots$$

dove $\lambda > 0$ è un parametro il cui significato è

$$\lambda = \lim_{n \rightarrow +\infty} \frac{N([0, n])}{n}$$

ossia coincide col numero medio di eventi osservati in un arco di tempo $[0, n]$ (per $n \rightarrow \infty$)

Remark 190. Non dimostriamo il risultato precedente

Remark 191. È utilizzata per modellare il numero di eventi registrati in un ambito circoscritto (*temporale* o *spaziale*), in cui vi è un largo numero di prove indipendenti (o quasi) caratterizzate ciascuna da una bassa probabilità di successo (per questa è chiamata legge degli eventi rari)

Proposition 5.10.3 (Paradigma di Poisson). *Siano E_1, \dots, E_n eventi con $p_i = \mathbb{P}(E_i)$, dove n è largo, p_i sono piccoli e gli E_i sono vc indipendenti o debolmente dipendenti. Sia*

$$X = \sum_{i=1}^n I_{E_i}$$

la somma di quanti eventi E_i siano accaduti. Allora X è abbastanza bene distribuita come una $\text{Pois}(\lambda)$ con $\lambda = \sum_i p_i$.

Dimostrazione. La prova dell'approssimazione di sopra è complessa, richiede definire la dipendenza debole e buona approssimazione; è omessa qui. \square

Remark 192 (Ruolo di λ). Il parametro λ è interpretato come *rate di occorrenza*: ad esempio $\lambda = 2$ mail di spam per giorno.

Remark 193. Nell'esempio sopra il numero di eventi X non è esattamente distribuito come Poisson perché una variabile di Poisson non ha limite superiore, mentre $I_{E_1} + \dots + I_{E_n}$ somma al più a n . Ma la distribuzione di Poisson da spesso una buona approssimazione e le condizioni per il verificarsi della situazione di sopra sono abbastanza flessibili: infatti i p_i non devono essere uguali e le prove non devono essere strettamente indipendenti. Questo fa sì che il modello di Poisson sia spesso un buon punto di partenza per dati che assumono valore intero non negativo (chiamati conteggi)

È comunque possibile quantificare l'errore commesso.

Proposition 5.10.4 (Errore di approssimazione). *Se E_i sono indipendenti e sia $N \sim \text{Pois}(\lambda)$, allora l'errore di approssimazione che si fa nell'utilizzare la poisson per stimare la probabilità di un dato set di interi non negativi $I \subset \mathbb{N}$, è dato dalla seguente:*

$$\mathbb{P}(X \in I) - \mathbb{P}(N \in I) \leq \min\left(1, \frac{1}{\lambda}\right) \sum_{i=1}^n p_i^2 \quad (5.43)$$

Dimostrazione. Anche questa è per ora complessa (necessita di una tecnica chiamata metodo di Stein). \square

Remark 194. La 5.43 fornisce un limite superiore dell'errore commesso nell'utilizzare una approssimazione di Poisson: non solo per l'intera distribuzione (se $I = \mathbb{N}$) ma per qualsiasi suo sottoinsieme. Altresì precisa quanto i p_i dovrebbero essere piccoli: vogliamo che $\sum_{i=1}^n p_i^2$ sia molto piccolo, o quanto meno lo sia rispetto a λ .

5.10.5.2 Legami con la binomiale

Remark 195. La relazione tra Poisson e Binomiale è simile a quella intercorrente tra Binomiale e Ipergeometrica: possiamo andare dalla Poisson alla binomiale condizionando, e viceversa dalla Binomiale alla Poisson prendendo un limite. Prima un risultato strumentale.

Dalla binomiale alla Poisson (Rigo)

Example 5.10.1 (Dalla binomiale alla poisson). Se $X_n \sim \text{Bin}(n, p_n)$, $X \sim \text{Pois}(\lambda)$ ed $np_n \rightarrow \lambda$ allora

$$\mathbb{P}(X = x) = \lim_{n \rightarrow +\infty} \mathbb{P}(X_n = x), \forall x = 0, 1, 2, \dots$$

Per fare i conti, occorre richiamare un risultato di analisi

$$\begin{aligned} \lim_{n \rightarrow +\infty} \left(1 + \frac{1}{n}\right)^n &= e, \\ \lim_{n \rightarrow +\infty} \left(1 + \frac{a}{n}\right)^n &= e^a \end{aligned}$$

Più in generale: se $\lim_{n \rightarrow +\infty} a_n = a$ allora $\lim_{n \rightarrow +\infty} \left(1 + \frac{a_n}{n}\right)^n = e^a$. Allora:

$$\begin{aligned} \mathbb{P}(X_n = x) &= \binom{n}{x} p_n^x (1 - p_n)^{n-x} \\ &= \frac{1}{x!} \underbrace{n(n-1) \cdots (n-x+1)}_{=1} \underbrace{(np_n)^x}_{=\lambda^x} \underbrace{(1-p_n)^{-x}}_{(1-0)^{-j}=1} \underbrace{\left(1 - \frac{np_n}{n}\right)^n}_{e^{-\lambda}} \end{aligned}$$

Quindi

$$\lim_{n \rightarrow \infty} P(X_n = x) = \frac{1}{x!} \cdot 1 \cdot \lambda^x \cdot 1 \cdot e^{-\lambda} = \frac{e^{-\lambda} \lambda^x}{x!} = \mathbb{P}(X = x)$$

Dalla binomiale alla Poisson (Blitzstein)

Remark 196. Se prendiamo il limite della $\text{Bin}(n, p)$ per $n \rightarrow \infty$ e $p \rightarrow 0$ con np fisso arriviamo alla Poisson.

Proposition 5.10.5 (Approssimazione Poissoniana della binomiale). *Se $X \sim \text{Bin}(n, p)$ e facciamo tendere $n \rightarrow \infty$, $p \rightarrow 0$ ma $\lambda = np$ rimane fisso, allora la PMF di X converge a $\text{Pois}(\lambda)$.*

La stessa conclusione si ha se $n \rightarrow \infty$, $p \rightarrow 0$ ed np converge ad una costante λ .

Remark 197. Questo è un *caso speciale* del paradigma di Poisson dove E_i sono indipendenti e hanno la stessa probabilità, quindi $\sum_{i=1}^n I_{E_i}$ ha distribuzione binomiale. In questo caso speciale possiamo dimostrare che l'approssimazione di Poisson ha senso limitandoci a prendere il limite della Binomiale.

Dimostrazione. Effettueremo la dimostrazione per $\lambda = np$ fisso (considerando $p = \lambda/n$), mostrando che la PMF $\text{Bin}(n, p)$ converge alla $\text{Pois}(\lambda)$. Per $0 \leq x \leq n$:

$$\begin{aligned} \mathbb{P}(X = x) &= \binom{n}{x} p^x (1 - p)^{n-x} \\ &= \frac{n(n-1) \cdots (n-x+1)}{x!} \left(\frac{\lambda}{n}\right)^x \left(1 - \frac{\lambda}{n}\right)^n \left(1 - \frac{\lambda}{n}\right)^{-x} \\ &= \frac{\lambda^x}{x!} \frac{n(n-1) \cdots (n-x+1)}{n^x} \left(1 - \frac{\lambda}{n}\right)^n \left(1 - \frac{\lambda}{n}\right)^{-x} \end{aligned}$$

Per $n \rightarrow \infty$ con k fisso

$$\begin{aligned} \frac{\overbrace{n(n-1) \cdot \dots \cdot (n-x+1)}^{x \text{ termini}}}{n^x} &\stackrel{(1)}{=} \frac{n \cdot n(1 - \frac{1}{n}) \cdot \dots \cdot n(1 - \frac{k-1}{n})}{n^x} \rightarrow 1 \\ \left(1 - \frac{\lambda}{n}\right)^n &\rightarrow e^{-\lambda} \\ \left(1 - \frac{\lambda}{n}\right)^{-k} &= \left[\left(1 - \frac{\lambda}{n}\right)^n\right]^{-\frac{k}{n}} \rightarrow e^{-\frac{k}{n}} = 1 \end{aligned}$$

dove in (1) abbiamo raccolto un n a partire dal secondo fattore, lasciando fuori parentesi k n che si moltiplicano. Pertanto

$$\mathbb{P}(X = x) \rightarrow \frac{e^{-\lambda} \lambda^x}{x!} = \text{Pois}(\lambda)$$

□

Remark 198. Il precedente risultato implica che se n è grande, p piccolo e np moderato, possiamo approssimare $\text{Bin}(n, p)$ con $\text{Pois}(np)$; come visto in precedenza l'errore nell'approssimare $\mathbb{P}(X \in I)$ con $\mathbb{P}(N \in I)$ per $X \sim \text{Bin}(n, p)$ e $N \sim \text{Pois}(np)$ è al massimo $\min(p, np^2)$.

Example 5.10.2. Il proprietario di un sito vuole studiare la distribuzione del numero di visitatori. Ogni giorno un milione di persone in maniera indipendente decide se visitare il sito o meno, con probabilità $p = 2 \times 10^{-1}$. Fornire una approssimazione della probabilità di avere almeno tre visitatori al giorno. Se $X \sim \text{Bin}(n, p)$ è il numero di visitatori con $n = 10^6$, fare i calcoli con la binomiale va incontro a difficoltà computazionali ed errori numerici del pc (dato che n è largo e p molto basso). Ma data la situazione con n largo p basso e $np = 2$ moderato, $\text{Pois}(2)$ è una buona approssimazione. Questo porta a

$$\mathbb{P}(X \geq 3) = 1 - \mathbb{P}(X < 3) \approx 1 - e^{-2} - e^{-2} \cdot 2 - e^{-2} \cdot \frac{2^2}{2!} = 1 - 5e^{-2} \approx 0.3233$$

che è una approssimazione molto accurata.

Dalla Poisson alla binomiale

Proposition 5.10.6. Se $X \sim \text{Pois}(\lambda_1)$ e $Y \sim \text{Pois}(\lambda_2)$ sono indipendenti, allora la distribuzione condizionata di X dato che $XY = n$ è $\text{Bin}(n, \lambda_1/(\lambda_1 + \lambda_2))$.

Dimostrazione. Utilizziamo la regola di Bayes per calcolare la PMF condizionata $\mathbb{P}(X = x | X + Y = n)$:

$$\begin{aligned} \mathbb{P}(X = x | X + Y = n) &= \frac{\mathbb{P}(X + Y = n | X = x) \cdot \mathbb{P}(X = x)}{\mathbb{P}(X + Y = n)} \\ &= \frac{\mathbb{P}(Y = n - x | X = x) \cdot \mathbb{P}(X = x)}{\mathbb{P}(X + Y = n)} \\ &\stackrel{(1)}{=} \frac{\mathbb{P}(Y = n - x) \cdot \mathbb{P}(X = x)}{\mathbb{P}(X + Y = n)} \end{aligned}$$

con (1) per indipendenza delle due. Ora sostituendo le PMF di X, Y e $X + Y$; questa al denominatore è distribuita come $\text{Pois}(\lambda_1 + \lambda_2)$ per proposizione 5.10.7. Si ha:

$$\begin{aligned}
 \mathbb{P}(X = k | X + Y = n) &= \frac{\left(\frac{e^{-\lambda_2} \lambda_2^{n-k}}{(n-k)!}\right) \left(\frac{e^{-\lambda_1} \lambda_1^k}{k!}\right)}{\frac{e^{-(\lambda_1 + \lambda_2)} (\lambda_1 + \lambda_2)^n}{n!}} = \frac{\frac{e^{-(\lambda_1 + \lambda_2)} \cdot \lambda_1^k \cdot \lambda_2^{n-k}}{k!(n-k)!}}{\frac{e^{-(\lambda_1 + \lambda_2)} \cdot (\lambda_1 + \lambda_2)^n}{n!}} \\
 &= \frac{e^{-(\lambda_1 + \lambda_2)} \cdot \lambda_1^k \cdot \lambda_2^{n-k}}{k!(n-k)!} \cdot \frac{n!}{e^{-(\lambda_1 + \lambda_2)} \cdot (\lambda_1 + \lambda_2)^n} \\
 &= \frac{n!}{k!(n-k)!} \cdot \frac{\lambda_1^k \cdot \lambda_2^{n-k}}{(\lambda_1 + \lambda_2)^n} \\
 &= \binom{n}{k} \left(\frac{\lambda_1^k}{\lambda_1 + \lambda_2}\right)^k \left(\frac{\lambda_2^{n-k}}{\lambda_1 + \lambda_2}\right)^{n-k} \\
 &= \text{Bin}\left(n, \frac{\lambda_1}{\lambda_1 + \lambda_2}\right)
 \end{aligned}$$

□

5.10.5.3 Somma di Poisson indipendenti

Proposition 5.10.7 (Somma di Poisson indipendenti). *Siano $X \sim \text{Pois}(\lambda_1)$ e $Y \sim \text{Pois}(\lambda_2)$ vc indipendenti. Allora $X + Y \sim \text{Pois}(\lambda_1 + \lambda_2)$*

Dimostrazione. Per ottenere la PMF di $X + Y$ condizioniamo su X e utilizziamo il teorema delle probabilità totali

$$\begin{aligned}
 \mathbb{P}(X + Y = k) &= \sum_{j=0}^k \mathbb{P}(X + Y = k | X = j) \cdot \mathbb{P}(X = j) \\
 &= \sum_{j=0}^k \mathbb{P}(Y = k - j | X = j) \cdot \mathbb{P}(X = j) \\
 &\stackrel{(1)}{=} \sum_{j=0}^k \mathbb{P}(Y = k - j) \cdot \mathbb{P}(X = j) \\
 &= \sum_{j=0}^k \frac{e^{-\lambda_2} \lambda_2^{k-j}}{(k-j)!} \frac{e^{-\lambda_1} \lambda_1^j}{(j)!} \\
 &= \frac{e^{-(\lambda_1 + \lambda_2)}}{k!} \sum_{j=0}^k \binom{k}{j} \lambda_1^j \lambda_2^{k-j} \\
 &\stackrel{(2)}{=} \frac{e^{-(\lambda_1 + \lambda_2)} (\lambda_1 + \lambda_2)^k}{k!} = \text{Pois}(\lambda_1 + \lambda_2)
 \end{aligned}$$

con (1) data l'indipendenza e in (2) si è utilizzato il teorema binomiale $(a+b)^n = \sum_{i=0}^n \binom{n}{i} a^i b^{n-i}$ □

Remark 199. A intuito se vi sono due tipi di eventi che accadono ai rate λ_1 e λ_2 indipendentemente, allora il rate complessivo di eventi è $\lambda_1 + \lambda_2$.

5.11 Discrete uniform

5.11.1 Definition

Remark 200. La prova che genera la vc Uniforme discreta si può assimilare all'estrazione di una pallina da un'urna che contiene n palline identiche numerate da 1 a n . Viene in genere utilizzata quanto tutti i risultati dell'esperimento sono equiprobabili

Definition 5.11.1 (Uniforme discreta). Il numero X della pallina estratta dall'urna contenente n palline numerate (da 1 a n) si distribuisce come Uniforme discreta $X \sim \text{DUnif}(n)$.

5.11.2 Functions

Remark 201 (Supporto e spazio parametrico).

$$\begin{aligned} R_X &= \{1, \dots, n\} \\ \Theta &= \{n \in \mathbb{N} \setminus \{0\}\} \end{aligned}$$

Proposition 5.11.1 (Funzione di massa di probabilità).

$$p_X(x) = \mathbb{P}(X = x) = \frac{1}{n} \cdot \mathbb{1}_{R_X}(x) \quad (5.44)$$

Definition 5.11.2 (Funzione di ripartizione).

$$F_X(x) = \mathbb{P}(X \leq x) = \begin{cases} 0 & \text{se } x < 1 \\ \frac{k}{n} & \text{se } k \leq x < k+1, (k = 1, 2, \dots, n-1) \\ 1 & \text{se } x \geq n \end{cases} \quad (5.45)$$

Remark 202. La funzione di ripartizione è nulla in $(-\infty; 1)$ ed è una funzione a gradini di altezza costante pari a $1/n$, in corrispondenza di ogni valore intero $1 \leq x \leq n$ e vale 1 in $[n; +\infty)$.

5.11.3 Moments

Proposition 5.11.2 (Momenti caratteristici).

$$\mathbb{E}[X] = \frac{n+1}{2} \quad (5.46)$$

$$\text{Var}[X] = \frac{n^2 - 1}{12} \quad (5.47)$$

$$\text{Asym}(X) = 0 \quad (5.48)$$

$$\text{Kurt}(X) = 1.8 \quad (5.49)$$

Dimostrazione.

$$\mathbb{E}[X] = \sum_{x=1}^n x \frac{1}{n} = \frac{1}{n} (1 + 2 + \dots + n) = \frac{1}{n} \frac{n(n+1)}{2} = \frac{n+1}{2}$$

□

Dimostrazione.

$$\begin{aligned}
 \text{Var}[X] &= \mathbb{E}[X^2] - [\mathbb{E}[x]]^2 = \left(\sum_{x=1}^n x^2 \frac{1}{n} \right) - \left(\frac{n+1}{2} \right)^2 \\
 &= \left(\frac{1}{n} (1^2 + 2^2 + \dots + n^2) \right) - \left(\frac{n+1}{2} \right)^2 \\
 &= \left(\frac{1}{n} \cdot \frac{n(n+1)(2n+1)}{6} \right) - \left(\frac{n^2 + 1 + 2n}{4} \right) \\
 &= \left(\frac{(n+1)(2n+1)}{6} \right) - \left(\frac{n^2 + 1 + 2n}{4} \right) \\
 &= \frac{2(2n^2 + 2n + n + 1) - 3(n^2 + 1 + 2n)}{12} \\
 &= \frac{4n^2 + 4n + 2n + 2 - 3n^2 - 3 - 6n}{12} \\
 &= \frac{n^2 - 1}{12}
 \end{aligned}$$

□

5.12 Exercises

Example 5.12.1 (Uniforme discreta: dado regolare “generalizzato” (esercizio pag 45 McColl)). Sia X una vc con supporto $R_X = \{k+1, k+2, \dots, k+n\}$ con k intero ed n intero positivo.

Si tratta di una VC discreta; per descriverla dobbiamo usare un PMF discreta che assegni a ciascun valore del supporto una probabilit 

$$p_X(x) = \frac{1}{n} \mathbb{1}_{R_X}(x)$$

Quindi il valore di p_X   sempre comune e pari a $1/n$ per tutti i valori del supporto (PMF a sticks verticali della stessa altezza).

Questa   l'**uniforme discreta** definita tra $k+1$ a $k+n$:   utile in alcuni casi concreti. Ad esempio X pu  essere il punteggio nel lancio di un dado regolare se si imposta $k=0$ ed $n=6$.

Calcoliamo il valore atteso nel caso generico

$$\mathbb{E}[X] = \sum_{x=k+1}^{k+n} x \cdot p_X(x) = \sum_{x=k+1}^{k+n} x \cdot \frac{1}{n} = \frac{1}{n} \sum_{x=k+1}^{k+n} x$$

Trasliamo l'indice per farlo variare tra 1 e n , sostituendo $y = x - k$ (avendo dunque $x = y + k$). Riscriviamo la sommatoria in funzione di y

$$\begin{aligned}
 \mathbb{E}[X] &= \frac{1}{n} \sum_{y=1}^n y + k = \frac{1}{n} \left[\sum_{y=1}^n y + \sum_{y=1}^n k \right] = \frac{\sum_{y=1}^n y}{n} + k \\
 &\stackrel{(1)}{=} \frac{n(n+1)}{2n} + k = \frac{n+1}{2} + k
 \end{aligned}$$

Dove in (1) abbiamo sostituito la somma dei primi n interi positivi. Questo ultimo è il valore atteso. Tornando al grafico il valore atteso dipende da n e da k ma si piazza al centro degli sticks verticali di $p_X(x)$ (essendo la distribuzione simmetrica).

Ad esempio se nel caso del dado standard con $k = 0$ ed $n = 6$, si ha che $\mathbb{E}[X] = 0 + 7/2 = 3.5$. Questo è il punteggio atteso risultante da il lancio di un dado regolare; la particolarità di questo risultato è che il valore non è nel supporto. Questo avviene spesso se la variabile casuale è discreta (il valore atteso non coincide con alcun risultato possibile della variabile casuale).

Example 5.12.2 (Poisson (crash course, day 1 es 3 pag 6)). Let $f(k) = \frac{c^k e^{-\lambda}}{k!}$ for $k \in \{0, 1, \dots\}$ be the pmf that X satisfies:

1. find c
2. find $\mathbb{E}[X]$
3. find $\text{Var}[X]$

we have

1.

$$\sum_{k=0}^{\infty} f(k) = 1 = e^{-\lambda} \underbrace{\sum_{k=0}^{\infty} \frac{c^k}{k!}}_{e^c} = e^{-\lambda} e^c = e^{c-\lambda} = 1 = e^0 \implies c = \lambda$$

2.

$$\begin{aligned} \mathbb{E}[X] &= \sum_{k=0}^{\infty} k f(k) = \sum_{k=0}^{\infty} k \frac{\lambda^k e^{-\lambda}}{k!} = \sum_{k=1}^{\infty} k \frac{\lambda^k e^{-\lambda}}{k!} = \sum_{k=1}^{\infty} \frac{\lambda^k e^{-\lambda}}{(k-1)!} = \lambda \sum_{k=1}^{\infty} \frac{\lambda^{k-1} e^{-\lambda}}{(k-1)!} \\ &\stackrel{(1)}{=} \lambda \underbrace{\sum_{u=0}^{\infty} \frac{\lambda^u e^{-\lambda}}{u!}}_{F(\infty)=1} = \lambda \end{aligned}$$

with (1) substituting $u = k - 1$. This is the poisson distribution, we say $X \sim \text{Pois}(\lambda)$

3. first we find $\mathbb{E}[X^2]$, but first consider the following

$$\begin{aligned} \mathbb{E}[X(X-1)] &= \mathbb{E}[X^2] - \mathbb{E}[X] = \sum_{k=0}^{\infty} k(k-1) f(k) = \sum_{k=2}^{\infty} k(k-1) f(k) \\ &= \sum_{k=2}^{\infty} k(k-1) \frac{\lambda^k e^{-\lambda}}{k!} = \sum_{k=2}^{\infty} \frac{\lambda^k e^{-\lambda}}{(k-2)!} = \lambda^2 \sum_{k=2}^{\infty} \frac{\lambda^{k-2} e^{-\lambda}}{(k-2)!} \\ &\stackrel{(1)}{=} \lambda^2 \underbrace{\sum_{u=0}^{\infty} \frac{\lambda^u e^{-\lambda}}{u!}}_{F(\infty)=1} = \lambda^2 = \mathbb{E}[X^2] - \mathbb{E}[X] \end{aligned}$$

where in (1) doin subst $u = k - 2$. Therefore

$$\mathbb{E}[X^2] = \lambda^2 + \lambda \implies \text{Var}[X] = \mathbb{E}[X^2] - \mathbb{E}[X]^2 = \lambda^2 + \lambda - \lambda^2 = \lambda$$

Capitolo 6

Absolute continuous random variables

6.1 Uniforme continua

6.1.1 Definition

Remark 203. È una vc continua X definita sul supporto (a, b) (con estremi esclusi), $a < b$, ed esiti aventi la medesima densità, indicata con $X \sim \text{Unif}(a, b)$.

Si usa nella pratica se il fenomeno si verifica entro un certo range, ed in questo intervallo tutte le realizzazioni sono per noi indifferenti (non vi sono aree con più densità di altre).

Un **esempio**: l'orario di atterraggio di un aereo previsto per le 14.00 è ragionevolmente uniforme tra le 13.50 e 14.10: $\text{Orario} \sim \text{Unif}(13.50, 14.10)$.

Una formulazione frequente si ha se $a = 0, b = 1$.

6.1.2 Functions

Remark 204 (Supporto e spazio parametrico).

$$\begin{aligned} R_X &= (a, b) \\ \Theta &= \{a, b \in \mathbb{R} : a < b\} \end{aligned}$$

Definition 6.1.1 (Funzione di densità). In figura 6.1

$$f_X(x) = \frac{1}{b-a} \cdot \mathbb{1}_{R_X}(x) \quad (6.1)$$

Proposition 6.1.1. *L'area è 1.*

Dimostrazione.

$$(b-a) \cdot \frac{1}{(b-a)} = 1$$

□

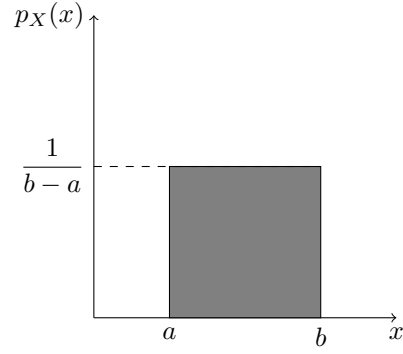


Figura 6.1: Uniforme continua

Definition 6.1.2 (Funzione di ripartizione).

$$F_X(x) = \begin{cases} 0 & \text{per } x \leq a \\ \frac{x-a}{b-a} & \text{se } a < x < b \\ 1 & \text{per } x \geq b \end{cases} \quad (6.2)$$

6.1.3 Moments

Proposition 6.1.2 (Momenti caratteristici).

$$\mathbb{E}[X] = \frac{a+b}{2} \quad (6.3)$$

$$\text{Var}[X] = \frac{(b-a)^2}{12} \quad (6.4)$$

$$\text{Asym}(X) = 0 \quad (6.5)$$

$$\text{Kurt}(X) = 1.8 \quad (6.6)$$

Dimostrazione.

$$\begin{aligned} \mathbb{E}[X] &= \int_a^b x \frac{1}{b-a} dx = \left[\frac{x^2}{2(b-a)} \right]_a^b \\ &= \left(\frac{b^2}{2(b-a)} + c \right) - \left(\frac{a^2}{2(b-a)} + c \right) \\ &= \frac{b^2 - a^2}{2(b-a)} = \frac{a+b}{2} \end{aligned}$$

□

Dimostrazione.

$$\begin{aligned}
 \text{Var}[X] &= \left(\int_a^b x^2 \frac{1}{b-a} dx \right) - \left(\frac{a+b}{2} \right)^2 \\
 &= \left[\frac{x^3}{3(b-a)} \right]_a^b - \left(\frac{a+b}{2} \right)^2 \\
 &= \left(\frac{b^3}{3(b-a)} + c \right) - \left(\frac{a^3}{3(b-a)} + c \right) - \left(\frac{a+b}{2} \right)^2 \\
 &= \frac{b^3 - a^3}{3(b-a)} - \frac{(a+b)^2}{4} \\
 &= \frac{(b-a)(a^2 + b^2 + ab)}{3(b-a)} - \frac{(a+b)^2}{4} \\
 &= \frac{a^2 + b^2 + ab}{3} - \frac{a^2 + b^2 + 2ab}{4} \\
 &= \frac{4a^2 + 4b^2 + 4ab - 3a^2 - 3b^2 - 6ab}{12} \\
 &= \frac{a^2 + b^2 - 2ab}{12} = \frac{(a-b)^2}{12} = \frac{(b-a)^2}{12}
 \end{aligned}$$

□

6.1.4 Shape

Remark 205. Si tratta di una variabile simmetrica e platicurtica (ovvero con una distribuzione molto piatta).

6.2 Esponenziale

6.2.1 Definition

Remark 206. L'esponenziale

- è generalmente usata per fenomeni di cui interessa un tempo t di attesa tra un istante iniziale e l'istante in cui si verifica un evento aleatorio di nostro interesse (tempo di vita, resistenza/funzionamento);
- può essere derivata se si ipotizza una funzione di rischio/azzardo costante $H(t) = \lambda > 0$, con λ tasso di occorrenza dell'evento (reciproco del numero di eventi per unità di tempo).

6.2.2 Functions

Remark 207 (Supporto e spazio parametrico).

$$\begin{aligned}
 R_X &= \{x \in \mathbb{R} : x > 0\} \\
 \Theta &= \{\lambda \in \mathbb{R} : \lambda > 0\}
 \end{aligned}$$

Definition 6.2.1 (Distribuzione esponenziale). Se $H(t) = \lambda > 0$ la funzione di ripartizione si ricava dalla 4.25 come

$$\begin{aligned} F_X(t) &= 1 - \exp\left(-\int_0^t H(w) \, dw\right) = 1 - \exp\left(-\int_0^t \lambda \, dw\right) \\ &= 1 - \exp(-\lambda t) \end{aligned}$$

Definition 6.2.2 (Funzione di ripartizione).

$$F_X(x) = \begin{cases} 1 - \exp(-\lambda x) & \text{per } x \geq 0 \\ 0 & \text{per } x < 0 \end{cases} \quad (6.7)$$

Remark 208. La funzione di densità si ottiene derivando dalla 6.7; pertanto una vc continua X si dice vc Esponenziale con parametro $\lambda > 0$, e si scrive $X \sim \text{Exp}(\lambda)$ se caratterizzata dalla seguente funzione di densità.

Definition 6.2.3 (Funzione di densità).

$$f_X(x) = \lambda \exp(-\lambda x) \cdot \mathbb{1}_{R_X}(x) \quad (6.8)$$

6.2.3 Moments

Proposition 6.2.1 (Momenti caratteristici).

$$\mathbb{E}[X] = \frac{1}{\lambda} \quad (6.9)$$

$$\text{Var}[X] = \frac{1}{\lambda^2} \quad (6.10)$$

$$\text{Asym}(X) = 2 \quad (6.11)$$

$$\text{Kurt}(X) = 9 \quad (6.12)$$

Dimostrazione. Per il valore atteso

$$\begin{aligned} \mathbb{E}[X] &= \int_{-\infty}^{+\infty} x f(x) \, dx = \int_0^{+\infty} x \lambda e^{-\lambda x} \, dx = \lambda \left\{ \left[\frac{e^{-\lambda x}}{-\lambda} x \right]_0^{+\infty} + \frac{1}{\lambda} \int_0^{+\infty} e^{-\lambda x} \, dx \right\} \\ &= \int_0^{+\infty} e^{-\lambda x} \, dx = \frac{1}{\lambda} \int_0^{+\infty} e^{-y} \, dy = \frac{1}{\lambda} \end{aligned}$$

□

6.2.4 Shape

Remark 209 (Forma distribuzione). La densità (figura 6.2) è:

- decrescente e asimmetrica positiva (tempi di attesa con densità più alte sono i più bassi)
- all'aumentare di λ la asimmetria positiva aumenta (con tempi più piccoli sempre più impattanti)

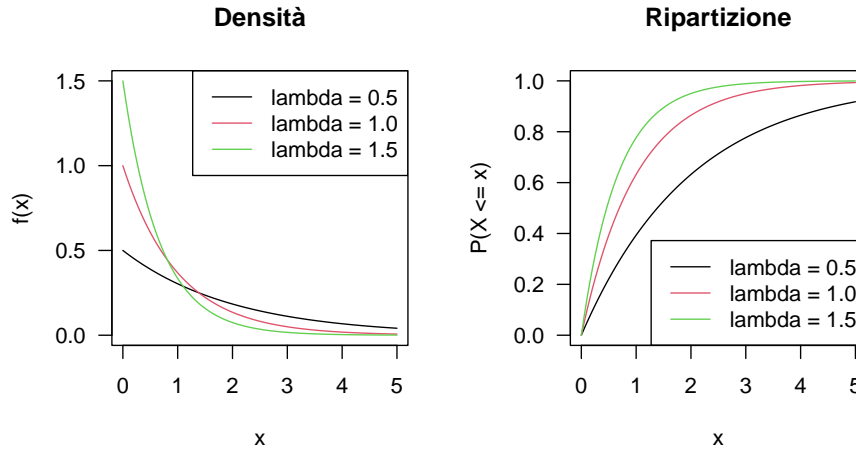


Figura 6.2: Distribuzione esponenziale

6.2.5 Extras

Proposition 6.2.2. *La vc esponenziale è l'unica rv assolutamente continua con mancanza di memoria*

$$\mathbb{P}(X > a + b | X > a) = \mathbb{P}(X > b), \quad \forall a, b > 0$$

Remark 210. Il tipico fenomeno che si descrive mediante la va esponenziale è la durata in vita (di qualcuno o qualcosa). In tale interpretazione

$$\begin{aligned} \mathbb{P}(X > a + b | X > a) &= \mathbb{P}(\text{sono vivo all'istante } a + b | \text{sono vivo all'istante } a) \\ &= \mathbb{P}(\text{sono vivo all'istante } b) \end{aligned}$$

Per inciso da questo segue che l'esponenziale NON è un modello adatto a descrivere la durata in vida di un essere vivente,

$$\mathbb{P}(X > 85 | X > 80) \neq \mathbb{P}(X > 5)$$

Dimostrazione. Dimostriamo che se X è esponenziale, allora vale la mancanza di memoria

$$\begin{aligned} \mathbb{P}(X > a + b | X > a) &= \frac{\mathbb{P}(X > a + b, X > a)}{\mathbb{P}(X > a)} = \frac{\mathbb{P}(X > a + b)}{\mathbb{P}(X > a)} = \frac{1 - F(a + b)}{1 - F(a)} \\ &= \frac{e^{-\lambda(a+b)}}{e^{-\lambda a}} = e^{-\lambda b} = \mathbb{P}(X > b) \end{aligned}$$

Dimostrare che se X è assolutamente continua e manca di memoria allora è esponenziale è molto più complicato \square

Remark 211. La vc Esponenziale presenta una struttura molto semplice ma rigida, per cui non si adatta facilmente a tutte le situazioni reali; infatti, talvolta non è realistico assumere che la funzione di rischio si costante rispetto al tempo. Pertanto si hanno almeno due generalizzazioni: la Weibull e la Gamma.

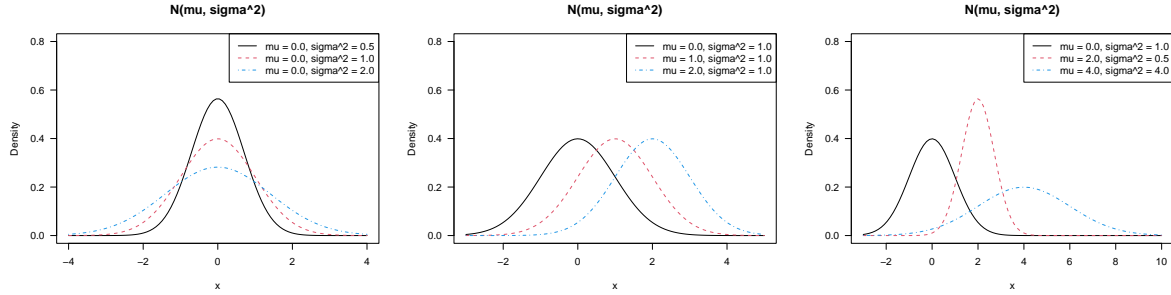


Figura 6.3: Distribuzione normale

6.3 Normale/Gaussiana

6.3.1 Definition

Remark 212.

Remark 213 (Situazioni di utilizzo). Viene utilizzata come prima approssimazione per descrivere variabili casuali a valori reali che tendono a concentrarsi attorno a un singolo valor medio. In particolare si usa:

- in inferenza
- quando si studiano gli errori di misura
- per la distribuzione di certe grandezze biometriche (es peso bambini alla nascita)
- se la grandezza aleatoria di interesse è la *somma* di un *numero* molto *elevato* di *fattori* aleatori tra loro *indipendenti*

6.3.2 Functions

Remark 214. Una vc continua si dice vc Normale con parametri μ e σ^2 , e la si indica con $X \sim N(\mu, \sigma^2)$ se è definita su tutto l'asse reale e presenta la seguente funzione di densità.

Remark 215 (Supporto e spazio parametrico).

$$R_X = \mathbb{R}$$

$$\Theta = \{\mu \in \mathbb{R}; \sigma^2 \in \mathbb{R} : \sigma^2 > 0\}$$

Definition 6.3.1 (Funzione di densità).

$$f_X(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \cdot \exp \left[-\frac{1}{2} \left(\frac{x - \mu}{\sigma} \right)^2 \right] \cdot \mathbb{1}_{R_X}(x) \quad (6.13)$$

In questa definizione potrei anche non mettere la funzione indicatrice del supporto, poiché essendo \mathbb{R} , assume sempre valore 1.

6.3.3 Shape

Remark 216 (Forma della distribuzione). Si ha che (figura 6.3)

- μ è parametro di locazione, σ^2 misura la dispersione attorno a μ ;
- la modifica di μ a parità di σ^2 implica una traslazione della funzione di densità lungo l'asse x ;
- al crescere di σ a parità di μ , i flessi si allontanano da μ e la funzione di densità attribuisce maggiore probabilità ai valori lontani dal valore centrale (e viceversa al diminuire di σ^2).

Inoltre:

- ha una forma campanulare ed simmetrica rispetto al parametro di locazione $x = \mu$

$$f_X(\mu - x) = f_X(\mu + x), \forall x$$

è crescente in $(-\infty, \mu)$ e decrescente in (μ, ∞) .

- in corrispondenza di μ , $f_X(x)$ ha il massimo (perché l'esponente negativo è minimo). Pertanto μ è il valore centrale la moda, mediana e valore medio della vc.
- si dimostra che $f_X(x)$ presenta due flessi in corrispondenza di $x = \mu \pm \sigma$. Ha come asintoto l'asse x

6.3.4 Normale standardizzata

Definition 6.3.2 (Normale standardizzata). Se $X \sim N(\mu, \sigma^2)$, la trasformazione lineare $Z = (X - \mu)/\sigma$ definisce la vc Normale standardizzata $Z \sim N(0, 1)$

Definition 6.3.3 (Funzione di densità (Normale standardizzata)).

$$\phi(z) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{z^2}{2}\right) \cdot \mathbb{1}_{R_X}(z) \quad (6.14)$$

Definition 6.3.4 (Funzione di ripartizione (Normale standardizzata)).

$$\Phi(z) = \int_{-\infty}^z \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{w^2}{2}\right) dw \quad (6.15)$$

Remark 217. La funzione di ripartizione della vc Z non ammette una formulazione esplicita ed è necessario predisporre delle tavole che per opportuni valori di z forniscano l'integrale con sufficiente accuratezza.

Remark 218. Sfruttando la simmetria della funzione di densità, è sufficiente conoscere $\Phi(z)$ per i soli valori di $z > 0$. Infatti $\Phi(0) = 0.5$ ed inoltre:

$$\Phi(-z) = 1 - \Phi(z) \quad \forall z \geq 0 \quad (6.16)$$

Remark 219. La conoscenza della funzione di ripartizione della vc $Z \sim N(0, 1)$ è sufficiente per calcolare la probabilità di qualsiasi vc $X \sim N(\mu, \sigma^2)$ mediante una semplice trasformazione:

$$\begin{aligned}\mathbb{P}(x_0 < X \leq x_1) &= \mathbb{P}\left(\frac{x_0 - \mu}{\sigma} < \underbrace{\frac{X - \mu}{\sigma}}_Z \leq \frac{x_1 - \mu}{\sigma}\right) \\ &= \Phi\left(\frac{x_1 - \mu}{\sigma}\right) - \Phi\left(\frac{x_0 - \mu}{\sigma}\right)\end{aligned}$$

In pratica per calcolare la probabilità che una vc normale assuma valori in un intervallo basta standardizzare gli estremi dell'intervallo ed utilizzare le tavole di $\Phi(z)$.

6.3.5 Moments

Proposition 6.3.1 (Momenti caratteristici (Normale)). *Da $X = \mu + \sigma Z$ si ha*

$$\mathbb{E}[X] = \mu \quad (6.17)$$

$$\text{Var}[X] = \sigma^2 \quad (6.18)$$

$$\text{Asym}(X) = 0 \quad (6.19)$$

$$\text{Kurt}(X) = 3 \quad (6.20)$$

6.3.6 Extras

Remark 220. La famiglia delle vc normali è chiusa rispetto ad ogni combinazione lineare: in particolare la combinazione lineare di vc normali e indipendenti è ancora una vc normale che ha per valore medio la combinazione lineare dei valori medi e per varianza la combinazione lineare delle varianze con i quadrati dei coefficienti (proprietà riproduttiva della vc normale).

Proposition 6.3.2. *Se $X_i \sim N(\mu_i, \sigma_i^2)$, allora:*

$$\sum_{i=1}^n a_i X_i \sim N\left(\sum_{i=1}^n a_i \mu_i, \sum_{i=1}^n a_i^2 \sigma_i^2\right)$$

6.4 Gamma

6.4.1 Definition

Remark 221. Viene utilizzata quando si deve verificare la lunghezza dell'intervallo di tempo fino all'istante in cui si verifica la n -esima manifestazione di un evento aleatorio di interesse.

Da notare che a livello matematico n può essere un numero reale positivo non solo naturale (Soffritti la chiama α)

Costituisce una generalizzazione dell'esponenziale, la quale si può ottenere se $n = 1$ Similmente alla Beta è chiamata così perché coinvolge l'omonima funzione matematica.

6.4.2 Functions

Remark 222 (Supporto e spazio parametrico).

$$\begin{aligned} R_X &= \{x \in \mathbb{R} : x > 0\} \\ \Theta &= \{n \in \mathbb{R} : n > 0; \lambda \in \mathbb{R} : \lambda > 0\} \end{aligned}$$

Definition 6.4.1 (Funzione di densità). Una vc continua X si distribuisce come una Gamma con parametri $n > 0, \lambda > 0$, indicata con $X \sim \text{Gamma}(n, \lambda)$, se presenta una funzione di densità come la:

$$f_X(x) = \frac{\lambda^n}{\Gamma(n)} \cdot x^{n-1} \exp(-\lambda x) \cdot \mathbb{1}_{R_X}(x) \quad (6.21)$$

Definition 6.4.2 (Funzione Gamma). È definita come

$$\Gamma(n) = \int_0^{+\infty} x^{n-1} e^{-x} dx \quad (6.22)$$

Important remark 42. Si ha che:

- il risultato della funzione 6.22 dipende solo da n e non da x perché sto integrando per quest'ultimo;
- questo integrale *non ammette soluzione* analitica e abbiamo bisogno di strumenti numerici per risolverlo. In R si usa l'omonima funzione **gamma** alla quale passiamo n
- tre proprietà della funzione gamma (appendice 2 McColl)

$$\Gamma(n) = (n-1) \cdot \Gamma(n-1), \quad n > 1 \text{ (è ricorsiva)} \quad (6.23)$$

$$\Gamma(n) = (n-1)! \quad \text{se } n \in \mathbb{N} \setminus \{0\} \text{ (link col fattoriale)} \quad (6.24)$$

$$\Gamma\left(\frac{1}{2}\right) = \sqrt{\pi} \quad \text{(un valore notevole)} \quad (6.25)$$

Remark 223 (Funzione di ripartizione). Non si può definire una funzione di ripartizione perché questa dipende dalla funzione Γ (a meno che n sia intero).

6.4.3 Moments

Proposition 6.4.1 (Momenti caratteristici).

$$\mathbb{E}[X] = \frac{n}{\lambda} \quad (6.26)$$

$$\text{Var}[X] = \frac{n}{\lambda^2} \quad (6.27)$$

$$\text{Asym}(X) = \frac{2}{\sqrt{n}} \quad (6.28)$$

$$\text{Kurt}(X) = 3 + \frac{6}{n} \quad (6.29)$$

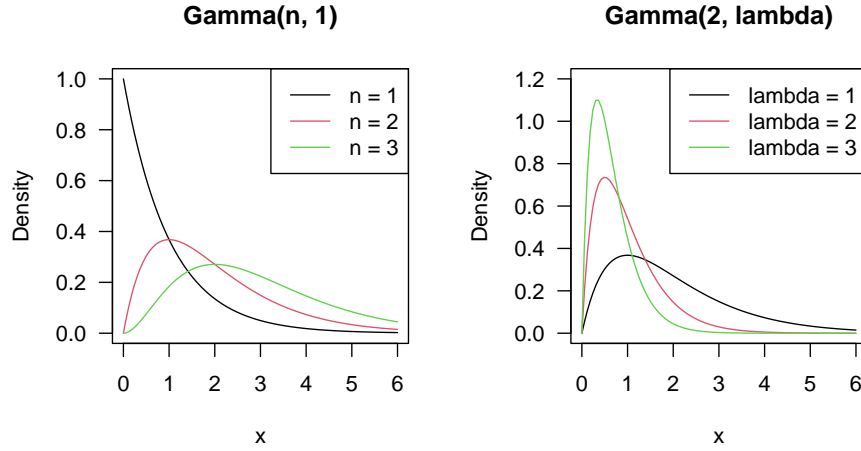


Figura 6.4: Distribuzione gamma

6.4.4 Shape

Remark 224 (Forma della distribuzione). λ è un parametro di scala mentre n determina la forma della distribuzione (figura 6.4):

- all'aumentare di n l'asimmetria positiva diminuisce perché se aumenta il numero di eventi di interesse i tempi di attesa si allungano; quando $n \rightarrow \infty$ la distribuzione diviene simmetrica e di forma campanulare (curtosi pari a 3);
- all'aumentare di λ la distribuzione si concentra sui valori/tempi di attesa più piccoli e l'asimmetria positiva aumenta (similmente a quanto visto per l'esponenziale)

6.4.5 Extras

Important remark 43 (Caso particolare: esponenziale). Si nota che se $n = 1$, la distribuzione gamma diviene una esponenziale in quanto la densità si semplifica a quella di una esponenziale

$$f_X(x) = \frac{\lambda^1}{1} x^0 e^{-\lambda x} \cdot \mathbb{1}_{R_X}(x) = \lambda e^{-\lambda x} \cdot \mathbb{1}_{R_X}(x)$$

che è la funzione di densità dell'esponenziale ovvero $\text{Gamma}(1, \lambda) \sim \text{Exp}(\lambda)$. Pertanto

- la gamma è una generalizzazione della esponenziale.
- il parametro λ ha lo stesso significato di quello dell'esponenziale (quindi ci aspettiamo di utilizzarlo per i tempi di attesa); in particolare usiamo la gamma per misurare l'intervallo di tempo tra un istante iniziale e l' n -esima manifestazione di un evento aleatorio di interesse (mentre per l'esponenziale è solo il primo)

Proposition 6.4.2. *La gamma gode della proprietà riproduttiva nel senso che la somma di gamma indipendenti ancora una gamma:*

$$\sum \text{Gamma}(n_i, \lambda) \sim \text{Gamma}\left(\sum_i n_i, \lambda\right) \quad (6.30)$$

6.5 Chi-quadrato

6.5.1 Definition

Remark 225. Si può ottenere:

1. dalla somma di ν vc normali standardizzate indipendenti ed elevate al quadrato (pag67 McColl); questa è una vc continua sul supporto $(0, +\infty)$ che si distribuisce come una Chi-quadrato con ν gradi di libertà

$$\sum_{i=1}^{\nu} Z_i^2 \sim \chi^2(\nu)$$

e come caso particolare

$$Z^2 \sim \chi^2(1)$$

2. come caso particolare della gamma: se $n = \frac{\nu}{2}$ (con $\nu \in \mathbb{N} \setminus \{0\}$, numero dei gradi di libertà) e $\lambda = \frac{1}{2}$

$$\text{Gamma}\left(\frac{\nu}{2}, \frac{1}{2}\right) \sim \chi^2(\nu)$$

Remark 226 (Utilizzi). Si usa in inferenza statistica.

6.5.2 Functions

Remark 227 (Supporto e spazio parametrico).

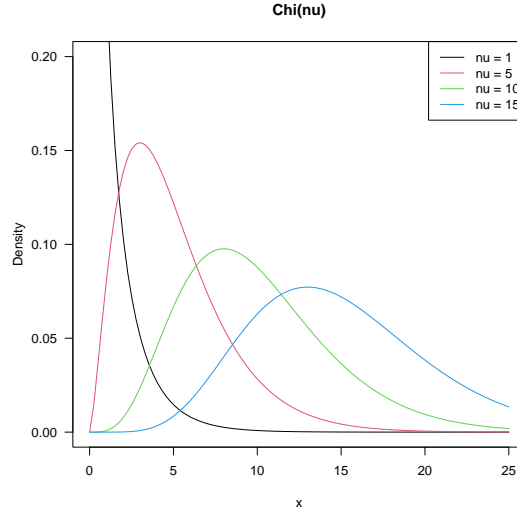
$$\begin{aligned} R_X &= \{x \in \mathbb{R} : x > 0\} \\ \Theta &= \{\nu \in \mathbb{N} \setminus \{0\}\} \end{aligned}$$

L'unico parametro della Chi-quadrato è ν detto *gradi di libertà*, intero positivo

Definition 6.5.1 (Funzione di densità).

$$f_X(x) = \frac{1}{2^{(\frac{\nu}{2})} \Gamma\left(\frac{\nu}{2}\right)} x^{(\frac{\nu}{2}-1)} e^{(-\frac{x}{2})} \cdot \mathbb{1}_{R_X}(x) \quad (6.31)$$

Remark 228. Anche se ν può esser qualsiasi numero reale positivo, in pratica le applicazioni hanno tipicamente ν intero positivo.

Figura 6.5: Distribuzione χ^2

6.5.3 Moments

Proposition 6.5.1 (Momenti caratteristici).

$$\mathbb{E}[X] = \nu \quad (6.32)$$

$$\text{Var}[X] = 2\nu \quad (6.33)$$

$$\text{Asym}(X) = \sqrt{\frac{8}{\nu}} \quad (6.34)$$

$$\text{Kurt}(X) = 3 + \frac{12}{\nu} \quad (6.35)$$

6.5.4 Shape

Remark 229 (Forma della distribuzione). La vc Chi-quadrato (figura 6.5)

- è asimmetrica positiva; al crescere di $\nu \rightarrow \infty$ l'asimmetria diminuisce e tende ad assumere una forma sempre più vicina alla Normale
- la forma della funzione di densità è monotona decrescente a zero se $\nu \leq 2$; se $\nu > 2$, presenta un picco intermedio in corrispondenza della moda (pari a $\nu - 2$).

6.5.5 Extras

Proposition 6.5.2. *Come la Gamma, anche la distribuzione Chi-quadrato gode della proprietà riproduttiva:*

$$\sum_{i=1}^n \chi_{\nu_i}^2 \sim \chi_{\sum_i \nu_i}^2$$

6.6 Beta

6.6.1 Definition

Remark 230. Viene utilizzata quando

- si studiano quantità aleatorie che stanno nel range 0-1 (ad esempio una probabilità): ad esempio nell'approccio bayesiano, dove si trattano le probabilità come una quantità aleatoria, si può usare la beta per definire la probabilità di successo per un esperimento bernoulliano;
- ciò che si analizza è destinato ad accadere in un intervallo delimitato (di tempo, costo, qualità) che può essere normalizzato nell'intervallo 0-1

6.6.2 Functions

Remark 231 (Supporto e spazio parametrico).

$$R_X = (0, 1)$$

$$\Theta = \{\alpha \in \mathbb{R} : \alpha > 0; \beta \in \mathbb{R} : \beta > 0\}$$

Definition 6.6.1 (Funzione di densità). Una vc continua X si definisce Beta con due parametri e la indichiamo con $X \sim \text{Beta}(\alpha, \beta)$ se la sua funzione di densità è:

$$f_X(x, \alpha, \beta) = \frac{x^{\alpha-1}(1-x)^{\beta-1}}{B(\alpha, \beta)} \cdot \mathbb{1}_{R_X}(x) \quad (6.36)$$

Definition 6.6.2 (Funzione Beta). La densità dipende dalla funzione beta, definita come

$$B(\alpha, \beta) = \int_0^1 x^{\alpha-1}(1-x)^{\beta-1} dx \quad (6.37)$$

Similmente a gamma

- dipende solamente da α, β , non da x per la quale integriamo
- non ammette soluzione analitica ma numerica (in R `beta(alpha, beta)`)
- presenta le seguenti proprietà

$$B(\alpha, \beta) = B(\beta, \alpha) \quad (6.38)$$

$$B(\alpha, \beta) = \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha + \beta)} \quad (6.39)$$

$$B(\alpha, \beta) = \frac{(\alpha-1)!(\beta-1)!}{(\alpha+\beta-1)!} \quad \text{se } \alpha, \beta \text{ sono interi positivi} \quad (6.40)$$

dove nell'ultima abbiamo solo sfruttato le proprietà viste in precedenza per la Gamma

6.6.3 Moments

Proposition 6.6.1 (Momenti caratteristici).

$$\mathbb{E}[X] = \frac{\alpha}{\alpha + \beta} \quad (6.41)$$

$$\text{Var}[x] = \frac{\alpha\beta}{(\alpha + \beta)^2(\alpha + \beta + 1)} \quad (6.42)$$

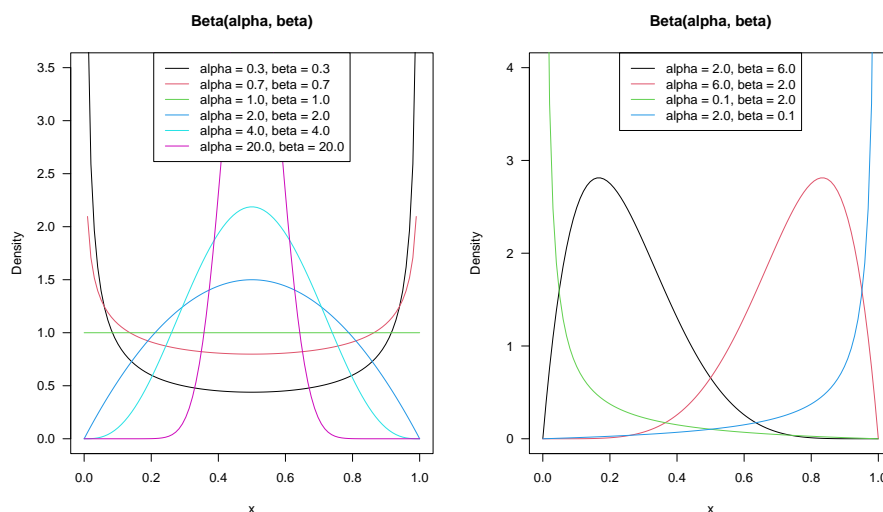


Figura 6.6: Distribuzione beta

6.6.4 Shape

Remark 232 (Forma della distribuzione). La forma (figura 6.6) dipende dai parametri α, β :

- se $\alpha = \beta$ la distribuzione è simmetrica rispetto al valore centrale $x = 1/2$;
 - nel caso particolare $\alpha = \beta = 1$, la distribuzione coincide con l'uniforme: $\text{Beta}(1, 1) \sim \text{Unif}(0, 1)$;
 - più $\alpha = \beta$ si avvicinano a 0, più la distribuzione assume la forma di U con code pesanti
 - più $\alpha = \beta$ si avvicinano a $+\infty$ maggiormente la distribuzione dà importanza ai valori centrali (tende a collapsare su $1/2$)
- se $\alpha \neq \beta$ il segno di $\beta - \alpha$ denota l'asimmetria (es se positivo l'asimmetria è positiva e la coda a destra, se negativo perché $\alpha > \beta$, così è l'asimmetria e la coda è a sinistra); scambiando α con β si inverte l'asse di simmetria.

6.6.5 Extras

Remark 233. Una vc Beta è definita nell'intervallo $[0, 1]$, ma effettuando la trasformazione $Y = X(b - a) + a$, la si può ricondurre all'intervallo $[a, b]$.

6.7 T di Student

6.7.1 Definition

Remark 234. Si ha che:

- il suo uso è prettamente teorico, in inferenza

- è la risultante di una trasformazione su due variabili, una normale e una chi quadrato.

6.7.2 Functions

Remark 235 (Supporto e spazio parametrico).

$$\begin{aligned} R_X &= \mathbb{R} \\ \Theta &= \{\nu \in \mathbb{N} \setminus \{0\}\} \end{aligned}$$

con ν detti *gradi di libertà*

Definition 6.7.1 (Funzione di densità).

TODO: vedi pag 205
mccoll

$$f_X(x) = \frac{\Gamma\left(\frac{\nu+1}{2}\right)}{\Gamma\left(\frac{\nu}{2}\right)\sqrt{\pi\nu}} \left(1 + \frac{x^2}{\nu}\right)^{-\frac{\nu+1}{2}} \cdot \mathbb{1}_{R_X}(x) \quad (6.43)$$

6.7.3 Moments

Proposition 6.7.1 (Momenti caratteristici).

$$\begin{aligned} \mathbb{E}[X] &= 0 \quad \text{se } \nu > 1 \\ \text{Var}[X] &= \frac{\nu}{\nu-2} \quad \text{se } \nu > 2 \\ \text{Kurt}(X) &= 3 + \frac{6}{\nu-4} \quad \text{se } \nu > 4 \end{aligned}$$

6.7.4 Shape

Remark 236 (Forma della distribuzione). Per (figura 6.7)

- ν piccolo si hanno code più alte rispetto alla normale; per questo detta variabile a code pesanti (considera più probabile rispetto alla normale cose che stanno nelle code)
- per $\nu \rightarrow \infty$ la distribuzione converge alla convergenza alla normale standardizzata; verso $\nu = 30$, l'approssimazione è già buona

6.7.5 Extras

Important remark 44 (Derivazione). Siano

$$\begin{aligned} Z &\sim N(0, 1) \\ C &\sim \chi^2(\nu) \\ Z &\perp\!\!\!\perp C \end{aligned}$$

Allora la seguente trasformazione definisce vc di Student con ν gradi di libertà

$$X = \frac{Z}{\sqrt{C/\nu}} \sim T(\nu) \quad (6.44)$$

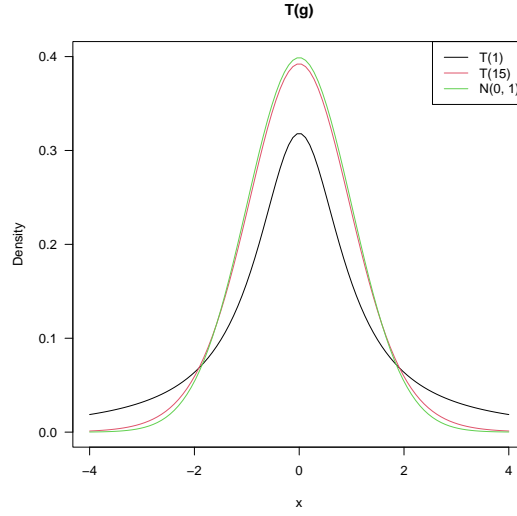


Figura 6.7: Distribuzione T

6.8 F di Fisher

6.8.1 Definition

Remark 237. Similmente alla T:

- il suo uso è prettamente teorico, in inferenza;
- è risultate di una trasformazione. È la distribuzione che deriva dal rapporto tra due vc Chi quadrato *indipendenti* tra loro e divise per i rispettivi gradi di libertà.

6.8.2 Functions

Remark 238 (Supporto e spazio parametrico).

$$R_X = \{x \in \mathbb{R} : x > 0\}$$

$$\Theta = \{\nu_1, \nu_2 \in \mathbb{N} \setminus \{0\}\}$$

Definition 6.8.1 (Funzione di densità).

$$f_X(x) = \frac{\Gamma\left(\frac{\nu_1 + \nu_2}{2}\right)}{\Gamma\left(\frac{\nu_1}{2}\right)\Gamma\left(\frac{\nu_2}{2}\right)} \cdot \left(\frac{\nu_1}{\nu_2}\right)^{\frac{\nu_1}{2}} \cdot \frac{x^{(\nu_1 - 2)/2}}{\left(1 + \frac{\nu_1}{\nu_2}x\right)^{\frac{\nu_1 + \nu_2}{2}}} \cdot \mathbb{1}_{R_X}(x) \quad (6.45)$$

Remark 239 (Funzione di ripartizione). Anche per la F non vi è una forma chiusa della ripartizione e ci si affida alle tavole.

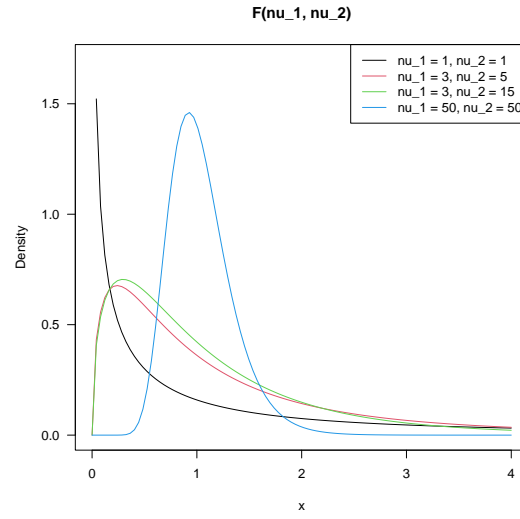


Figura 6.8: Distribuzione F

6.8.3 Moments

Proposition 6.8.1 (Momenti caratteristici).

$$\mathbb{E}[X] = \frac{\nu_2}{\nu_2 - 2} \quad \text{se } \nu_2 > 2$$

$$\text{Var}[X] = \frac{2\nu_2^2(\nu_1 + \nu_2 - 2)}{\nu_1(\nu_2 - 2)^2(\nu_2 - 4)} \quad \text{se } \nu_2 > 4$$

6.8.4 Shape

Remark 240 (Forma della distribuzione). Si nota che (figura 6.8)

- se $\nu_1 = \nu_2 = 1$ la funzione è monotona decrescente (asimmetria positiva)
- se $\nu_1, \nu_2 \neq 1$ la funzione è asimmetrica positiva.
- la distribuzione converge a quella di una normale solo se contemporaneamente $\nu_1 \rightarrow \infty$ e $\nu_2 \rightarrow \infty$.

6.8.5 Extras

Important remark 45 (Derivazione). Siano

$$C_1 \sim \chi^2(\nu_1)$$

$$C_2 \sim \chi^2(\nu_2)$$

$$C_1 \perp\!\!\!\perp C_2$$

Allora

$$X = \frac{C_1/\nu_1}{C_2/\nu_2} \sim F(\nu_1, \nu_2) \quad (6.46)$$

detto X si distribuisce come una F con ν_1 e ν_2 gradi di libertà (detti anche ordinatamente gradi di libertà del numeratore e del denominatore).

6.9 Logistica

6.9.1 Definition

Remark 241. Viene utilizzata per *modelli di crescita di grandezze aleatorie nel tempo*. La logistica si usa se la grandezza aleatoria è caratterizzata da una crescita che può essere divisa in tre fasi:

- la prima di crescita esponenziale,
- la seconda di saturazione (dove la crescita diminuisce)
- la terza di arresto (maturità, dove non vi è più crescita).

Un buon modello per rappresentare fenomeni di questo tipo è rappresentato dalla funzione di ripartizione logistica.

È matematicamente semplice e ci permette di focalizzarci su aspetti non numerici; è altresì importante nella regressione logistica.

6.9.2 Functions

Remark 242 (Supporto e spazio parametrico).

$$R_X = \mathbb{R}$$

$$\Theta = \{\mu \in \mathbb{R}, s \in \mathbb{R} : s > 0\}$$

Definition 6.9.1 (Funzione di densità). La funzione di densità di una vc $X \sim \text{Logistic}(\mu, s)$ è

$$f_X(x) = \frac{e^{-\frac{x-\mu}{s}}}{s \left(1 + e^{-\frac{x-\mu}{s}}\right)^2} \cdot \mathbb{1}_{R_X}(x) \quad (6.47)$$

Definition 6.9.2 (Funzione di ripartizione). La funzione di ripartizione di una vc $X \sim \text{Logistic}(\mu, s)$ è

$$F_X(x) = \frac{e^{\frac{x-\mu}{s}}}{1 + e^{\frac{x-\mu}{s}}} \cdot \mathbb{1}_{R_X}(x) \quad (6.48)$$

6.9.3 Moments

Proposition 6.9.1 (Momenti caratteristici).

$$\mathbb{E}[X] = \mu$$

$$\text{Var}[X] = \frac{\pi^2}{3} s^2$$

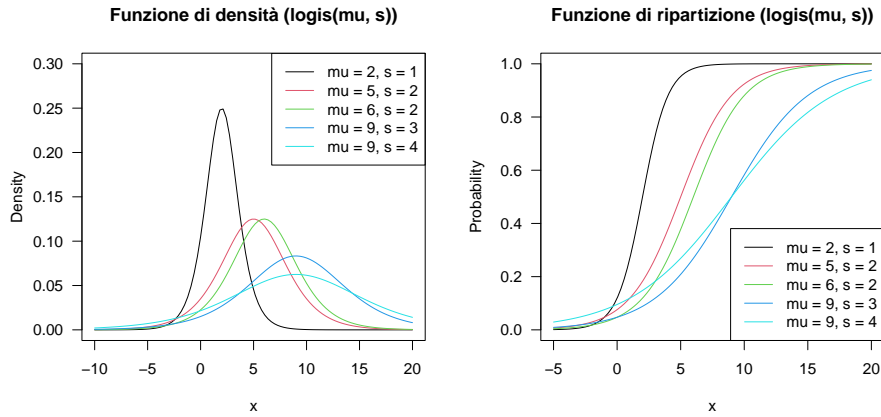


Figura 6.9: Distribuzione logistica

6.9.4 Shape

Remark 243. Per la forma (6.9):

- la distribuzione è simmetrica e risente dei valori di μ ed s
- μ è effettivamente il valore atteso/centro della distribuzione, mentre s non è la deviazione standard (radice della varianza) ma comunque un parametro proporzionale ad essa

6.10 Lognormale

6.10.1 Definition

Remark 244. Questa distribuzione:

- si chiama lognormale, ha un legame con la gaussiana e presenta gli stessi parametri μ, σ^2 (che possono assumere gli stessi valori);
- a differenza della gaussiana ammette solo realizzazioni positive;
- viene utilizzata quando la grandezza aleatoria oggetto di studio è/può essere visto come il risultato del *prodotto* di un *numero* molto *elevato* di *fattori* tra loro *indipendenti*, che agiscono tra loro in maniera moltiplicativa (mentre per la gaussiana è la *somma*)

6.10.2 Functions

Remark 245 (Supporto e spazio parametrico).

$$R_X = \{x \in \mathbb{R} : x > 0\}$$

$$\Theta = \{\mu \in \mathbb{R}, \sigma^2 \in \mathbb{R} : \sigma^2 > 0\}$$

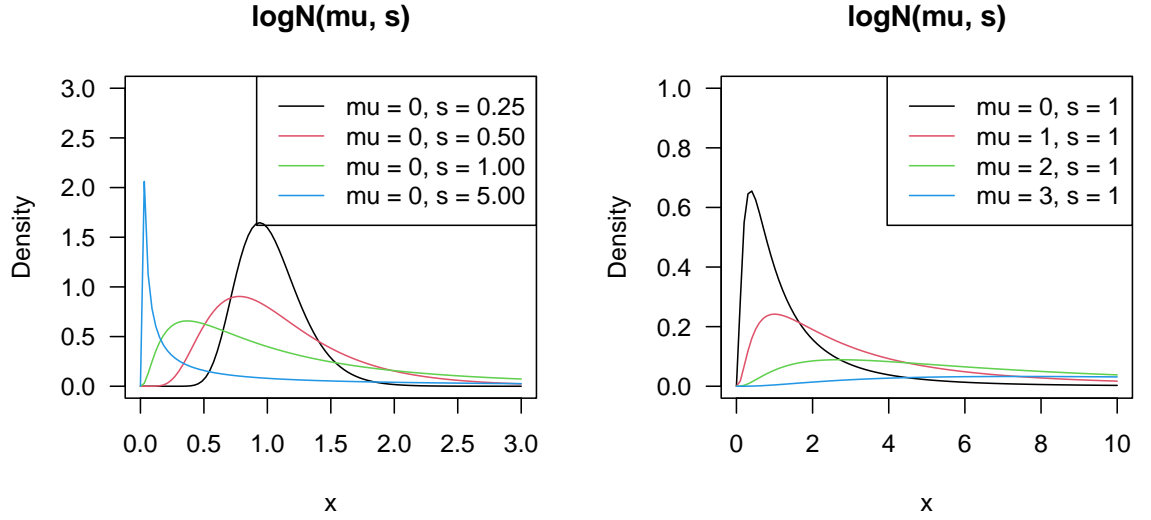


Figura 6.10: Distribuzione lognormale

Definition 6.10.1 (Funzione di densità).

$$f_X(x) = \frac{1}{x\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2}\left(\frac{\log x - \mu}{\sigma}\right)^2} \cdot \mathbb{1}_{R_X}(x) \quad (6.49)$$

6.10.3 Moments

Proposition 6.10.1 (Momenti caratteristici).

$$\begin{aligned} \mathbb{E}[X] &= e^{\mu + \frac{\sigma^2}{2}} \\ \text{Var}[X] &= e^{2\mu + 2\sigma^2} - e^{2\mu + \sigma^2} \end{aligned}$$

6.10.4 Shape

Remark 246 (Forma della distribuzione). In figura 6.10

- nella prima con $\mu = 0$ all'aumentare di σ l'asimmetria si incrementa (aumenta la probabilità di valori vicino a zero)
- nella seconda se a parità di σ aumentiamo μ chiaramente la distribuzione si sposta più in alto ma (e non si vede dal grafico) rimane comunque asimmetrica con coda a destra

6.10.5 Extras

Important remark 46 (Legame con la gaussiana). Si ha che se applico trasformazioni passo da una distribuzione all'altra:

$$\begin{aligned} X \sim \text{LogN}(\mu, \sigma^2) &\implies \log X \sim N(\mu, \sigma^2) \\ X \sim N(\mu, \sigma^2) &\implies e^X \sim \text{LogN}(\mu, \sigma^2) \end{aligned}$$

6.11 Weibull

6.11.1 Definition

Remark 247. Viene utilizzata per studiare l'affidabilità dei sistemi di produzione nei processi industriali, in particolare per valutare i tempi di rottura (tempi ad un evento) ed è generalizzazione dell'esponenziale

6.11.2 Functions

Remark 248 (Supporto e spazio parametrico).

$$\begin{aligned} R_X &= \{x \in \mathbb{R} : x > 0\} \\ \Theta &= \{a \in \mathbb{R} : a > 0, b \in \mathbb{R} : b > 0\} \end{aligned}$$

Definition 6.11.1 (Funzione di densità).

$$f_X(x) = \frac{a}{b} \left(\frac{x}{b}\right)^{a-1} e^{-(\frac{x}{b})^a} \cdot \mathbb{1}_{R_X}(x) \quad (6.50)$$

6.11.3 Moments

Proposition 6.11.1 (Momenti caratteristici).

$$\begin{aligned} \mathbb{E}[X] &= \frac{\Gamma(1 + \frac{1}{b})}{a^{1/b}} \\ \text{Var}[X] &= \frac{\Gamma(1 + \frac{2}{b}) - \Gamma^2(1 + \frac{1}{b})}{a^{2/b}} \end{aligned}$$

6.11.4 Shape

Remark 249 (Forma della funzione). Come detto è usata nei processi produttivi industriali dove si vuole studiare l'affidabilità dei sistemi di produzione. La distribuzione utile per analizzare le rotture che si possono verificare in un processo produttivo. Le situazioni che questa distribuzione riesce a gestire sono 3 in funzione del valore di a , che ne determina la forma (figura 6.11):

- se $a < 1$ il sistema è caratterizzato da *mortalità infantile*: le rotture si verificano tipicamente all'inizio quando il processo produttivo viene avviato. Questo può essere dovuto al fatto che nel processo produttivo ci sono dei pezzi difettosi; una volta che vengono sostituite il processo produttivo va a regime e le rotture non si verificano più all'inizio del processo produttivo

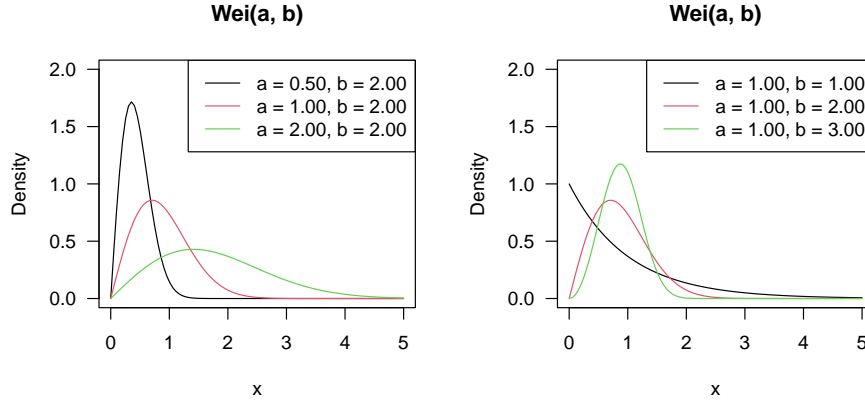


Figura 6.11: Distribuzione Weibull

- se $a = 1$ il tasso di rottura è *costante nel tempo*: le rotture si verificano in maniera casuale, ovvero regolari nel tempo (e la distribuzione coincide con una esponenziale di parametro $1/b$, ossia $\text{Weibull}(1, b) \sim \text{Exp}(\frac{1}{b})$)
- se $a > 1$ il sistema è affetto da *logoramento* e il tasso di rottura è *crescente nel tempo*: all'inizio va tutto bene ma più il tempo passa più le parti coinvolte nel processo produttivo sono sottoposte a usura e quindi le rotture crescono col passare del tempo

Apparently anche con l'incremento di b a parità di a la distribuzione diviene meno asimmetrica.

Se si aumentano i valori delle due costanti si passa da una forte asimmetria positiva ad una maggiore simmetria.

6.11.5 Extras

Example 6.11.1 (Legame con l'esponenziale). Se $a = 1$:

$$\begin{aligned} f_X(x) &= \frac{1}{b} \left(\frac{x}{b}\right)^{1-1} e^{-\left(\frac{x}{b}\right)^1} \cdot \mathbb{1}_{R_X}(x) \\ &= \frac{1}{b} e^{-\frac{x}{b}} \end{aligned}$$

Ricordando che l'esponenziale ha densità

$$f_X(x) \lambda e^{-\lambda x}$$

la Weibull diviene una esponenziale con parametro $\lambda = \frac{1}{b}$. Ossia $\text{Weibull}(1, b) \sim \text{Exp}(\frac{1}{b})$.

6.12 Pareto

6.12.1 Definition

Remark 250. Viene utilizzata quando si studiano distribuzioni di variabili che hanno un minimo (ad esempio come, con x_m = reddito minimo).

Storicamente ideata da Pareto nel contesto della distribuzione della ricchezza/reddito all'interno di popolazioni umane. La distribuzione è sviluppata a partire dal principio dell'80/20: il 20% della popolazione (più ricca) detiene l'80% della ricchezza. La distribuzione è utile per descrivere la distribuzione della ricchezza in una popolazione umana.

6.12.2 Functions

Remark 251 (Supporto e spazio parametrico).

$$\begin{aligned} R_X &= (x_m, +\infty) \\ \Theta &= \{x_m, k \in \mathbb{R} : x_m, k > 0\} \end{aligned}$$

Definition 6.12.1 (Funzione di densità). Se $X \sim \text{Pa}(k, x_m)$

$$f_X(x) = k \frac{x_m^k}{x^{k+1}} \cdot \mathbb{1}_{R_X}(x) \quad (6.51)$$

6.12.3 Moments

Proposition 6.12.1 (Momenti caratteristici).

$$\begin{aligned} \mathbb{E}[X] &= \frac{kx_m}{k-1} \quad \text{per } k > 1 \\ \text{Var}[X] &= \left(\frac{x_m}{k-1}\right)^2 \frac{k}{k-2} \quad \text{per } k > 2 \end{aligned}$$

6.12.4 Shape

Remark 252 (Forma della distribuzione e significato parametri). Vediamo che

- il supporto è delimitato inferiormente da x_m , uno dei due parametri della distribuzione, che dunque assume significato di valore più dei valori possibili per la variabile casuale (es se X è il reddito, x_m il reddito della persona più povera).
- da notare che se ci fosse ripartizione equa del reddito/ricchezza la curva sarebbe piatta (tutti i valori di ricchezza equivalentemente probabili). Notiamo in 6.12 che fissando $x_m = 1$ e incrementando k la disuguaglianza (assimmetria positiva) aumenta perchè aumenta la probabilità di redditi bassi. Dunque k può essere interpretato come *indicatore della disuguaglianza* della distribuzioni dei redditi a parità di x_m ; al crescere di k la distribuzione diviene più disuguale (molto probabile trovare valori vicini al limite inferiore x_m , meno valori molto grandi).

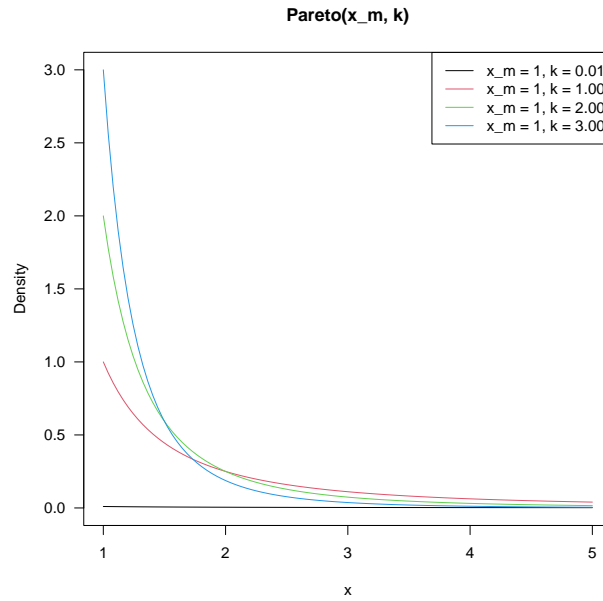


Figura 6.12: Distribuzione di Pareto

6.13 Exercises

Example 6.13.1 (Exponential (crash course, giorno 1)). Let X be a rv that has the density

$$f(x) = \begin{cases} ce^{-\lambda x} & \text{if } x \geq 0 \\ 0 & \text{if } x < 0 \end{cases}$$

Find:

1. c
2. $\mathbb{E}[X]$
3. $\text{Var}[X]$
4. $F(X)$

We have

1. it must be that

$$\begin{aligned} 1 &= \int_{-\infty}^{+\infty} f(x) \, dx = \int_0^{+\infty} ce^{-\lambda x} \, dx = c \int_0^{+\infty} e^{-\lambda x} \, dx = c \left[\left(\frac{1}{-\lambda} e^{-\lambda x} \right) \right]_0^{+\infty} \\ &= 0 - \frac{c}{-\lambda} \cdot 1 \end{aligned}$$

therefore $c = \lambda$ (this is the exponential distribution)

2. we have,

$$\mathbb{E}[X] = \int_0^{+\infty} x \cdot \lambda e^{-\lambda x} dx = \lambda \int_0^{+\infty} x \cdot e^{-\lambda x} dx$$

using integration by parts we have

$$\begin{aligned} \int x e^{-\lambda x} dx &= x \left(-\frac{1}{\lambda} e^{-\lambda x} \right) - \int -\frac{1}{\lambda} e^{-\lambda x} \\ &= \left(-\frac{x}{\lambda} e^{-\lambda x} \right) + \frac{1}{\lambda} \int e^{-\lambda x} \\ &= \left(-\frac{x}{\lambda} e^{-\lambda x} \right) + \frac{1}{\lambda} \left(-\frac{1}{\lambda} e^{-\lambda x} \right) \end{aligned}$$

che opportunamente valutato

$$\left[\left(-\frac{x}{\lambda} e^{-\lambda x} \right) + \frac{1}{\lambda} \left(-\frac{1}{\lambda} e^{-\lambda x} \right) \right]_0^{+\infty} = 0 + 0 - \left(0 - \frac{1}{\lambda^2} \right)$$

Per cui tornando al valore atteso

$$\mathbb{E}[X] = \lambda \left(\frac{1}{\lambda^2} \right) = \frac{1}{\lambda}$$

3. first we find $\mathbb{E}[X^2]$

$$\begin{aligned} \mathbb{E}[X^2] &= \lambda \int_{-\infty}^{+\infty} x^2 e^{-\lambda x} dx \stackrel{(1)}{=} \lambda \left[\left(x^2 \frac{-1}{\lambda} e^{-\lambda x} \right) \Big|_0^{\infty} + \frac{2}{\lambda} \int_0^{+\infty} x e^{-\lambda x} dx \right] \\ &= 2 \int_0^{+\infty} x e^{-\lambda x} dx = \frac{2}{\lambda^2} \end{aligned}$$

where in (1) again by integration by parts. So

$$\text{Var}[X] = \mathbb{E}[X^2] - \mathbb{E}[X]^2 = \frac{2}{\lambda^2} - \frac{1}{\lambda^2} = \frac{1}{\lambda^2}$$

4. we have

$$\begin{aligned} F(x) &= \int_0^x f(s) ds = \lambda \int_0^x e^{-\lambda s} ds = \lambda \left(-\frac{1}{\lambda} e^{-\lambda s} \right) \Big|_0^x \\ &= 1 - e^{-\lambda x}, \quad \text{for } x \geq 0 \end{aligned}$$

for $x < 0$, $F(x) = \int_{-\infty}^x f(s) ds = 0$ so

$$F(x) = \begin{cases} 0 & x < 0 \\ 1 - e^{-\lambda x} & x \geq 0 \end{cases}$$

Example 6.13.2 (Standard normal (crashcourse, day 1 es 5 pag 8)). Let $f(x) = ce^{-\frac{x^2}{2}}$, $x \in \mathbb{R}$:

1. find c
2. $\mathbb{E}[X]$
3. $\text{Var}[X]$

Respectively

1. we know $\int_{-\infty}^{+\infty} cf(x) dx = 1$ so we can do this trick

$$\begin{aligned} 1 &= \underbrace{\int_{-\infty}^{+\infty} cf(x) dx}_1 \underbrace{\int_{-\infty}^{+\infty} cf(y) dy}_1 \\ &= c^2 \int_{-\infty}^{+\infty} f(x)f(y) dx dy = c^2 \int_{-\infty}^{+\infty} e^{-\frac{x^2+y^2}{2}} dx dy \end{aligned}$$

Now transforming variable to polar coordinates that is applying

$$\begin{cases} x = r \cos \theta \\ y = r \sin \theta \end{cases}, \quad r \in [0, \infty), \theta \in [0, 2\pi)$$

so that $x^2 + y^2 = r^2$ and $dx dy = r dr d\theta$ we have

$$\begin{aligned} 1 &= c^2 \int_0^{2\pi} \int_0^{\infty} e^{-\frac{r^2}{2}} r dr d\theta \stackrel{(1)}{=} c^2 \int_0^{2\pi} \underbrace{\int_0^{\infty} e^{-u} du}_{=1} d\theta \\ &= c^2 \int_0^{2\pi} d\theta = c^2 2\pi = 1 \implies c = \frac{1}{\sqrt{2\pi}} \end{aligned}$$

where in (1) we substitute $u = \frac{r^2}{2}$ so $du = r dr$.

2. $\mathbb{E}[X] = \int_{-\infty}^{+\infty} xf(x) dx$. We have that $f(x)$ is an even function:

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

it's symmetric. However we are interested in $\mathbb{E}[X] = \int_{-\infty}^{+\infty} xf(x) dx$ that is trying to find the area under an odd function. Now in general if we're trying to find

- odd functions: given that it's symmetric around origin, positive areas compensates with negative areas so it's integral (over \mathbb{R}) is 0 (this holds for any odd function).
- even functions: since symmetric around y axis to calculate integral on region $(-\infty, \infty)$ we can double the integral on region $(0, \infty)$

Therefore our $\mathbb{E}[X] = 0$.

3. for the variance first get

$$\begin{aligned}\mathbb{E}[X^2] &= \int_{-\infty}^{+\infty} \overbrace{x^2}^{\text{even}} \underbrace{f(x)}_{\text{even}} dx \\ &= \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{+\infty} x^2 e^{-\frac{x^2}{2}} dx \stackrel{(1)}{=} \frac{2}{\sqrt{2\pi}} \int_0^{+\infty} x^2 e^{-\frac{x^2}{2}} dx \stackrel{(2)}{=} \frac{2}{\sqrt{2\pi}} \int_0^{+\infty} \sqrt{2} u^{1/2} e^{-u} du \\ &= \frac{2}{\sqrt{\pi}} \int_0^{+\infty} u^{1/2} e^{-u} du = \frac{2}{\sqrt{\pi}} \Gamma\left(\frac{3}{2}\right) = 1\end{aligned}$$

where in (1) since it's even and (2) using the variable change $u = \frac{x^2}{2}$ therefore $du = x dx$ and $x = \sqrt{2u}$.

$\Gamma(x)$ is called gamma function, we will be familiar with it in the next years, for now trust me that $\Gamma(x) = (x-1)\Gamma(x-1)$ and $\Gamma(1) = 1$ $\Gamma(1/2) = \sqrt{\pi}$ so for integers n $\Gamma(n) = (n-1)!$ but for our case $\Gamma(3/2) = \frac{1}{2}\Gamma(\frac{1}{2}) = \frac{\sqrt{\pi}}{2}$

Therefore in the end for $X \sim N(0, 1)$, $\mathbb{E}[X] = 0$ and $\text{Var}[X] = 1$. This is called a standard rv. But in general normal rvs can have different mean and variance: the general case is denoted as $X \sim N(\mu, \sigma^2)$, $\mu \in \mathbb{R}$, $\sigma^2 \in \mathbb{R}^+$ and this correspond to translation of a standard normal rv and then scaling it.

Let $Z \sim N(0, 1)$ and $X = \sigma Z + \mu$ then

$$\begin{aligned}\mathbb{E}[X] &= \mathbb{E}[\sigma Z + \mu] = \sigma \underbrace{\mathbb{E}[Z]}_{=0} + \mu = \mu \\ \text{Var}[X] &= \text{Var}[\sigma Z + \mu] = \sigma^2 \underbrace{\text{Var}[Z]}_{=1} = \sigma^2\end{aligned}$$

Example 6.13.3 (Beta expected value). Vediamo il valore atteso di $X \sim \text{Beta}(\alpha, \beta)$

$$\begin{aligned}\mathbb{E}[X] &= \int_0^1 x \frac{1}{B(\alpha, \beta)} x^{\alpha-1} (1-x)^{\beta-1} dx \\ &\stackrel{(1)}{=} \frac{1}{B(\alpha, \beta)} \int_0^1 x^\alpha (1-x)^{\beta-1} dx \\ &= \frac{1}{B(\alpha, \beta)} \int_0^1 B(\alpha+1, \beta) dx = \frac{B(\alpha+1, \beta)}{B(\alpha, \beta)} \\ &\stackrel{(2)}{=} \frac{1}{\frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha+\beta)}} \frac{\Gamma(\alpha+1)\Gamma(\beta)}{\Gamma(\alpha+\beta+1)} \\ &= \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)} \cdot \frac{\Gamma(\alpha+1)}{\Gamma(\alpha+\beta+1)} \\ &\stackrel{(3)}{=} \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha)} \frac{\alpha\Gamma(\alpha)}{(\alpha+\beta)\Gamma(\alpha+\beta)} \\ &= \frac{\alpha}{\alpha+\beta}\end{aligned}$$

dove

- in (1) abbiamo portato fuori quello che non dipende da x ($B(\alpha, \beta)$ dipende solo dalle due costanti e non da x dato che nella definizione della funzione

speciale B si calcola l'integrale su x che varia) e notato che entro parentesi si ha una espressione simile a $B(\alpha, \beta) = \int_0^1 x^{\alpha-1} (1-x)^{\beta-1} dx$

- in (2) usiamo eq 6.39
- in (3) usiamo eq 6.23

In questo caso non abbiamo avuto bisogno di calcolare l'integrale, ma abbiamo sfruttato le proprietà delle funzioni speciali.

Nel particolare caso $\alpha = \beta$ si ha che $\mathbb{E}[X] = 1/2$ e la distribuzione è simmetrica per cui $f_X(1/2 - c) = f_X(1/2 + c)$

Example 6.13.4 (Esame vecchio viroli). A random variable X is distributed according to $N(0, 2)$ where 2 is the variance. What is the distribution of $Y = 2X$? Il risultato è $Y \sim N(0, 8)$ (come confermato dal Bigo).

Example 6.13.5 (Esame vecchio viroli). A random variable X is distributed according to $N(-1, 1)$. What is the distribution of $Y = -2X + 1$. Correct answer is $Y \sim N(3, 4)$

Example 6.13.6 (Esame vecchio viroli). Let $X \sim N(0, 2)$ and $Y \sim N(1, 1)$ be independent random variables where the parameters in the bracket are the expectation and the variace. What is the distribution of $Z = 2X + Y$

1. $Z \sim N(1, 9)$
2. $Z \sim N(1, 5)$
3. not possible to determine
4. $Z \sim N(1, 2)$

should be the first

Capitolo 7

Misc topics

7.1 Quantili

Definition 7.1.1. Sia X una va univariata e sia $\alpha \in (0, 1)$. Un numero $b \in \mathbb{R}$ si dice quantile di ordine α di X se: NB: Rigo, dalla triennale

$$\mathbb{P}(X \leq b) \geq \alpha \geq \mathbb{P}(X < b)$$

Equivalentemente. detta F la funzione di ripartizione di X

$$F(b) \geq \alpha \geq F(b^-)$$

Dove $F(b^-) = \lim_{x \rightarrow b^-} F(x)$

Remark 253. Se

- F è continua, la condizione precedente diviene $F(b) = \alpha$
- F è continua e strettamente crescente $\forall \alpha \in (0, 1)$ esiste uno e un solo quantile di ordine α , ovvero $\exists! b \in \mathbb{R}$ tale che $F(b) = \alpha$. In particolare, ciò è vero se X è assolutamente continua con densità strettamente positiva

Remark 254. Se

- $\alpha = 1/2$ un quantile di ordine $\frac{1}{2}$ si chiama mediana
- $\alpha = 1/4$ un quantile di ordine $\frac{1}{4}$ si chiama primo quartile
- $\alpha = 3/4$ un quantile di ordine $\frac{3}{4}$ si chiama terzo quartile

Remark 255. In ogni caso, al di là della definizione formale, detto in parole povere, un quantile di ordine α è un qualsiasi valore b che lasci alla propria sinistra un'area (della densità) pari a α . Questo si vede bene se X è assolutamente continua con densità > 0

Example 7.1.1. Sia $X \sim N(\mu, \sigma^2)$ e sia b il quantile di ordine α , allora

$$\alpha = \mathbb{P}(X \leq b) = \mathbb{P}\left(\frac{X - \mu}{\sigma} \leq \frac{b - \mu}{\sigma}\right) = \Phi\left(\frac{b - \mu}{\sigma}\right)$$

Dove Φ è la funzione di ripartizione della $N(0, 1)$. Due commenti:

- se b è quantile di ordine α per una $N(\mu, \sigma^2)$ allora $\frac{b-\mu}{\sigma}$ è quantile di ordine α per una $N(0, 1)$. Naturalmente vale anche il viceversa, ovvero se c è quantile di ordine α per una $N(0, 1)$, allora $\mu + \sigma c$ è quantile di ordine α per una $N(\mu, \sigma^2)$. Infatti se $X \sim N(\mu, \sigma^2)$

$$\mathbb{P}(X \leq \mu + \sigma c) = \mathbb{P}\left(\frac{X - \mu}{\sigma} \leq c\right) = \Phi(c) = \alpha$$

- in queste considerazioni abbiamo usato il fatto che

$$X \sim N(\mu, \sigma^2) \implies \alpha + \beta X \sim N(\cdot, \cdot), \quad \text{for all } \alpha \in \mathbb{R}, \forall \beta \neq 0$$

Ovvero ad eccezione del caso banale $\beta = 0$, una qualsiasi trasformazione lineare di X è ancora normale.

7.2 Order statistics

Remark 256. Together with *rank* statistics, *order* statistics are fundamental tools in non-parametric statistics and inference.

Definition 7.2.1 (Order statistics). Let $X = (X_1, \dots, X_n)^\top$ be any n -variate random variable. The corresponding order statistics are the element of the vector $Y = (X_{(1)}, \dots, X_{(n)})^\top$ where $X_{(1)} \leq \dots \leq X_{(n)}$ are the elements of X arranged in non decreasing order, that is the following random variables

$$\begin{aligned} X_{(1)} &= \min\{X_1, \dots, X_n\} \\ X_{(2)} &= \min\{\{X_1, \dots, X_n\} \setminus \{X_{(1)}\}\} \\ &\dots \\ X_{(n)} &= \max\{X_1, \dots, X_n\} \end{aligned}$$

Example 7.2.1. Se $n = 4$, $X_1 = 2, X_2 = 0, X_3 = 5, X_4 = 1$ allora $X_{(1)} = -1, X_{(2)} = 0, X_{(3)} = 2, X_{(4)} = 5$.

Con $n = 3$, $X_1 = 1, X_2 = -2, X_3 = 1$, $X_{(1)} = -2, X_{(2)} = X_{(3)} = 1$.

Remark 257. La terminologia deriva dalla statistica: basta pensare ad X come a un campione e ad Y come al campione ordinato.

Remark 258. Here the random vector can be conceptualized as measurement on the same variable for different units (not several measurement within one unit).

Definition 7.2.2 (k -th order statistic). The k -th order statistic of the sample is equal to its k -th smallest value.

Example 7.2.2 (Minimum and maximum). Important special cases of the order statistics are the *minimum* $X_{(1)}$, the *maximum* $X_{(n)}$, the sample *median* and other sample *quantiles*.

Example 7.2.3. Throwing a dice 6 times, having the sequence X_1, \dots, X_6 . To study the distribution of the minimum $X_{(1)}$, we can say that

$$\begin{aligned} \mathbb{P}(X_{(1)} = 6) &= \frac{1}{6} \cdot \dots \cdot \frac{1}{6} = \left(\frac{1}{6}\right)^6 \\ \mathbb{P}(X_{(1)} = 1) &= 1 - \left(\frac{5}{6}\right)^6 \end{aligned}$$

Important remark 47 (Our focus). We:

- are interested in studying distribution/properties of these newly defined random variables, or in general, given the distribution of X , find the distribution of Y (note this another example of the general problem of transforming variable)
- deal with the simplest case where X_1, \dots, X_n are iid

7.2.1 Minimum

Proposition 7.2.1 (Distribution function). *We have that*

$$F_{(1)}(x) = 1 - [1 - F_X(x)]^n \quad (7.1)$$

NB: Direi sia roba di Viroli, Rigo l'ha fatto come caso particolare di $X_{(i)}$

Dimostrazione.

$$\begin{aligned} F_{(1)}(x) &= \mathbb{P}(X_{(1)} \leq x) = 1 - \mathbb{P}(X_{(1)} > x) \\ &= 1 - \mathbb{P}(X_1 > x, X_2 > x, \dots, X_n > x) \stackrel{(1)}{=} 1 - \prod_{i=1}^n \mathbb{P}(X_i > x) \\ &\stackrel{(2)}{=} 1 - \prod_{i=1}^n \mathbb{P}(X > x) = 1 - [\mathbb{P}(X > x)]^n = 1 - [1 - \mathbb{P}(X \leq x)]^n \\ &= 1 - [1 - F_X(x)]^n \end{aligned}$$

with (1) we considered independent rvs and (2) identically distributed. \square

Remark 259. Interpretazione affinché il minimo sia al più x si fa il complemento in cui si guarda la probabilità che siano tutte contemporaneamente $> x$

Proposition 7.2.2 (Density function).

$$f_{(1)}(x) = n f_X(x) \cdot [1 - F_X(x)]^{n-1}$$

Dimostrazione.

$$\begin{aligned} f_{(1)}(x) &= \frac{\partial F_{(1)}(x)}{\partial x} = -n [1 - F_X(x)]^{n-1} (-f_X(x)) \\ &= n f_X(x) \cdot [1 - F_X(x)]^{n-1} \end{aligned}$$

\square

Example 7.2.4. A room is lit by 5 light bulbs, each bulb lifetime has a distribution $X \sim \text{Exp}(\lambda = \frac{1}{100})$. What is the probability that after 200 days *all the bulbs are still working*?

We can setup this as $\mathbb{P}(X_{(1)} > 200)$, therefore:

$$\mathbb{P}(X_{(1)} > 200) = 1 - \mathbb{P}(X_{(1)} \leq 200) = 1 - F_{(1)}(200)$$

we have that, being X an exponential

$$F_{(1)}(200) = 1 - (1 - F_X(200))^5 = 1 - \left(1 - 1 + e^{-200/100}\right)^5 = 1 - \frac{1}{e^{10}}$$

Therefore

$$\mathbb{P}(X_{(1)} > 200) = 1 - 1 + \frac{1}{e^{10}} = \frac{1}{e^{10}}$$

Example 7.2.5 (Viols eserciziario 1, es 6). Let X_1, \dots, X_n be a random sample from a Weibull (α, β) distribution, that is

$$f(x) = \alpha\beta x^{\beta-1} e^{-\alpha x^\beta}, \quad x > 0, \alpha, \beta > 0$$

Derive the probability density function of $X_{(1)}$ and recognize it. The distribution function of a Weibull rv is

$$F_X(x) = 1 - e^{-\alpha x^\beta}$$

therefore

$$F_{(1)}(x) = 1 - [1 - F_X(x)]^n = 1 - [e^{-\alpha x^\beta}]^n = 1 - e^{-n\alpha x^\beta}$$

which is a weibull with parameters $n\alpha$ and β

7.2.2 Maximum

NB: Direi sia roba di Virolì, Rigo l'ha fatto come caso particolare di $X_{(i)}$

Proposition 7.2.3 (Distribution function).

$$F_{(n)}(x) = [F_X(x)]^n \quad (7.2)$$

Dimostrazione.

$$\begin{aligned} F_{(n)}(x) &= \mathbb{P}(X_{(n)} \leq x) = \mathbb{P}(X_1 \leq x, \dots, X_n \leq x) \\ &\stackrel{(iid)}{=} [\mathbb{P}(X \leq x)]^n = [F_X(x)]^n \end{aligned}$$

□

Remark 260. Il massimo sia $\leq x$ se tutte le vc sono $\leq x$

Proposition 7.2.4 (Density function).

$$f_{(n)}(x) = n [F_X(x)]^{n-1} f_X(x) \quad (7.3)$$

Dimostrazione.

$$f_{(n)}(x) = \frac{\partial}{\partial x} F_{(n)}(x) = n [F_X(x)]^{n-1} f_X(x)$$

□

Example 7.2.6. Considering again a room lit by 5 light bulbs, each bulb life-time has a distribution $X \sim \text{Exp}(\lambda = \frac{1}{100})$. What is the probability that after 200 days *at least a bulb will be working*?

This can be setup with

$$\begin{aligned} \mathbb{P}(X_{(n)} > 200) &= 1 - \mathbb{P}(X_{(n)} \leq 200) = 1 - F_{(n)}(200) \\ &= 1 - [F_X(200)]^5 = 1 - (1 - e^{-2})^5 \simeq 0.52 \end{aligned}$$

Example 7.2.7. Draw randomly 12 numbers between from $X \sim \text{Unif}(0, 1)$. What is the probability that at least a number > 0.9 ?

If $X \sim \text{Unif}(0, 1)$, $F_X(x) = x$. We have

$$\mathbb{P}(X_{(n)} > 0.9) = 1 - \mathbb{P}(X_{(n)} \leq 0.9) = 1 - [F_X(0.9)]^{12} = 1 - 0.9^{12} = 0.718$$

Example 7.2.8 (Esame vecchio viroli). A random variable X has density function

$$f(x, \theta) = \frac{3x^2}{\theta^3}$$

with $X \in [0, \theta]$. Compute the cumulative distribution function of the maximum $X_{(n)}$.

Per ottenerla occorre sviluppare la cumulata della funzione di partenza

$$F_X(x) = \int \frac{3x^2}{\theta^3} = \frac{3}{\theta^3} \int x^2 = \frac{3}{\theta} \frac{x^3}{3} = \frac{x^3}{\theta^3}$$

Da cui

$$F_{X_{(i)}}(x) = [F_X(x)]^n = \left(\frac{x}{\theta}\right)^{3n}$$

come confermato da taluni

Example 7.2.9 (Esame vecchio viroli). A random variable X has density function

$$f(x, \theta) = \frac{2x}{\theta^2}$$

with $X \in [0, \theta]$. Compute the probability distribution function of the maximum $X_{(n)}$

1. $F_n(x) = \frac{x^{2n}}{\theta^n}$
2. $F_n(x) = \frac{x^{n-1}}{\theta^n}$
3. $F_n(x) = \frac{x^{3n-1}}{\theta^{3n}}$
4. $F_n(x) = \frac{x^{2n}}{\theta^{2n}}$; taluni suggeriscono questa

Analogamente

$$F_X(x) = \int \frac{2x}{\theta^2} = \frac{2}{\theta^2} \int x = \frac{2}{\theta} \frac{x^2}{2} = \frac{x^2}{\theta^2}$$

da cui

$$F_{X_{(i)}}(x) = [F_X(x)]^n = \left(\frac{x}{\theta}\right)^{2n}$$

7.2.3 Generalized $X_{(i)}$

Important remark 48. If we write $X_{(i)} \sim F_{(i)}(x)$, with $i = 1, \dots, n$ we mean that $X_{(i)}$ is distributed following the i -th ordered statistic.

Proposition 7.2.5 (Distribution function of i -th ordered statistics). *We have*

$$F_{(i)}(x) = \mathbb{P}(X_{(i)} \leq x) = \sum_{j=i}^n \binom{n}{j} F_X(x)^j \cdot (1 - F_X(x))^{n-j} \quad (7.4)$$

Dimostrazione. To find the distribution function of $X_{(i)}$, it is convenient to think that a success occurs at trial i if $X_i \leq x$ (here is just comparing the unsorted sequence of realization with a threshold x of interest). One obtains

$$\begin{aligned}
 F_{(i)}(x) &= \mathbb{P}(X_{(i)} \leq x) \\
 &\stackrel{(1)}{=} \mathbb{P}(\text{at least } i \text{ successes/observation below } x) \\
 &= \sum_{j=i}^n \mathbb{P}(\text{exactly } j \text{ successes occur}) \\
 &\stackrel{(2)}{=} \sum_{j=i}^n \binom{n}{j} p^j (1-p)^{n-j} \\
 &= \sum_{j=1}^n \binom{n}{j} F(x)^j (1-F(x))^{n-j}
 \end{aligned}$$

where in

- (1) to have the i -th ordered observation under a certain threshold x means that we need *at least* i observations under that threshold (could be more as well, no problem, we're just focusing on the first i);
- in (2) where p is the probability of a success in a single trial, $\mathbb{P}(X \leq x)$ and it coincides with the distribution function F common to X_1, \dots, X_n , that is $p = F(x)$

□

Example 7.2.10. Imagine $n = 3$ with $x_{(1)} = 3$, $x_{(2)} = 5$, $x_{(3)} = 7$. We have that $\mathbb{P}(X_{(2)} \leq x)$ is the probability that 2 rvs are $\leq x$ OR the probability that 3 random variables are $\leq x$.

Example 7.2.11 (Maximum). As a sepecial case, for $i = n$, one obtains

$$\mathbb{P}(X_{(n)} \leq x) = \binom{n}{n} F(x)^n (1-F(x))^{n-n} = F(x)^n$$

The above result may be also obtained arguing as follows

$$\mathbb{P}(X_{(n)} \leq x) = \mathbb{P}(X_i \leq x, \forall i) = \prod_{i=1}^n \mathbb{P}(X_i \leq x) = \mathbb{P}(X_1 \leq x)^n = F(x)^n$$

Example 7.2.12 (Minimum). For $i = 1$

$$\begin{aligned}
 \mathbb{P}(X_{(1)} \leq x) &= \sum_{j=1}^n \binom{n}{j} F(x)^j (1-F(x))^{n-j} \\
 &= \left[\sum_{j=0}^n \binom{n}{j} F(x)^j (1-F(x))^{n-j} \right] - \binom{n}{0} F(x)^0 (1-F(x))^{n-0} \\
 &= (F(x) + 1 - F(x))^n - (1 - F(x))^n \\
 &= 1 - (1 - F(x))^n
 \end{aligned}$$

Once again this result can also be obtained as follows

$$\begin{aligned}\mathbb{P}(X_{(1)} \leq x) &= 1 - \mathbb{P}(X_{(1)} > x) = 1 - \mathbb{P}(X_i > x, \forall i) = 1 - \prod_{i=1}^n \mathbb{P}(X_i > x) \\ &= 1 - \mathbb{P}(X_1 > x)^n = 1 - [1 - F(x)]^n\end{aligned}$$

Remark 261. Next we want more: we want the distribution of the ordered vector $Y = \{X_{(1)}, \dots, X_{(n)}\}$.

To this end, it is convenient to make a further assumption: not only the element of X are iid, but their common distribution is absolutely continuous

Theorem 7.2.6. *If X_1, \dots, X_n are iid and their common distribution is absolutely continuous, then Y is still absolutely continuous and the joint density of Y is*

$$g(x_1, x_2, \dots, x_n) = \begin{cases} n! \prod_{i=1}^n f(x_i) & \text{if } x_1 < x_2 < \dots < x_n \\ 0 & \text{otherwise} \end{cases}$$

where f denotes the density common to X_1, \dots, X_n

Remark 262. So if we have a random vector composed by iid absolutely continuous random variables, the vector of order statistics is still absolutely continuous and the joint density is described above.

Intuitively the productory of f is due to the original vector components (iid), then we have $n!$ permutations to produce the same arrangement.

Example 7.2.13. For instance if $n = 2$ then $\mathbb{P}(X_1 = X_2) = 0$ since X_1, X_2 are absolutely continuous and $\mathbb{P}(X_{(1)} < X_{(2)}) = 1$.

The density g of $(X_{(1)}, X_{(2)})^\top$ is null on the part under main bisector, $(\{(x, y) \in \mathbb{R}^2 : x \geq y\})$, we have that the higher ordered element is y while lowest is x .

On the set $\{(x, y) \in \mathbb{R}^2 : y > x\}$ (above bisettrice) we have that the density is given by

$$g(x, y) = 2f(x)f(y)$$

Example 7.2.14. Let X_1 and X_2 be iid with $X_1 \sim \text{Unif}(0, 1)$. Then $Y = (X_{(1)}, X_{(2)})^\top$ is absolutely continuous with density

$$g(x, y) = \begin{cases} 2!f(x)f(y) & \text{if } x < y \\ 0 & \text{otherwise} \end{cases}$$

Since

$$f(x) = \begin{cases} 1 & \text{if } x \in (0, 1) \\ 0 & \text{otherwise} \end{cases}$$

one finally obtains

$$g(x, y) = \begin{cases} 2 \cdot 1 \cdot 1 = 2 & \text{if } 0 < x < y < 1 \\ 0 & \text{otherwise} \end{cases}$$

Intuition (area below have to sum always to 1 so):

- the density of $X = (X_1, X_2)^\top$ is 1 on the square between $(0, 0)$ and $(1, 1)$;
- the density of $X = (X_{(1)}, X_{(2)})^\top$ is 2 on half of the above square, that is on triangle between $(0, 0)$, $(1, 1)$ and $(1, 0)$

Proposition 7.2.7 (Density function for i -th order statistic).

NB: Da qui in poi della Viols direi

$$f_{(i)}(x) = \mathbb{P}(X_{(i)} = x) = \binom{n}{i} \cdot i \cdot F_X(x)^{i-1} \cdot f_X(x)(1 - F_X(x))^{n-i} \quad (7.5)$$

Important remark 49. Eg when $i = 1$ we obtain the formula for minimum

$$\begin{aligned} f_{(i)}(x) &= \binom{n}{1} 1 F_X(x)^0 \cdot f_X(x)(1 - F_X(x))^{n-1} \\ &= n f_X(x) \cdot [1 - F_X(x)]^{n-1} \end{aligned}$$

while for $i = n$ the maximum

$$f_{(n)}(x) = \binom{n}{n} n F_X(x)^{n-1} \cdot f_X(x)(1 - F_X(x))^0 = n [F_X(x)]^{n-1} f_X(x)$$

Example 7.2.15. Let $X_1, \dots, X_n \sim \text{Unif}(0, 1)$ be n iid uniforms, therefore having

$$f_X(x) = \begin{cases} 1 & \text{if } 0 \leq x < 1 \\ 0 & \text{elsewhere} \end{cases}, \quad F_X(x) = \begin{cases} 0 & \text{if } x < 0 \\ x & \text{if } 0 < x \leq 1 \\ 1 & \text{if } x > 1 \end{cases}$$

The k -th ordered statistic is distributed as a beta. Let's see it:

$$f_{(k)}(x) = k \binom{n}{k} x^{k-1} (1-x)^{n-k}$$

Now we have that

$$k \binom{n}{k} = \frac{n!}{(k-1)!(n-k)!} = \frac{\Gamma(n+1)}{\Gamma(k)\Gamma(n-k+1)} = \frac{1}{B(k, n-k+1)}$$

Therefore

$$f_{(k)}(x) = \frac{1}{B(k, n-k+1)} x^{k-1} (1-x)^{n-k}$$

or $X_{(k)} \sim \text{Beta}(k, n-k+1)$. As special cases

$$\begin{aligned} X_{(1)} &\sim \text{Beta}(1, n) \\ X_{(n)} &\sim \text{Beta}(n, 1) \end{aligned}$$

7.3 Inequalities

7.3.1 Tchebychev (Rigo)

Remark 263. One reason for Tchebychev inequality is so useful is that it holds for any rv X without any further assumption.

However, just for this reason it usually does not provide a precise estimate of $\mathbb{P}(|X| > c)$, just an upper margin.

Theorem 7.3.1. For any real random variable X , $\forall c > 0, \forall \alpha > 0$ (eg $\alpha = 1, 2, \dots$)

$$\mathbb{P}(|X| \geq c) \leq \frac{\mathbb{E}[|X|^\alpha]}{c^\alpha} \quad (7.6)$$

Dimostrazione. In general, given an event A in \mathcal{A} we have the indicator random variable

$$I_A = \begin{cases} 1 & \text{if } \omega \in A \\ 0 & \text{if } \omega \notin A \end{cases}$$

Then we have that

$$\mathbb{E}[I_A] = 0 \cdot \mathbb{P}(I_A = 0) + 1 \cdot \mathbb{P}(I_A = 1) = \mathbb{P}(I_A = 1) = \mathbb{P}(A)$$

To prove Tchebychev lets define

$$A = \{\omega : |X(\omega)| \geq c\} = \{|X| \geq c\}$$

then

$$\mathbb{E}[|X|^\alpha] \stackrel{(1)}{\geq} \mathbb{E}[I_A \cdot |X|^\alpha] \stackrel{(2)}{\geq} \mathbb{E}[I_A \cdot c^\alpha] = c^\alpha \mathbb{E}[I_A] = c^\alpha \mathbb{P}(A)$$

where:

- (1) because $|X|^\alpha \geq I_A \cdot |X|^\alpha$
- (2) we have $|X|^\alpha \geq c^\alpha$ since $|X| \geq c$ when we select with the indicator I_A (otherwise inside parenthesis is 0)

Therefore we conclude that

$$\mathbb{P}(A) = \mathbb{P}(|X| \geq c) \leq \frac{\mathbb{E}[|X|^\alpha]}{c^\alpha}$$

□

Remark 264. An important special case is when $X = Y - \mathbb{E}[Y]$ and $\alpha = 2$, in this case the inequality goes to

$$\mathbb{P}(|Y - \mathbb{E}[Y]| \geq c) \leq \frac{\mathbb{E}[(Y - \mathbb{E}[Y])^2]}{c^2} = \frac{\text{Var}[Y]}{c^2}$$

But to apply Tchebychev in this form we need to know that the variance exists. Some books call this special case Chebichev inequality and call the general case Markov inequality.

7.3.2 Jensen (Rigo)

Definition 7.3.1 (Convex function (conca tipo $y = x^2$)). $f : I \rightarrow \mathbb{R}$ is a convex function if

$$\begin{cases} f[\alpha x + (1 - \alpha)y] \leq \alpha f(x) + (1 - \alpha)f(y) \\ \forall \alpha \in [0, 1], x, y \in I \end{cases}$$

where:

- $f[\alpha x + (1 - \alpha)y]$ can be seen as the value given by the function at the mean point between x and y
- $\alpha f(x) + (1 - \alpha)f(y)$ the mean of the value assumed by the function in the two extremes

Important remark 50. If f is twice differentiable:

$$f \text{ is convex} \iff f'' \geq 0$$

Example 7.3.1. For instance $f(x) = x^2$, $f(x) = e^x$, $f(x) = |x|$ are convex.

Definition 7.3.2 (Strictly convex function). Same definition as above but instead of \leq we have $<$: $f : I \rightarrow \mathbb{R}$ is strictly convex

$$\begin{cases} f[\alpha x + (1 - \alpha)y] < \alpha f(x) + (1 - \alpha)f(y) \\ \forall \alpha \in [0, 1], x, y \in I \end{cases}$$

Important remark 51. If f is twice differentiable:

$$f \text{ is strictly convex} \iff f'' > 0$$

Example 7.3.2. For instance $f(x) = x^2$, $f(x) = e^x$ are strictly convex. Similarly if $I = (0, \infty)$, $f(x) = \frac{1}{x}$ is strictly convex. In fact $f''(x) = 2x^{-3} > 0$, $\forall x > 0$

Proposition 7.3.2 (Jensen inequality). Let X be a real random variable and $f : I \rightarrow \mathbb{R}$ a function defined on interval I . If

1. f is a convex function
2. $\mathbb{P}(X \in I) = 1$
3. $\mathbb{E}[|X|] < +\infty$, $\mathbb{E}[|f(X)|] < +\infty$

Then:

$$\mathbb{E}[f(X)] \geq f(\mathbb{E}[X])$$

Moreover, if f is strictly convex and X is not degenerate, then

$$\mathbb{E}[f(X)] > f(\mathbb{E}[X])$$

NB: quest'anno non fatte?

Example 7.3.3. Let's see some application of Jensen inequality.

- $f(x) = x^2$ is strictly convex (second derivative = 2 > 0). If we apply Jensen we find out that

$$\mathbb{E}[X^2] > [\mathbb{E}[X]]^2 \quad (7.7)$$

This was already known since variance (for non degenerate variables as per the theorem) is ≥ 0 (by computational formula of variance).

- absolute value $f(x) = |x|$ (second derivative = 0); applying Jensen we discover something new

$$\mathbb{E}[|X|] \geq |\mathbb{E}[X]| \quad (7.8)$$

- $f(x) = x^{b/a}$ for any $x \geq 0$ with $(0 < a < b)$. Applying Jensen

$$\mathbb{E} \left[|X|^b \right] = \mathbb{E} \left[(|X|^a)^{\frac{b}{a}} \right] \geq [\mathbb{E} [|X|^a]]^{\frac{b}{a}} \quad (7.9)$$

thus Jensen implies that

$$\mathbb{E} \left[(|X|^a)^{\frac{1}{a}} \right] \leq \mathbb{E} \left[(|X|^b)^{\frac{1}{b}} \right]$$

Remark 265. Now we use Jensen to prove that the rv is degenerate iff its variance is 0.

Proposition 7.3.3.

$$X \sim \delta_a \iff \text{Var} [X] = 0$$

Dimostrazione. Respectively:

- supposing $X = a$ almost surely ($\mathbb{P}(X = a) = 1$), then $\mathbb{E}[X] = a$ and also $\mathbb{E}[X^2] = a^2$, thus

$$\text{Var} [X] = \mathbb{E} [X^2] - (\mathbb{E} [X])^2 = a^2 - a^2 = 0$$

- otherwise suppose $\text{Var} [X] = 0$: we prove that by contradiction. By applying Jensen inequality with $f(x) = x^2$, strictly convex, if X is *non degenerate* we get:

$$\mathbb{E} [X^2] = \mathbb{E} [f(X)] > f(\mathbb{E} [X]) = (\mathbb{E} [X])^2$$

this happens if and only if $\text{Var} [X] = \mathbb{E} [X^2] - (\mathbb{E} [X])^2 > 0$: but we assumed $\text{Var} [X] = 0$ so we found a contradiction (thus X must be degenerate). \square

7.3.3 Markov (Viroli)

Theorem 7.3.4. Given $X \in \mathbb{R}^+$, $D_X = \mathbb{R}^+$, $\lambda > 0$

$$\mathbb{P}(X \geq \lambda \cdot \mathbb{E} [X]) \leq \frac{1}{\lambda} \quad (7.10)$$

Dimostrazione. Let

$$\mathbb{E} [X] = m = \int_{D_X} x \cdot f(x) \, dx = \int_0^{+\infty} x \cdot f(x) \, dx$$

Now

$$m \geq \int_{\lambda m}^{+\infty} x \cdot f(x) \, dx \geq \int_{\lambda m}^{+\infty} x \cdot m \cdot f(x) \, dx = \lambda m \underbrace{\int_{\lambda m}^{+\infty} f(x) \, dx}_{\mathbb{P}(X \geq \lambda \cdot m)}$$

therefore

$$m \geq \lambda m \mathbb{P}(X \geq \lambda \cdot m) \iff \frac{1}{\lambda} \geq \mathbb{P}(X \geq \lambda \cdot m)$$

\square

Example 7.3.4 (Esame vecchio viroli). Let $\{X_n\}$ be a sequence of independent exponential random variables with parameter $\lambda_n = \frac{n}{2}$. Find the value of n such that $\mathbb{P}(X_n > 0.25) \leq 0.8$.

According to markov inequality

$$\mathbb{P}(X \geq c \mathbb{E}[X]) \leq \frac{1}{c}$$

We have that

$$\mathbb{E}[X_n] = \frac{1}{\lambda_n} = \frac{2}{n}$$

So

$$\mathbb{P}\left(X_n \geq c \cdot \frac{2}{n}\right) \leq \frac{1}{c}$$

TODO: boh qui non mi è chiarissimo

Now if $\frac{1}{c} = 0.8$ then $c = 1.25$ and we have:

$$\mathbb{P}\left(X_n \geq 1.25 \frac{2}{n}\right) \leq 0.8$$

$$\mathbb{P}\left(X_n \geq \frac{2.5}{n}\right) \leq 0.8$$

$$\frac{2.5}{n} = 0.25$$

$$n = 10$$

Risposta $n = 10$

7.3.4 Tchebychev (Viroli)

Important remark 52. We have two equivalent formulations

Theorem 7.3.5 (Tchebychev inequality). *Respectively*

$$\mathbb{P}(|X - \mathbb{E}[X]| \geq \lambda \cdot \sigma_X) \leq \frac{1}{\lambda^2} \quad (7.11)$$

$$\mathbb{P}(|X - \mathbb{E}[X]| < \lambda \cdot \sigma_X) \geq 1 - \frac{1}{\lambda^2} \quad (7.12)$$

where σ_X is the standard deviation of X

Dimostrazione. We do by applying Markov inequality to $Y = (X - \mathbb{E}[X])^2$. We have that $\mathbb{E}[Y] = \sigma_X^2$ (by definition of variance), so by Markov

$$\mathbb{P}(Y \geq \lambda \mathbb{E}[Y]) \leq \frac{1}{\lambda}$$

$$\mathbb{P}\left((X - \mathbb{E}[X])^2 \geq \lambda \sigma_X^2\right) \leq \frac{1}{\lambda}$$

$$\mathbb{P}\left(|X - \mathbb{E}[X]| \geq \sqrt{\lambda} \sigma_X\right) \leq \frac{1}{\lambda}$$

Then by setting $\lambda^* = \sqrt{\lambda}$ we conclude

$$\mathbb{P}(|X - \mathbb{E}[X]| \geq \lambda^* \sigma_x) \leq \frac{1}{(\lambda^*)^2}$$

□

Example 7.3.5 (esame viroli). Let X_n be a sequence of independent poisson random variable with parameter 9 and $\bar{x}_n = \sum_{i=1}^n X_i/n$ is the partial mean. By the chebychev inequality find the value of n such that

$$\mathbb{P}(|\bar{x} - 9| < 15) \geq 0.99$$

- n = 36
- n = 10
- n = 4 taluni suggeriscono questa, confermata sotto
- n = 40

Qui effettivamente si ha che 9 è il valore atteso della somma di queste poissoniane poiché

$$\mathbb{E} \left[\frac{\sum_{i=1}^n X_i}{n} \right] = \frac{1}{n} \mathbb{E} \left[\sum_{i=1}^n X_i \right] = \frac{n \mathbb{E}[X_i]}{n} = \mathbb{E}[X_i] = 9$$

Il setup è dunque giusto e data la richiesta dobbiamo applicare la disuguaglianza nella seconda versione; abbiamo che

$$1 - \frac{1}{\lambda^2} = 0.99 \iff \lambda = 10$$

Dunque si ha che

$$10 \sqrt{\text{Var} \left[\sum_{i=1}^n X_i/n \right]} = 15$$

Ora per ricavare n (ricordando che $\text{Var}[X_i] = \mathbb{E}[X_i] = 9$)

$$\begin{aligned} \text{Var} \left[\frac{\sum_{i=1}^n X_i}{n} \right] &= \frac{\text{Var} [\sum_{i=1}^n X_i]}{n^2} = \frac{n \text{Var} [X_i]}{n^2} = \frac{\text{var} X_i}{n} = \frac{9}{n} \\ \sqrt{\text{Var} \left[\frac{\sum_{i=1}^n X_i}{n} \right]} &= \frac{3}{\sqrt{n}} \end{aligned}$$

Dunque

$$10 \frac{3}{\sqrt{n}} = 15 \iff \sqrt{n} = 2 \iff n = 4$$

7.4 Characteristic and moment generating function

7.4.1 Characteristic function

Definition 7.4.1 (Characteristic function). If $\mathbf{X} = \begin{bmatrix} X_1 \\ \dots \\ X_n \end{bmatrix}$ is a n -variate random vector, the characteristic function $\varphi : \mathbb{R}^n \rightarrow \mathbb{C}$ of \mathbf{X} is

$$\begin{aligned} \varphi_{\mathbf{X}}(\mathbf{t}) &= \mathbb{E} \left[e^{i\mathbf{t}^\top \mathbf{X}} \right] = \mathbb{E} \left[e^{i \sum_i t_i X_i} \right] \quad \forall \mathbf{t} = \begin{bmatrix} t_1 \\ \dots \\ t_n \end{bmatrix} \in \mathbb{R}^n \\ &= \mathbb{E} \left[\cos \mathbf{t}^\top \mathbf{X} + i \sin \mathbf{t}^\top \mathbf{X} \right], \end{aligned}$$

where

- $i \in \mathbb{I} : i^2 = -1$
- being both \mathbf{t} and \mathbf{X} vectors the $\mathbf{t}^\top \mathbf{X} = \sum_{i=1}^n t_i X_i$ is a scalar and Euler's formula ($e^{ix} = \cos(x) + i \sin(x)$) applies.

Remark 266. However, from now on we assume single variable (because it's more convenient) not n -variate random vector. The definition above simplifies to the following

Definition 7.4.2 (Characteristic function). Let X be a random variable, the characteristic function $\varphi_X(t) : \mathbb{R} \rightarrow \mathbb{C}$, existing $\forall t \in \mathbb{R}$ is defined as

$$\begin{aligned} \varphi_X(t) &= \mathbb{E} [e^{itX}] = \int_{-\infty}^{+\infty} e^{itx} f(x) \, dx \\ &= \int_{-\infty}^{+\infty} \cos(tx) f(x) \, dx + i \int_{-\infty}^{+\infty} \sin(tx) f(x) \, dx \end{aligned}$$

NB: esempio viroliano credo

Example 7.4.1 (Characteristic function of a binomial). Let $X \sim \text{Bin}(n, p)$, $D_x = \{0, 1, \dots, n\}$, the characteristic function is

$$\begin{aligned} \varphi_X(t) &= \mathbb{E} [e^{itX}] = \sum_{x=0}^n e^{itx} \cdot \binom{n}{x} p^x (1-p)^{n-x} = \sum_{x=0}^n \binom{n}{x} (\underbrace{pe^{it}}_a)^x (\underbrace{1-p}_b)^{n-x} \\ &\stackrel{(1)}{=} (1-p + pe^{it})^n \end{aligned}$$

where in (1) we applied binomial formula $(a+b)^n = \sum_{x=0}^n \binom{n}{x} a^x b^{n-x}$

Example 7.4.2 (Characteristic function of $X \sim N(0, 1)$).

$$\begin{aligned} \varphi_X(t) &= \mathbb{E} [e^{itX}] = \int_{-\infty}^{+\infty} e^{itx} f(x) \, dX = \int_{-\infty}^{+\infty} e^{itx} \frac{e^{-\frac{x^2}{2}}}{\sqrt{2\pi}} \, dX \\ &= e^{-\frac{t^2}{2}} \end{aligned}$$

l'ultimo passaggio l'ha giusto detto a-la "trust me" presumo (integrale di funzione complessa)

Important remark 53 (Usefulness). Despite being complicated/complex functions, they are useful for several reasons (both theoretical and practical):

1. they *determine the distribution* of the random variable: this is the reason this stuff is so important to statistic (**important for Rigo**);
2. they provide a *link with the moment* of order k of the variable via *differentiation* (with respect to t evaluated at $t = 0$);
3. they provide a *link with the distribution function* via the *inversion formula*.

Theorem 7.4.1 (Link with distribution). *Supposing we have two random vectors \mathbf{X}, \mathbf{Y} , these have the same distribution iff they share the characteristic function:*

$$X \sim Y \iff \varphi_{\mathbf{X}}(\mathbf{t}) = \varphi_{\mathbf{Y}}(\mathbf{t}), \quad \forall \mathbf{t} \in \mathbb{R}^n \quad (7.13)$$

Proposition 7.4.2 (Link with the moments). *We have:*

NB: credo viroliano

$$\left[\frac{\partial^k}{\partial t^k} \varphi_X(t) \right]_{t=0} = i^k \mathbb{E}[X^k]$$

and therefore

$$\mathbb{E}[X^k] = \frac{\left[\frac{\partial^k}{\partial t^k} \varphi_X(t) \right]_{t=0}}{i^k}$$

Dimostrazione. We have

$$\frac{\partial^k}{\partial t^k} \varphi_X(t) = \frac{\partial^k}{\partial t^k} \mathbb{E}[e^{itX}] = \mathbb{E}\left[\frac{\partial^k}{\partial t^k} e^{itX}\right] = \mathbb{E}[i^k X^k e^{itX}]_{t=0} \stackrel{(1)}{=} i^k \mathbb{E}[X^k]$$

where in (1) we evaluated for $t = 0$. □

Proposition 7.4.3 (Link with density (inversion formula)).

NB: credo viroliano

$$f_X(x) = \frac{1}{2\pi} \int_{-\infty}^{+\infty} e^{-itx} \varphi_X(t) dt$$

Dimostrazione. Virols skips it. □

Important remark 54 (Important properties (Rigo)). We have:

1. **link with random variables independence:** if $X \perp\!\!\!\perp Y$, the characteristic function of the sum is equal to the product of the single characteristic functions

$$\varphi_{X+Y}(t) = \varphi_X(t) \cdot \varphi_Y(t), \quad \forall t \in \mathbb{R} \quad (7.14)$$

This because

$$\begin{aligned} \varphi_{X+Y}(t) &= \mathbb{E}[e^{it(X+Y)}] = \mathbb{E}[e^{itX} e^{itY}] \stackrel{(1)}{=} \mathbb{E}[e^{itX}] \mathbb{E}[e^{itY}] \\ &= \varphi_X(t) \cdot \varphi_Y(t), \quad \forall t \in \mathbb{R} \end{aligned}$$

where in (1), since $X \perp\!\!\!\perp Y$, any combination is independent as well, and so we apply the expected value property for product of independent variables $Z \perp\!\!\!\perp W \implies \mathbb{E}[ZW] = \mathbb{E}[Z] \mathbb{E}[W]$.

Because of 7.14, characteristic function becomes *very handy* when working with sums of independent rvs.

2. **characteristic function and moments:** if the random variable has the moment of order j , that is $\mathbb{E}[|X|^j] < +\infty$, then:

- (a) the characteristic function $\varphi_X(t) \in C^j$, where C^j is the collection of functions which have the derivative of order j and such derivative is continuous;
- (b) the derivative of order $r \leq j$ is:

$$\begin{aligned}\varphi_X(t)^{(r)} &= \frac{\partial^r}{\partial t^r} \varphi_X(t) = \frac{\partial^r}{\partial t^r} \mathbb{E}[e^{itX}] = \mathbb{E}\left[\frac{\partial^r}{\partial t^r} e^{itX}\right] \\ &= \mathbb{E}[(iX)^r e^{itX}]\end{aligned}$$

This latter means that in each derivative up to order j we can interchange the operator of derivative and the operator of expectation. For instance, suppose we want to calculate the first derivative; by setting $r = 1$

$$\varphi_X(t)' = \mathbb{E}[iX e^{itX}]$$

Actually the **interesting fact** is that if we evaluate the r -th derivative for $t = 0$ we have a direct interpretation/link with the r -th moment

$$\varphi_X(0)^{(r)} = \mathbb{E}[(iX)^r e^{i0X}] = i^r \mathbb{E}[X^r]$$

Before we said that $\mathbb{E}[|X|^j] < +\infty \implies \varphi_X(t) \in C^j$. The converse implication does not generally holds (only if j is even). We have

$$\begin{cases} \text{If } j \text{ is odd, } \mathbb{E}[|X|^j] < +\infty \implies \varphi_X(t) \in C^j \\ \text{If } j \text{ is even, } \mathbb{E}[|X|^j] < +\infty \iff \varphi_X(t) \in C^j \end{cases}$$

As counterexample of the first missing counterimplication, we will see a case where (with $j = 1$ odd) where $\varphi_X(t) \in C^1$ ma $\mathbb{E}[|X|] = +\infty$;

3. **inversion theorem** gives a closed formula for determining the distribution function starting from characteristic function.
Let F be the distribution function of X If/forall $a < b$ such that $\mathbb{P}(X = a) = \mathbb{P}(X = b) = 0$ then:

$$F(b) - F(a) = \mathbb{P}(a < X \leq b) = \frac{1}{2\pi i} \lim_{c \rightarrow +\infty} \int_{-c}^c \frac{e^{-ita} - e^{-itb}}{t} \varphi_X(t) dt$$

4. **continuity theorem:** if X_n and X are real rvs then

$$X_n \xrightarrow{d} X \iff \varphi_X(t) = \lim_{n \rightarrow +\infty} \varphi_{X_n}(t), \quad \forall t \in \mathbb{R}$$

We'll see later \xrightarrow{d} means convergence in distribution (an important type of convergence): point is that any time we want to prove convergence in distribution we can, if convenient, prove the limit of characteristic function.

NB: For inversion thm, the important fact to recall for the exam is that characteristic function can be inverted: if you know the characteristic function, there exists a formula that allows to write down the distribution function (no need to memorize it for the exam).

Proposition 7.4.4 (Altre proprietà utili trovate su wikipedia). *Si ha:*

1. If X_1, \dots, X_n are independent random variables, and $a_1, \dots, a_n \in \mathbb{R}$, the characteristic function of the linear combination

$$\varphi_{a_1 X_1 + \dots + a_n X_n}(t) = \varphi_{X_1}(a_1 t) \cdot \dots \varphi_{X_n}(a_n t).$$

2. Let the random variable $Y = aX + b$ be the linear transformation of a random variable X . The characteristic function of Y is $\varphi_Y(t) = e^{itb} \varphi_X(at)$.

3. For random vectors \mathbf{X} and $\mathbf{Y} = \mathbf{A}\mathbf{X} + \mathbf{B}$ (where \mathbf{A} is a constant matrix and \mathbf{B} a constant vector), we have

$$\varphi_{\mathbf{Y}}(\mathbf{t}) = e^{i\mathbf{t}^\top \mathbf{B}} \varphi_{\mathbf{X}}(\mathbf{A}^\top \mathbf{t})$$

Example 7.4.3 (Characteristic function of $N(\mu, \sigma^2)$). As example of the second Rigo considered that if $X \sim \mu + \sigma Z$ where $Z \sim N(0, 1)$, then $X \sim N(\mu, \sigma^2)$. For its characteristic function we have

$$\begin{aligned} \varphi_X(t) &= \mathbb{E}[e^{itX}] = \mathbb{E}[e^{it(\mu + \sigma Z)}] = \mathbb{E}[e^{it\mu} e^{it\sigma Z}] = e^{it\mu} \mathbb{E}[e^{i(t\sigma)Z}] = e^{it\mu} \varphi_Z(t\sigma) \\ &= e^{it\mu} e^{-\frac{t^2 \sigma^2}{2}} \end{aligned}$$

where in the last passage we used results from example 7.4.2.

Example 7.4.4 (Example by Rigo, fatto solo il mio anno? forse da postporre alle convergenze: weak law of large number). In this example we show that if X_n is iid and the characteristic function has the first derivative at 0, $\exists \varphi_X(0)'$, then the sample mean converges (in distribution and probability) to a constant/degenerate rv.

Let $\{X_n\}_{n \in \mathbb{N}}$ be a sequence of iid rvs; we define the sample mean as:

$$\bar{X}_n = \frac{\sum_{i=1}^n X_i}{n}$$

The characteristic function of the sample mean is

$$\varphi_{\bar{X}_n}(t) = \varphi_{\sum_i X_i}\left(\frac{t}{n}\right) \stackrel{(\text{ii})}{=} \prod_{i=1}^n \varphi_{X_i}\left(\frac{t}{n}\right) \stackrel{(\text{id})}{=} \left[\varphi_{X_i}\left(\frac{t}{n}\right) \right]^n$$

Suppose now that the first derivative of the characteristic function of X_i exists in 0, that is $\exists \varphi_{X_i}(0)'$; then by Taylor expansion formula

$$\varphi_{\bar{X}_n}(t) = \left[\varphi_{X_i}\left(\frac{t}{n}\right) \right]^n = \left[\varphi_{X_i}(0) + \frac{t}{n} \varphi_{X_i}(0)' + o\left(\frac{t}{n}\right) \right]^n = \left[1 + \frac{t \varphi_{X_i}(0)' + no\left(\frac{t}{n}\right)}{n} \right]^n$$

where $o\left(\frac{t}{n}\right)$ is the Peano rest. In general $g = o(f)$ if $\lim_{x \rightarrow x_0} \frac{g(x)}{f(x)} = 0$.

Now, what is the limit of the formula above for $n \rightarrow +\infty$? Using the fact that

$$\text{if } a_n \rightarrow a \implies \left(1 + \frac{a_n}{n}\right)^n \rightarrow e^a$$

we have (with $a_n = t\varphi_{X_i}(0)' + no(\frac{t}{n})$ and noted that $a_n \rightarrow t\varphi_{X_i}(0)' + 0$)

$$\varphi_{\bar{X}_n}(t) \rightarrow e^{t\varphi_{X_i}(0)'}$$

Now it can be shown (we won't) that the first derivative in 0 is

$$\varphi_{X_i}(0)' = i\alpha, \quad \alpha \in \mathbb{R}$$

and thus we our characteristic function converges to

$$\varphi_{\bar{X}_n}(t) \rightarrow e^{it\alpha}, \forall t \in \mathbb{R}$$

Is $e^{it\alpha}$ a characteristic function? Yes the δ_α has this characteristic function since if $X \sim \delta_\alpha$

$$\varphi_X(t) = \mathbb{E}[e^{itX}] = \mathbb{E}[e^{it\alpha}] = e^{it\alpha}$$

Hence $\bar{X}_n \xrightarrow{d} \alpha$, by continuity theorem, and since the limit is a degenerate rv, we have not only convergence in distribution but also convergence in probability $\bar{X}_n \xrightarrow{p} \alpha$.

Important remark 55. The above should be *weak law of large number* (convergence not a.s. but only in probability, check with Viols).

Furthermore, if the sequence is not only iid, but also the mean exists, $\mathbb{E}[|X_i|] < +\infty$, then $\bar{X}_n \xrightarrow{a.s.} \mathbb{E}[X_i]$ then the sample mean converges almost surely to the mean (this is the *strong law of large number*).

But as noted above, it may be that $\exists \varphi_{X_i}(0)'$ even if $\mathbb{E}[|X_i|] = +\infty$.

7.4.2 Moment generating function

Definition 7.4.3 (Moment generating function (mgf)). It's obtained from the characteristic function by evaluating it at $-it$, $\varphi_X(-it)$, so that there are no complex number:

$$\varphi_X(-it) = \mathbb{E}[e^{-iitX}] = \mathbb{E}[e^{tX}] = M_X(t), \quad \forall t \in \mathbb{R}$$

so

$$M_X(t) = \mathbb{E}[e^{tX}], \quad \forall t \in \mathbb{R}$$

Poiché $e^{tX} > 0$ si ha che $M_X(t) > 0$

Important remark 56. It's simpler than characteristic function (no i here) but has its drawbacks. Differently from characteristic function (always exists, we have inversion thm):

- MGF always exists for $t = 0$

$$M_X(0) = 1 < +\infty$$

MGF may *fail to exist* for $t \neq 0$. If for some $t \neq 0$ (or for $\forall t \in \mathbb{R}$) it is $M_X(t) = +\infty$, in those case MGF is not useful/does not exist;

- we don't have an inversion theorem, so it's useful only for the moments (not distribution).

Important remark 57 (Random variable with MGF). If we know that the moment generating function is finite in a neighborhood of $t = 0$, that is

$$M_X(t) < +\infty, \quad \forall t \in (-\varepsilon, \varepsilon)$$

we say that X has *moment generating function*. In that case it may be convenient to use it instead of the characteristic function, since it can be proven that:

- the random variable has *moments of every order*: $\mathbb{E}[|X|^n] < +\infty, \forall n$
- the probability distribution of X is *determined* by its moments that is $X \sim Y$ for any rv Y such that $\mathbb{E}[X^n] = \mathbb{E}[Y^n], \forall n$.
The sequence of moments $\mathbb{E}[X^n]$, with $n = 1, 2, \dots$, determines the distribution, in the sense that if X and Y does *not* have the same distribution then *either* one of them have some moments not finite or moments both are finite but different for some n :

$$X \approx Y \implies \begin{cases} \mathbb{E}[|X|^n] = +\infty, \text{ for some } n, \text{ OR} \\ \mathbb{E}[X^n] \neq \mathbb{E}[Y^n], \text{ for some } n \end{cases}$$

Important remark 58. Consider two random variables X, Y , with moments of every order (mean, variance, third moment etc) *existing* and *coinciding*:

$$\begin{cases} \mathbb{E}[|X|^n] < +\infty, \mathbb{E}[|Y|^n] < +\infty \\ \mathbb{E}[X^n] = \mathbb{E}[Y^n] \end{cases} \quad \forall n$$

Can we conclude that the two random variables have the same distribution? **No** we cannot conclude that (this is contrary to intuition).

Eg if X is lognormal, one can build a suitable rv Y such that X and Y have the same moments of every order and yet $X \approx Y$.

However this annoying fact doesn't occur if one between X and Y has moment generating function. In that case we can say they have the same distribution.

$$\begin{cases} \mathbb{E}[|X|^n] < +\infty, \mathbb{E}[|Y|^n] < +\infty \\ \mathbb{E}[X^n] = \mathbb{E}[Y^n] \\ X \text{ or } Y \text{ has finite moment generating function} \end{cases} \quad \forall n \implies X \sim Y$$

Example 7.4.5. An important special case where $M_X(t) < +\infty, \forall t \in \mathbb{R}$ is

$$|X| \leq c \text{ a.s. for some constant } c$$

In fact

$$M_X(t) = \mathbb{E}[e^{tX}] \leq \mathbb{E}[e^{|tX|}] \leq \mathbb{E}[e^{|t|c}] = e^{|t|c} < +\infty, \quad \forall t \in \mathbb{R}$$

NB: da qui in poi roba della viroli direi

Proposition 7.4.5 (Properties).

$$\left[\frac{\partial^k}{\partial t^k} M_X(t) \right]_{t=0} = \mathbb{E} [X^k] \quad (7.15)$$

$$M_X(0) = \mathbb{E} [e^{0X}] = \mathbb{E} [1] = 1 \quad (7.16)$$

$$M_X(t) = M_Y(t), \forall t \iff F_X(x) = F_Y(y) \quad (\text{uniqueness}) \quad (7.17)$$

$$M_{aX+b}(t) = e^{tb} M_X(at), \quad a, b \in \mathbb{R} \quad (7.18)$$

$$X \perp\!\!\!\perp Y \implies M_{X+Y}(t) = M_X(t) \cdot M_Y(t) \quad (7.19)$$

TODO: l'implicazione per l'indipendenza è anche coimplicazione?

Dimostrazione. For 7.18

$$M_{aX+b}(t) = \mathbb{E} [e^{t(aX+b)}] = \mathbb{E} \left[e^{taX} \cdot \underbrace{e^{tb}}_{\text{constant}} \right] = e^{tb} \cdot \mathbb{E} [e^{taX}] = e^{tb} M_X(at)$$

For 7.19

$$M_{X+Y}(t) = \mathbb{E} [e^{t(X+Y)}] = \mathbb{E} [e^{tX} e^{tY}]$$

Now note that:

- first

TODO: questo andrebbe portato più in vista nella sezione indipendenza o prop v. atteso

$$\begin{aligned} \mathbb{E} [g(X)h(Y)] &= \int_{D_x} \int_{D_y} g(x)h(y)f(x,y) \, dx \, dy \stackrel{(1)}{=} \int_{D_x} \int_{D_y} g(x)h(y)f_X(x)f_Y(y) \, dx \, dy \\ &= \int_{D_x} g(x)f_X(x) \, dx \int_{D_y} h(y)f_Y(y) \, dy = \mathbb{E} [g(X)] \mathbb{E} [h(Y)] \end{aligned}$$

where (1) due to be $X \perp\!\!\!\perp Y$.

- furthermore

$$\begin{aligned} \mathbb{E} [g(X) + h(Y)] &= \int_{D_x} \int_{D_y} (g(x) + h(y))f(x,y) \, dx \, dy \\ &= \int_{D_x} \int_{D_y} g(x)f(x,y) \, dx \, dy + \int_{D_x} \int_{D_y} h(y)f(x,y) \, dx \, dy \\ &= \int_{D_x} g(x) \underbrace{\int_{D_y} f(x,y) \, dy}_{f(x)} \, dx + \int_{D_x} \int_{D_y} h(y)f(x,y) \, dx \, dy \\ &= \int_{D_x} g(x)f(x) \, dx + \int_{D_y} h(y)f(y) \, dy = \mathbb{E} [g(X)] + \mathbb{E} [h(Y)] \end{aligned}$$

Therefore coming back to our focus, under independence and using the first one

$$M_{X+Y}(t) = \mathbb{E} [e^{tX} e^{tY}] \stackrel{(1)}{=} \mathbb{E} [e^{tX}] \mathbb{E} [e^{tY}] = M_X(t) M_Y(t)$$

in (1) because of $\perp\!\!\!\perp$

□

Example 7.4.6 (Esempio rigo triennale). Se $X \sim N(\mu, \sigma^2)$ allora X possiede la MGF che è (non dimostrato)

$$M_X(t) = \mathbb{E}[e^{tX}] = e^{t\mu + \frac{\sigma^2 t^2}{2}}, \quad \forall t \in \mathbb{R}$$

Inoltre si ha che la derivata prima

$$M_X(t)' = e^{t\mu + \frac{\sigma^2 t^2}{2}}(\mu + \sigma^2 t)$$

e valutandola per $t = 0$

$$M_X(0)' = e^0(\mu + \sigma^2 \cdot 0) = \mu \implies \mathbb{E}[X^1] = \mathbb{E}[X] = \mu$$

Proseguendo con la derivata seconda

$$M_X(t)'' = e^{t\mu + \frac{\sigma^2 t^2}{2}}(\mu + \sigma^2 t)^2 + e^{t\mu + \frac{\sigma^2 t^2}{2}}\sigma^2$$

e valutandola per $t = 0$ si ha.

$$M_X(0)'' = e^0(\mu + \sigma^2 \cdot 0)^2 + e^0\sigma^2 = \mu^2 + \sigma^2 \implies \mathbb{E}[X^2] = \mu^2 + \sigma^2$$

Da cui possiamo ricavare la varianza considerando

$$\text{Var}[X] = \mathbb{E}[X^2] - (\mathbb{E}[X])^2 = \mu^2 + \sigma^2 - \mu^2 = \sigma^2$$

Example 7.4.7 (Mgf of bernoulli and binomial). If $X \sim \text{Bern}(p)$, $p(x) = p^x(1-p)^{1-x}$, $D_x = \{0, 1\}$. Its mgf is:

$$M_X(t) = \mathbb{E}[e^{tX}] = e^{t \cdot 0} \cdot (1-p)p^0 + e^{t \cdot 1}p^1(1-p)^0 = 1 - p + pe^t$$

Being the binomial $Y = X_1 + \dots + X_n$ with X_i iid, by properties of mgfs, the mgf of a binomial is

$$M_Y(t) = \prod_{i=1}^n (1 - p + pe^t) = (1 - p + pe^t)^n$$

Example 7.4.8 (Mgf of poisson). Let $X \sim \text{Pois}(\lambda)$, let's determine $M_X(t)$

$$\begin{aligned} M_X(t) &= \mathbb{E}[e^{tX}] = \sum_{x=0}^{\infty} e^{tx} \frac{1}{x!} e^{-\lambda} \lambda^x = \sum_{x=0}^{\infty} (e^t \lambda)^x \frac{1}{x!} e^{-\lambda} \\ &\stackrel{(1)}{=} e^{-\lambda} \cdot e^{\lambda e^t} = e^{-\lambda(1-e^t)} = e^{\lambda(e^t-1)} \end{aligned}$$

where in (1) we used $\sum_{x=0}^{\infty} \frac{c^x}{x!} = e^c$.

Example 7.4.9 (Esercizio richiesto Viroli). By using 7.19 find $M_Y(t)$, with $Y = \sum_{i=1}^n X_i$, $X_i \sim \text{Pois}(\lambda_i)$, and $X_i \perp\!\!\!\perp X_j$.

La mgf di una poisson con parametro λ è $M_X(t) = e^{\lambda(e^t-1)}$, da cui per l'indipendenza possiamo applica la produttoria

$$M_Y(t) = \prod_{i=1}^n e^{\lambda_i(e^t-1)} = e^{\sum_{i=1}^n \lambda_i(e^t-1)} = e^{(e^t-1) \cdot \sum_{i=1}^n \lambda_i}$$

che è la mgf di una poisson con parametro lambda la somma delle lambda componenti (come atteso).

Therefore $\implies Y \sim \text{Pois}(\sum_{i=1}^n \lambda_i)$.

Example 7.4.10 (Esame vecchio viroli). Let X be a bernoulli rv with parameter $\frac{1}{2}$. Find the moment generating functions of $Y = \frac{1}{2} + \frac{X}{2}$. We have that for the bernoulli

$$M_X(t) = 1 - p + pe^t$$

and consider

$$M_{aX+b}(t) = e^{bt}M_X(at)$$

Now here we have $Y = \frac{1}{2} + \frac{X}{2}$ so $a = b = \frac{1}{2}$, therefore:

$$M_Y(t) = e^{t/2}M_X\left(\frac{t}{2}\right) = e^{\frac{t}{2}}(1 - p + pe^{t/2})$$

Finally, if $p = \frac{1}{2}$

$$M_Y(t) = e^{t/2}\left(\frac{1}{2} + \frac{e^{t/2}}{2}\right) = \frac{1}{2}(e^{t/2} + e^t)$$

Therefore we have that $M_Y(t) = \frac{1}{2}(e^t + e^{\frac{t}{2}})$

Example 7.4.11 (Esame vecchio viroli). Let X_1 and X_2 be two independent Bernoulli rv with parameters $1/2$. find the moment generating function of $Z = X_1 - X_2$.

If $X \sim \text{Bern}(p)$, its $M_X(t) = (1 - p + pe^t)$. Here for the difference of two bernoulli we apply the mgf properties

$$M_{X_1 - X_2}(t) = M_{X_1 + (-X_2)}(t) \stackrel{(1)}{=} M_{X_1}(t) \cdot M_{-X_2}(t) \stackrel{(2)}{=} M_{X_1}(t) + M_{X_2}(-t)$$

with 1 for independence and 2 for linear transformation properties. So considering both as bernoulli with $p = 1/2$

$$\begin{aligned} M_{X_1 - X_2}(t) &= (1 - p + pe^t)(1 - p + pe^{-t}) = \left(\frac{1}{2} + \frac{1}{2}e^t\right)\left(\frac{1}{2} + \frac{1}{2}e^{-t}\right) \\ &= \frac{1}{4} + \frac{1}{4}e^{-t} + \frac{1}{4}e^t + \frac{1}{4} = \frac{1}{2} + \frac{1}{4}(e^t + e^{-t}) \end{aligned}$$

so $M_{X_1 - X_2}(t) = 1/2 + 1/4(e^t + e^{-t})$. And Bigo confirms.

Remark 267. The following is a result that become useful sometimes (eg clt)

Proposition 7.4.6 (Mc Laurin expansion of mgf).

$$M_X(t) = 1 + t \mathbb{E}[X] + \frac{t^2}{2!} \mathbb{E}[X^2] + \frac{t^3}{3!} \mathbb{E}[X^3] + \dots \quad (7.20)$$

Dimostrazione. In general decomposition of $M_X(t)$ is like the following. Considered that mclaurin expansion of e^{tx}

$$e^{tx} = 1 + tx + \frac{(tx)^2}{2!} + \frac{(tx)^3}{3!} + \dots$$

then

$$\begin{aligned}
 M_X(t) &= \mathbb{E}[e^{tX}] = \int_{D_X} e^{tx} f(x) \, dx \\
 &= \underbrace{\int_{D_X} 1 f(x) \, dx}_{=1} + \int_{D_X} tx f(x) \, dx + \int_{D_X} \frac{(tx)^2}{2!} f(x) \, dx + \int_{D_X} \frac{(tx)^3}{3!} f(x) \, dx + \dots \\
 &= 1 + t \int_{D_X} x f(x) \, dx + \frac{t^2}{2!} \int_{D_X} x^2 f(x) \, dx + \frac{t^3}{3!} \int_{D_X} x^3 f(x) \, dx + \dots \\
 &= 1 + t \mathbb{E}[X] + \frac{t^2}{2!} \mathbb{E}[X^2] + \frac{t^3}{3!} \mathbb{E}[X^3] + \dots
 \end{aligned}$$

□

Remark 268. Now we see an example where mgf does not always exists

Example 7.4.12 (Mgf of Gamma). Let $X \sim \text{Gamma}(\alpha, \beta)$, $\alpha, \beta > 0$

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

with $D_x = [0, +\infty)$ and

$$\begin{aligned}
 \Gamma(x) &= \int_0^{+\infty} x^{\alpha-1} e^{-x} \, dx, \quad \forall \alpha > 0 \\
 \alpha \in \mathbb{N} &\implies \Gamma(\alpha) = (\alpha - 1)!
 \end{aligned}$$

Let's evaluate $M_X(t)$

$$\begin{aligned}
 M_X(t) &= \mathbb{E}[e^{tX}] = \int_0^{+\infty} e^{tx} \cdot \frac{\beta^\alpha}{\Gamma(\alpha)} x^{\alpha-1} e^{-\beta x} \, dx \\
 &= \frac{\beta^\alpha}{\Gamma(\alpha)} \int_0^{+\infty} e^{-(\beta-t)x} \cdot x^{\alpha-1} \, dx \\
 &= \frac{\beta^\alpha}{\Gamma(\alpha)} \int_0^{+\infty} e^{-(\beta-t)x} \cdot x^{\alpha-1} \frac{(\beta-t)^\alpha}{(\beta-t)^\alpha} \, dx \\
 &= \frac{\beta^\alpha}{(\beta-t)^\alpha} \underbrace{\int_0^{+\infty} \frac{(\beta-t)^\alpha}{\Gamma(\alpha)} \cdot e^{-(\beta-t)x} x^{\alpha-1} \, dx}_{=1, \text{ since } f(x) \text{ of a Gamma } (\alpha, \beta-t)}
 \end{aligned}$$

Therefore

$$M_X(t) = \frac{\beta^\alpha}{(\beta-t)^\alpha} = \left(\frac{\beta}{(\beta-t)} \right)^\alpha = \left(\frac{\beta-t}{\beta} \right)^{-\alpha} = \left(1 - \frac{t}{\beta} \right)^{-\alpha}$$

where, since $\alpha > 0$ (and it's an exponent), $M_X(t)$ is well defined only if the base is positive

$$1 - \frac{t}{\beta} > 0 \iff t < \beta$$

Example 7.4.13 (Esercizio richiesto Viroli). For this exercise:

1. compute the second moment $\mathbb{E}[X^2]$ of the binomial distribution using the second derivative of mgf evaluated in 0;
2. for the binomial, verify property 2 of mgf, that is $M_X(0) = 1$;
3. eval $\mathbb{E}[X]$ where X is Gamma by using first derivative of mgf

We have

1. per la prima deriviamo due volte e valutiamo in 0 la mgf della binomiale che è $(1 - p + pe^t)^n$. Si ha

$$\begin{aligned} [(1 - p + pe^t)^n]' &= n(1 - p + pe^t)^{n-1}(pe^t) \\ [(1 - p + pe^t)^n]'' &= n[(n-1)(1 - p + pe^t)^{n-2}(pe^t)^2 + (pe^t)(1 - p + pe^t)^{n-1}] \end{aligned}$$

che valutata per $t = 0$ da

$$n(n-1)p^2 + np = n^2p^2 - np^2 + np$$

Possiamo verificare il risultato applicando la formula di calcolo della varianza (dato che della binomiale si conoscono varianza e valore atteso)

$$\begin{aligned} \text{Var}[X] &= \mathbb{E}[X^2] - \mathbb{E}[X]^2 \\ np(1-p) &= \mathbb{E}[X^2] - n^2p^2 \\ \mathbb{E}[X^2] &= np(1-p) + n^2p^2 = np - np^2 + n^2p^2 \end{aligned}$$

2. per $t = 0$, si ha $(1 - p + pe^t)^n = (1 - p + p)^n = 1$
3. la mgf della gamma è $\left(\frac{\lambda}{\lambda - t}\right)^\alpha$ la sua derivata prima

$$\alpha \left(\frac{\lambda}{\lambda - t}\right)^{\alpha-1} \left(-\frac{\lambda(-1)}{(\lambda - t)^2}\right) = \alpha \left(\frac{\lambda}{\lambda - t}\right)^{\alpha-1} \left(\frac{\lambda}{(\lambda - t)^2}\right)$$

che valutata in $t = 0$ da α/λ , il valore atteso della gamma

Example 7.4.14 (Normal distributions). Let $X \sim N(0, 1)$, then lets derive the mgf

$$\begin{aligned} M_X(t) &= \mathbb{E}[e^{tX}] = \int_{-\infty}^{+\infty} e^{tx} \cdot \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}x^2} dx \\ &= \int_{-\infty}^{+\infty} \frac{1}{\sqrt{2\pi}} e^{tx - \frac{1}{2}x^2} dx \end{aligned}$$

Now we apply this substitution trick

$$tx - \frac{1}{2}x^2 = \frac{t^2 - (x - t)^2}{2}$$

because of the expansion

$$\frac{t^2 - (x - t)^2}{2} = \frac{t^2 - x^2 - t^2 + 2xt}{2} = tx - \frac{x^2}{2}$$

So

$$\begin{aligned}
 M_X(t) &= \int_{-\infty}^{+\infty} \frac{1}{\sqrt{2\pi}} e^{\frac{t^2 - (x-t)^2}{2}} dx \\
 &= \int_{-\infty}^{+\infty} \frac{1}{\sqrt{2\pi}} e^{\frac{t^2}{2}} e^{\frac{-(x-t)^2}{2}} dx \\
 &= e^{\frac{t^2}{2}} \underbrace{\int_{-\infty}^{+\infty} \frac{1}{\sqrt{2\pi}} e^{\frac{-(x-t)^2}{2}} dx}_{=1 \text{ since integral of } N(t,1)} \\
 &= e^{\frac{t^2}{2}}
 \end{aligned}$$

Therefore

$$X \sim N(0, 1) \iff M_X(t) = e^{t^2/2}$$

while applying properties of mgf it turns out that, if $X \sim N(0, 1)$

$$\sigma X + \mu \sim N(\mu, \sigma^2) \iff M_{\sigma X + \mu}(t) = e^{\mu t} M_X(\sigma t) = e^{\mu t} e^{\frac{\sigma^2 t^2}{2}}$$

Example 7.4.15 (Esercizio richiesto Viroli). Regarding the normal (consider $X \sim N(0, 1)$):

- prove $\frac{\partial M_{\sigma X + \mu}(t)}{\partial t} = \mu$
- derive $\mathbb{E}[X^2]$ by mgf
- check that $\text{Var}[\sigma X + \mu] = \sigma^2$ (applying $\mathbb{E}[X^2] - \mathbb{E}[X]^2$)

If the mgf of the general normal is $e^{\mu t} e^{\frac{\sigma^2 t^2}{2}}$

1. we derive it one time and evaluate for $t = 0$ to find μ

$$e^{\mu t} \cdot \mu \cdot e^{\frac{1}{2}\sigma^2 t^2} + e^{\frac{1}{2}\sigma^2 t^2} \cdot \left(\frac{1}{2}\sigma^2 2t\right) \cdot e^{\mu t} = e^{\mu t + \frac{1}{2}\sigma^2 t^2} \left(\mu + \frac{1}{2}\sigma^2 2t\right)$$

che valutata per $t = 0$ restituisce $e^0(\mu + 0) = 1 \cdot \mu = \mu$

2. la derivata seconda è

$$\begin{aligned}
 &e^{\mu t + \frac{1}{2}\sigma^2 t^2} \left(\mu + \frac{1}{2}\sigma^2 2t\right)^2 + \left(\frac{1}{2}\sigma^2 2t\right) e^{\mu t + \frac{1}{2}\sigma^2 t^2} \\
 &e^{\mu t + \frac{1}{2}\sigma^2 t^2} \left[\left(\mu + \frac{1}{2}\sigma^2 2t\right)^2 + \sigma^2\right]
 \end{aligned}$$

se $t = 0$

$$e^0 [(\mu + 0)^2 + \sigma^2] = \mu^2 + \sigma^2$$

3. abbiamo

$$\text{Var}[Y] = \mathbb{E}[Y^2] - \mathbb{E}[Y]^2 = \mu^2 + \sigma^2 - \mu^2 = \sigma^2$$

Example 7.4.16 (Esercizio viroli, primo set). Let $\mathbf{X} = \begin{bmatrix} X_1 \\ X_2 \end{bmatrix}$ be a bivariate vector with joint density $f_{\mathbf{X}}(x_1, x_2) = 2e^{-(x_1+x_2)}$ where $X_1 > X_2 > 0$

1. find $M_{\mathbf{X}}(t)$
2. compute $\mathbb{E}[X_1]$ by $M_{\mathbf{X}}(t)$
3. compute $\mathbb{E}[X_1]$ by definition
4. are $X_1 \perp\!\!\!\perp X_2$, both by density and by moment generating function

We have:

1.

$$\begin{aligned}
 M_{\mathbf{X}}(t) &= 2 \int_0^{+\infty} \int_{x_2}^{\infty} e^{tx_1} e^{tx_2} e^{-(x_1+x_2)} dx_1 dx_2 \\
 &= 2 \int_0^{+\infty} e^{-x_2(1-t_2)} \cdot \int_{x_2}^{+\infty} dx_1 dx_2 \\
 &= 2 \int_0^{\infty} e^{-x_2(1-t_2)} \cdot \left[-\frac{e^{x_1(1-t_1)}}{1-t_1} \right]_{x_2}^{\infty} dx_2 \\
 &= 2 \frac{1}{1-t_1} \int_0^{+\infty} e^{-x_2(2-t_1-t_2)} dx_2 \\
 &= \frac{2}{(1-t_1)(2-t_1-t_2)}
 \end{aligned}$$

2.

$$\begin{aligned}
 \frac{\partial M_{\mathbf{X}}(\mathbf{t})}{\partial t_1} \Big|_{\mathbf{t}=\mathbf{0}} &= 2(1-t_1)^{-2}(2-t_1-t_2)^{-1} + 2(1-t_1)^{-1}(2-t_1-t_2)^{-2} \Big|_{\mathbf{t}=\mathbf{0}} \\
 &= \frac{2}{2} + \frac{2}{4} = \frac{3}{2}
 \end{aligned}$$

3. it's longer, we have:

$$\mathbb{E}[X_1] = \int_{D_{X_1}} x_1 f_{X_1}(x_1) dx_1$$

where

$$\begin{aligned}
 f_{X_1}(x_1) &= \int_{D_{X_2}} f_{\mathbf{X}}(x_1, x_2) dx_2 \\
 &= \int_0^{x_1} 2e^{-(x_1+x_2)} dx_2 = \int_0^{x_1} 2e^{-x_1} e^{-x_2} dx_2 \\
 &= 2e^{-x_1} \cdot \int_0^{x_1} e^{-x_2} dx_2 = 2e^{-x_1} [-e^{-x_2}]_0^{x_1} \\
 &= 2e^{-x_1}(1 - e^{-x_1}) = 2e^{-x_1} - 2e^{-2x_1}
 \end{aligned}$$

therefore

$$\begin{aligned}\mathbb{E}[X_1] &= \int_0^{+\infty} x_1 (2e^{-x_1} - 2e^{-2x_1}) dx_1 \\ &= 2 \underbrace{\int_0^{+\infty} x_1 e^{-x_1} dx_1}_{\text{expected value of Exp (1)}} - \underbrace{\int_0^{+\infty} x_1 2e^{-2x_1} dx_1}_{\text{expected value of Exp (2)}} \\ &= 2 \cdot 1 - \frac{1}{2} = \frac{3}{2}\end{aligned}$$

4. by the density

$$f_{X_2}(x_2) = \int_{x_2}^{+\infty} 2e^{-(x_1+x_2)} dx_1 = 2e^{-x_1} \cdot [-e^{-x_1}]_{x_2}^{+\infty} = e^{-x_2} e^{-x_2} = e^{-2x_2}$$

Now we check if $f_{X_1}(x_1) \cdot f_{X_2}(x_2) = f_{\mathbf{X}}(x_1, x_2)$:

$$2e^{-x_1}(1 - e^{-x_1})e^{-2x_2} \neq 2e^{-(x_1+x_2)}$$

therefore they are not independent.

Now let's check according to the moment generating function; we observe that:

$$M_{X_1}(t_1) = M_{\mathbf{X}}(t_1, 0) = \frac{2}{(1-t_2)^{\frac{1}{2-t_1}}} \quad M_{X_2}(t_2) = M_{\mathbf{X}}(0, t_2) = \frac{2}{2-t_2}$$

Since $M_{\mathbf{X}}(\mathbf{t}) \neq M_{X_1}(t_1)M_{X_2}(t_2)$ are not independent.

Note: in case of mutually independent rvs:

$$\begin{aligned}f_{\mathbf{X}}(\mathbf{x}) &= \prod_{i=1}^p f_{X_i}(x_i) \\ F_{\mathbf{X}}(\mathbf{x}) &= \prod_{i=1}^p F_{X_i}(x_i) \\ M_{\mathbf{X}}(\mathbf{t}) &= \prod_{i=1}^p M_{X_i}(t_i)\end{aligned}$$

Example 7.4.17 (Mgf of Geometric and Negative binomial). Let $X_1, \dots, X_n \sim \text{Geom}(p)$ iid rvs. Find $M_Y(t)$ where $Y = \sum_{i=1}^n X_i$. What can you say about the distribution of Y ?

For a geometric rv we have

$$\mathbb{P}(X = x) = p(1-p)^{x-1}, \quad D_X = \{1, 2, \dots\}$$

so

$$\begin{aligned}M_X(t) &= \sum_{x=1}^{\infty} e^{tx} p(1-p)^{x-1} = \sum_{x=1}^{\infty} e^{tx} p \frac{1-p}{1-p} (1-p)^{x-1} \\ &= \frac{p}{1-p} \cdot \sum_{x=1}^{\infty} [e^t(1-p)]^x = \frac{p}{1-p} \cdot \left(\sum_{x=0}^{\infty} [e^t(1-p)]^x - 1 \right)\end{aligned}$$

Now we define $q = 1 - p$; if $|e^t(1 - p)| < 1$ the previous series converges to $\frac{1}{1 - qe^t}$. Therefore the $M_X(t)$ exists only for $e^t < \frac{1}{1-p}$, that is $t < -\log(1 - p)$. For such values we have

$$M_X(t) = \frac{p}{q} \left(\frac{1}{1 - qe^t} - 1 \right) = \frac{p}{q} \left(\frac{qe^t}{1 - qe^t} \right) = \frac{pe^t}{1 - qe^t}$$

Now

$$M_Y(t) = \prod_{i=1}^n M_{X_i}(t) = \left[\frac{pe^t}{1 - qe^t} \right]^n$$

with the last being the moment generating function of a negative binomial distribution with parameters n and p

7.5 Conditional distribution

7.5.1 Definition and examples

Remark 269. Roughly speaking the problem is: given 2 real random variable Y, X we aim to evaluate the distribution of Y given that $X = x$.

Ad esempio se $X = \text{altezza}$ e $Y = \text{peso}$ vorrei conoscere la distribuzione del peso nell'ipotesi che l'altezza sia 1.70

Definition 7.5.1 (Conditional distribution). Let $\begin{bmatrix} X \\ Y \end{bmatrix}$ be a bivariate rv; the conditional distribution of Y given X is any function

$$\mathbb{P} \left(\begin{bmatrix} X \\ Y \end{bmatrix} \in C | X = x \right)$$

defined $\forall x \in \mathbb{R}, \forall C \in \beta(\mathbb{R}^n)$ satisfying the following properties:

1. for each fixed $x \in \mathbb{R}$, the map

$$C \rightarrow \mathbb{P} \left(\begin{bmatrix} X \\ Y \end{bmatrix} \in C | X = x \right)$$

is a probability measure on $\beta(\mathbb{R}^2)$

2. for each fixed $C \in \beta(\mathbb{R}^2)$, the map

$$x \rightarrow \mathbb{P} \left(\begin{bmatrix} X \\ Y \end{bmatrix} \in C | X = x \right)$$

is Borel measurable and satisfies

$$\mathbb{P} \left(\begin{bmatrix} X \\ Y \end{bmatrix} \in C \right) = E_X \left\{ \mathbb{P} \left(\begin{bmatrix} X \\ Y \end{bmatrix} \in C | X = x \right) \right\}$$

where E_X means expectation with respect to X .

Important remark 59 (Important remarks). Some important remarks:

NB: direi gli input siano
due, C e x

1. it can be shown that the conditional distribution of Y given X (namely a function satisfying definition) *always exists* and is *almost surely unique*. This remark is important because looking at the defn it's not sure that any function such as that defined exists. Here a.s uniqueness is meant with respect to the probability distribution of X

2. the notation $\mathbb{P}\left(\begin{bmatrix} X \\ Y \end{bmatrix} \in C | X = x\right)$ is very useful but also quite dangerous; it should be regarded as “the conditional probability that $\begin{bmatrix} X \\ Y \end{bmatrix} \in C$ given $X = x$. Be careful however: by definition $\mathbb{P}\left(\begin{bmatrix} X \\ Y \end{bmatrix} \in C | X = x\right)$ is only a function satisfying the two properties of the definition. In particular it's *not necessarily equal* to ratio between probability intersection divided by probability $P(X = x)$:

$$\text{it's not necessarily } \frac{\mathbb{P}\left(X = x, \begin{bmatrix} X \\ Y \end{bmatrix} \in C\right)}{\mathbb{P}(X = x)}$$

This should be obvious since it may be that $P(X = x) = 0, \forall x \in \mathbb{R}$ (which is possible eg in continuous distribution).

For instance suppose $P(X = x) = 0, \forall x \in \mathbb{R}$ (or equivalently the distribution function is continuous) then by the previous remark $\mathbb{P}\left(\begin{bmatrix} X \\ Y \end{bmatrix} \in C | X = x\right)$ exists, but it certainly does not coincide with the ratio above (not defined)

3. if X is **discrete** the operator $E_X(\cdot)$ means

$$E_X \left\{ \mathbb{P} \left(\begin{bmatrix} X \\ Y \end{bmatrix} \in C | X = x \right) \right\} = \sum_{x \in B} \mathbb{P}(X = x) \mathbb{P} \left(\begin{bmatrix} X \\ Y \end{bmatrix} \in C | X = x \right)$$

where B is any set, finite or countable such that $\mathbb{P}(X \in B) = 1$.

Similarly if X is **absolutely continuous** the operator $E_X(\cdot)$ means

$$E_X \left\{ \mathbb{P} \left(\begin{bmatrix} X \\ Y \end{bmatrix} \in C | X = x \right) \right\} = \int_{-\infty}^{+\infty} f(x) \mathbb{P} \left(\begin{bmatrix} X \\ Y \end{bmatrix} \in C | X = x \right)$$

where f is the density of X

Important remark 60 (Some useful properties of conditional distribution). Some properties:

- any time we aim to evaluate the conditional probability we can substitute as follows:

$$\mathbb{P} \left(\begin{bmatrix} X \\ Y \end{bmatrix} \in C | X = x \right) = \mathbb{P} \left(\begin{bmatrix} x \\ Y \end{bmatrix} \in C | X = x \right) \quad (7.21)$$

This is intuitively obvious, since we're conditioning on $X = x$ we know that $X = x$ and can substitute it within parenthesis;

- if $X \perp\!\!\!\perp Y$, then

$$\mathbb{P}\left(\begin{bmatrix} X \\ Y \end{bmatrix} \in C | X = x\right) = \mathbb{P}\left(\begin{bmatrix} x \\ Y \end{bmatrix} \in C | X = x\right) = \mathbb{P}\left(\begin{bmatrix} x \\ Y \end{bmatrix} \in C\right)$$

where in the last passage i can drop the conditioning because X and Y are independent

- if we know that $\mathbb{P}(X \in A) = 1$ for some $A \in \beta(\mathbb{R}^n)$, then is enough to evaluate

$$\mathbb{P}\left(\begin{bmatrix} X \\ Y \end{bmatrix} \in C | X = x\right), \quad \forall x \in A$$

and not necessarily $\forall x \in \mathbb{R}$.

For example if $X \sim \text{Unif}(0, 1)$, its enough to evaluate

$$\mathbb{P}\left(\begin{bmatrix} X \\ Y \end{bmatrix} \in C | X = x\right), \quad \forall x \in (0, 1)$$

Remark 270. Unfortunately, in general there is not an intuitive formula to evaluate conditional distributions (there is in some cases as we'll see later).

Remark 271. Let's see some examples, the first of which is fundamental

NB: qui ho solo invertito $\mathbb{P}(Y = X)$ per facilitare la memorizzazione

Example 7.5.1 ($\mathbb{P}(Y = X)$: a usual question at the Rigo exam). Suppose $X \perp\!\!\!\perp Y$ and Y has a continuous distribution function. What is the $\mathbb{P}(X = Y)$? This should be 0. Let's show it.

To answer let's define $C = \{(x, y) \in \mathbb{R}^2 : x = y\}$ which is the set of points constituting the diagonal

$$\begin{aligned} \mathbb{P}(Y = X) &= \mathbb{P}\left(\begin{bmatrix} X \\ Y \end{bmatrix} \in C\right) \stackrel{(1)}{=} E_X \left\{ \mathbb{P}\left(\begin{bmatrix} X \\ Y \end{bmatrix} \in C | X = x\right) \right\} \\ &\stackrel{(2)}{=} E_X \left\{ \mathbb{P}\left(\begin{bmatrix} x \\ Y \end{bmatrix} \in C | X = x\right) \right\} \stackrel{(3)}{=} E_X \left\{ \mathbb{P}\left(\begin{bmatrix} x \\ Y \end{bmatrix} \in C\right) \right\} \\ &= E_X \underbrace{(\mathbb{P}(Y = x))}_0 \stackrel{(4)}{=} E_X(0) = 0 \end{aligned}$$

with:

- (1) by property 2 of defn,
- (2) by 7.21 (since we're conditioning I can write x instead of X)
- (3) since they are independent I can drop the conditioning
- (4) since being Y continuous, the probability that $Y = x$ (aka a single value) is zero

NB: fatto solo il mio anno?

Remark 272. Note that: in statistical inference the elements of the sample are often assumed to be iid. Under this assumption, if the distribution of the character in the population is *continuous* what is the prob of having the sample with all different observation?

It's 1 (almost sure event). This because $\mathbb{P}(X_i = X_j) = 0, \forall i \neq j$, so that the probability that $\mathbb{P}(X_1, \dots, X_n \text{ are all distinct}) = 1$

in questo ho invertito i ruoli di X e Y se mi incartavo. Nota che nel caso precedente Y era continua, qui è X ad essere continua.

Example 7.5.2 ($\mathbb{P}(Y = \sin(X))$). Suppose $X \perp\!\!\!\perp Y$, and X has a continuous distribution function. Let's prove that $\mathbb{P}(Y = \sin(X)) = 0$.

A quick way to do it is the following: since Y is independent of X then is still independent of any transformation (and thus $\sin(X)$).

Hence by exercise 7.5.1 it suffices to prove either/equivalently that:

- $\sin(X)$ has a continuous distribution function (because if it's continuous we can repeat the argument of the previous exercise)
- $\mathbb{P}(\sin(X) = a) = 0, \forall a \in \mathbb{R}$ (this is trivially true if $a \notin [-1, 1]$)

We follow the second way, supposing $a \in [-1, 1]$ and define:

$$I_a = \{x \in \mathbb{R} : \sin x = a\}$$

we have that I_a is countable (pensa sull'asse delle x , ci sono infiniti punti di seno che hanno altezza a se questo è tra -1 e 1). Thus the probability:

$$\mathbb{P}(\sin(X) = a) = \mathbb{P}(X \in I_a) \stackrel{(1)}{=} \sum_{x \in I_a} \mathbb{P}(X = x) \stackrel{(2)}{=} \sum_{x \in I_a} 0 = 0$$

with (1) since I_a is countable and (2) because X is a continuous distribution function

Example 7.5.3 ($\mathbb{P}(Y = X)$ with independent discrete distribution). What can be said about $\mathbb{P}(Y = X)$ if $X \perp\!\!\!\perp Y$ and they are both discrete?

Since X is discrete, there's a set $B \subset \mathbb{R}$ such that $\mathbb{P}(X \in B) = 1$. Hence the $\mathbb{P}(X = Y)$ can be written as

$$\begin{aligned} \mathbb{P}(X = Y) &= \sum_{x \in B} \mathbb{P}(X = x, Y = X) = \sum_{x \in B} \mathbb{P}(X = x, Y = x) \\ &\stackrel{(\perp)}{=} \sum_{x \in B} \mathbb{P}(X = x) \mathbb{P}(Y = x) \end{aligned}$$

and this may be > 0 .

Example 7.5.4. Let A, B, C be iid with all $\sim \text{Exp}(1)$ ¹. Let's define the random parabola (named like this because coefficient a, b, c are rvs):

$$f(x) = Ax^2 + Bx + C, \quad \forall x \in \mathbb{R}$$

What is the probability that f has real roots?

It is the probability $\mathbb{P}(B^2 - 4AC \geq 0)$. To evaluate it, we have to choose one

¹Always positive rv with density function e^{-x} and distribution $1 - e^{-x}$

of the three variable and condition on it; eg let's condition on C :

$$\begin{aligned}
 \mathbb{P}(B^2 - 4AC \geq 0) &= E_C \{ \mathbb{P}(B^2 \geq 4AC | C = c) \} \\
 &= E_C \{ \mathbb{P}(B^2 \geq 4Ac | C = c) \} \\
 &\stackrel{(\perp)}{=} E_C \{ \mathbb{P}(B^2 \geq 4Ac) \} \\
 &= \int_{-\infty}^{+\infty} \mathbb{P}(B^2 \geq 4Ac) f(c) \, dc \\
 &\stackrel{(1)}{=} \int_0^{+\infty} \mathbb{P}(B^2 \geq 4Ac) e^{-c} \, dc \\
 &\stackrel{(2)}{=} \int_0^{+\infty} E_A \{ \mathbb{P}(B^2 \geq 4Ac) | A = a \} e^{-c} \, dc \\
 &= \int_0^{+\infty} E_A \{ \mathbb{P}(B^2 \geq 4ac) \} e^{-c} \, dc \\
 &= \int_0^{+\infty} \int_0^{+\infty} \mathbb{P}(B^2 \geq 4ac) e^{-a} e^{-c} \, da \, dc \\
 &= \int_0^{+\infty} \int_0^{+\infty} \mathbb{P}(B \geq 2\sqrt{ac}) e^{-a} e^{-c} \, da \, dc \\
 &\stackrel{(3)}{=} \int_0^{+\infty} \int_0^{+\infty} e^{-2\sqrt{ac}} e^{-a} e^{-c} \, da \, dc
 \end{aligned}$$

where in

- (1) since C is exponential (and with density only on positive side)
- (2) we have to evaluate $\mathbb{P}(B^2 \geq 4Ac)$ and its convenient to do it conditioning further on A
- (3) considered that B is exponential (if $Z \sim \text{Exp}(\lambda)$ then $\mathbb{P}(Z > z) = 1 - \mathbb{P}(Z \leq z) = 1 - (1 - e^{-\lambda z}) = e^{-\lambda z}$).

NB: fatto solo il mio anno?

Example 7.5.5. Suppose $X \perp\!\!\!\perp Y$, $X \sim \text{Unif}(0, 1)$ and $Y \sim N(0, 1)$. We want to evaluate the distribution function of the product XY . Here conditional distribution become handy. For all $a \in \mathbb{R}$, by definition the distribution function is

$$\begin{aligned}
 \mathbb{P}(XY \leq a) &\stackrel{(1)}{=} E_X(\mathbb{P}(XY \leq a | X = x)) = E_X(\mathbb{P}(xY \leq a | X = x)) \\
 &\stackrel{(\perp)}{=} E_X(\mathbb{P}(xY \leq a)) \stackrel{(2)}{=} \int_{-\infty}^{+\infty} \mathbb{P}(xY \leq a) f(x) \, dx \\
 &= \int_0^1 \mathbb{P}(xY \leq a) 1 \, dx \stackrel{(3)}{=} \int_0^1 \mathbb{P}\left(N(0, 1) \leq \frac{a}{x}\right) \, dx
 \end{aligned}$$

with (1) by definition, (2) since X is uniform (absolutely continuous), (3) because Y is normal.

After this we go to our friend mathematician asking for help.

7.5.2 Formula to calculate it?

Important remark 61. How to calculate $\mathbb{P}\left(\begin{bmatrix} X \\ Y \end{bmatrix} \in C | X = x\right)$? We know this object exists and in many problem its enough to know it.

Unfortunately *there is not* a general formula which allows to calculate the probability above in every situation. Such a formula exists in *two special cases*:

1. X discrete
2. $\begin{bmatrix} X \\ Y \end{bmatrix}$ absolutely continuous

Definition 7.5.2 (Discrete case). If X is discrete, there is a set $B \subset \mathbb{R}$, B finite or countable, $\mathbb{P}(X \in B) = 1$, and $\mathbb{P}(X = x) > 0$, $\forall x \in B$ (true by definition of discreteness). Hence it suffices to let

$$\mathbb{P}\left(\begin{bmatrix} X \\ Y \end{bmatrix} \in C | X = x\right) = \frac{\mathbb{P}\left(X = x, \begin{bmatrix} X \\ Y \end{bmatrix} \in C\right)}{\mathbb{P}(X = x)}, \forall x \in B, \forall C \in \beta(\mathbb{R}^2)$$

This is just the base definition of conditional probability with positive denominator, being the distribution discrete and focusing on $x \in B$.

Definition 7.5.3 (Continuous case). If $(X, Y)^\top$ is absolutely continuous with joint density $f(x, y)$, then the conditional distribution of Y given $X = x$ is still absolutely continuous with *conditional density*:

$$h(y|x) = \frac{f(x, y)}{f_1(x)}, \quad \text{where } f_1(x) > 0$$

where $f(x, y)$ is the joint density of $(X, Y)^\top$ and f_1 is the marginal density of X (namely the integral $f(x, y)$ in dy)

In this case we have an explicit formula for *conditional distribution* function of Y given $X = x$ as:

$$\mathbb{P}(Y \leq y | X = x) = \int_{-\infty}^y \frac{f(x, t)}{f_1(x)} dt, \quad \forall x, y \in \mathbb{R} : f_1(x) > 0$$

Example 7.5.6. Suppose $X \sim \text{Unif}(0, 1)$ and $Y|X = x \sim \text{Bin}(n, x)$, $\forall x \in (0, 1)$. Find the conditional distribution of X given Y .

To this end we first note that, being Y discrete

TODO: esempio fatto solo quest'anno, non nel mio

$$\mathbb{P}(Y \in \{0, 1, \dots, n\}) = E_X \left\{ \underbrace{\mathbb{P}(Y \in \{0, 1, \dots, n\} | X = x)}_{=1} \right\} = E_X(1) = 1$$

Again Y is discrete so that

$$\begin{aligned}
 \mathbb{P}(X \in A | Y = y) &= \frac{\mathbb{P}(X \in A, Y = y)}{\mathbb{P}(Y = y)} \\
 &= \frac{E_x \{\mathbb{P}(X \in A, Y = y | X = x)\}}{E_x \{\mathbb{P}(Y = y | X = x)\}} \\
 &= \frac{E_x \{\mathbb{P}(X \in A, Y = y | X = x)\}}{E_x \left\{ \binom{n}{y} x^y (1-x)^{n-y} \right\}} \\
 &= \frac{E_x \left\{ 1_A(x) \binom{n}{y} x^y (1-x)^{n-y} \right\}}{E_x \left\{ \binom{n}{y} x^y (1-x)^{n-y} \right\}} \\
 &= \frac{\int_0^1 1_A(x) \binom{n}{y} x^y (1-x)^{n-y} dx}{\int_0^1 x^y (1-x)^{n-y} dx}
 \end{aligned}$$

Example 7.5.7. If $n = 2$, $X_{(1)} = \min(X_1, X_2)$ and $X_{(2)} = \max(X_1, X_2)$

Example 7.5.8 (Example with order statistics). Let S and T be iid with $S \sim \text{Unif}(0, 1)$. Define $X = \min(S, T)$ and $Y = \max(S, T)$. We want to write the conditional distribution of Y given $X = x$.

To this end we first note that for theorem 7.2.6 $\begin{bmatrix} X \\ Y \end{bmatrix}$ is absolutely continuous being the order statistic of corresponding to $\begin{bmatrix} S \\ T \end{bmatrix}$, and S, T are iid and absolutely continuous.

In addition the joint density f of $\begin{bmatrix} X \\ Y \end{bmatrix}$ is

$$f(x, y) = \begin{cases} 2! \cdot g(x)g(y) & \text{if } x < y \\ 0 & \text{otherwise} \end{cases}$$

where g is density common to S, T , namely $\text{Unif}(0, 1)$

$$g_1(x) = \begin{cases} 1 & \text{if } x \in (0, 1) \\ 0 & \text{otherwise} \end{cases}$$

Hence

$$f(x, y) = \begin{cases} 2 & \text{if } 0 < x < y < 1 \\ 0 & \text{otherwise} \end{cases}$$

the marginals are (eg)

$$g(x) = \int f(x, y) dy = \int_x^1 2 dy = 2(1-x), \quad \forall x \in (0, 1)$$

Hence since $\begin{bmatrix} X \\ Y \end{bmatrix}$ is absolutely continuous, we have the formula for the *conditional distribution* of Y given X , which is still absolutely continuous with

density:

$$h(y|x) = \begin{cases} \frac{f(x,y)}{g(x)} = \frac{2}{2(1-x)} = \frac{1}{1-x}, & \text{if } 0 < x < y < 1 \\ 0, & \text{elsewhere} \end{cases}$$

So $\forall x \in (0,1)$ (select first the min in $(0,1)$) $Y|X = x \sim \text{Unif}(x,1)$ (the max is uniformly distributed between the min and 1, distribution max). A bits of *interpretation* on the results: Since S and T are iid $\text{Unif}(0,1)$ observing the pair (S,T) is like to select "at random" a point from the unit square.

Suppose now that $X = \min(S,T)$; what can be said about $Y = \max(S,T)$? Certainly $Y > X$ so if we fix a point $x \in [0,1]$, y is above the diagonal $y = x$, that is it is in the $[x,1]$. In fact the distribution of $Y|X \sim \text{Unif}(x,1)$: and this is why we obtained $1/(1-x)$ as density (coming from that distribution)

7.6 Multivariate normal

Remark 273. Let's start from univariate and see that multivariate formula are univariate generalization

NB: quest'anno è partito secco dalla multivariate e pace

Proposition 7.6.1 (Characteristic functions of univariate normal). *If $Z \sim N(0,1)$ and $X \sim N(\mu, \sigma^2)$ then $\forall t \in \mathbb{R}$:*

$$\varphi_Z(t) = e^{-t^2/2} \quad (7.22)$$

$$\varphi_X(t) = e^{it\mu - \frac{1}{2}(t\sigma)^2} \quad (7.23)$$

Dimostrazione. If $Z \sim N(0,1)$: then its characteristic function is

$$\varphi_Z(t) = \mathbb{E}[e^{itZ}] = \int_{-\infty}^{+\infty} e^{itx} \frac{\exp(-\frac{1}{2}x^2)}{\sqrt{2\pi}} dx \stackrel{(1)}{=} \dots = e^{-t^2/2}, \quad \forall t \in \mathbb{R}$$

(in (1) after doing calculation). If $X \sim N(\mu, \sigma^2)$ then X can be written as $X = \mu + \sigma Z$ with $Z \sim N(0,1)$ and thus we can derive the formula given the definition (in the univariate case) as:

$$\begin{aligned} \varphi_X(t) &= \mathbb{E}[e^{it(\mu + \sigma Z)}] = \mathbb{E}\left[\underbrace{e^{it\mu}}_{\text{constant}} e^{it\sigma Z}\right] = e^{it\mu} \mathbb{E}[e^{i(t\sigma)Z}] \\ &= e^{it\mu} \cdot \varphi_Z(t\sigma) = e^{it\mu - \frac{1}{2}(t\sigma)^2} \quad \forall t \in \mathbb{R} \end{aligned}$$

□

Remark 274. MVN is not so important for this course: it's very important for statistician, but from point of view of probability it's just a special distribution among the others.

Definition 7.6.1. An n -dimensional random vector $X = (X_1, \dots, X_n)^T$ is normally distributed with parameters $\mu \in \mathbb{R}^n$ and Σ (a $n \times n$ symmetric non-negative definite/ positive semi-definite (≥ 0) matrix ²), if the characteristic

²Recappino da wikipedia: For any real invertible matrix X the product $X^T X$ is a positive definite matrix (if the means of the columns of X are 0, then this is also called the covariance matrix). A simple proof is that for any non-zero vector \mathbf{z} , the condition $\mathbf{z}^T X^T X \mathbf{z} = (X\mathbf{z})^T (X\mathbf{z}) = \|X\mathbf{z}\|^2 > 0$, since the invertibility of matrix X means that $X\mathbf{z} \neq 0$.

function of X is given by:

$$\varphi_X(t) = \mathbb{E} \left[e^{it^T \boldsymbol{\mu} - \frac{1}{2} t^T \boldsymbol{\Sigma} t} \right], \quad \forall t \in \mathbb{R}^n$$

Important remark 62 (Meaning of parameters). $\boldsymbol{\mu}$ is the vector of the mean, $\boldsymbol{\Sigma}$ is the so called covariance matrix which have variances on the diagonal, covariance out of main diagonal

$$\boldsymbol{\mu} = \begin{bmatrix} \mathbb{E}[X_1] \\ \vdots \\ \mathbb{E}[X_n] \end{bmatrix}, \quad \boldsymbol{\Sigma} = \begin{bmatrix} \text{Var}[X_1] & \text{Cov}(X_1, X_2) & \dots & \text{Cov}(X_1, X_n) \\ \text{Cov}(X_2, X_1) & \text{Var}[X_2] & \dots & \text{Cov}(X_2, X_n) \\ \vdots & \vdots & \ddots & \vdots \\ \text{Cov}(X_n, X_1) & \dots & \dots & \text{Var}[X_n] \end{bmatrix}$$

Remark 275. Our definition includes not only absolutely continuous normal vector, but also other (eg degenerate in some cases)

Important remark 63 (Main properties). Some remarks:

- a **linear transformation** of a normal random variable is still normal: if $\mathbf{X} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ then

$$\mathbf{Y} = \mathbf{a} + \mathbf{B}\mathbf{X} \sim N(\mathbf{a} + \mathbf{B}\boldsymbol{\mu}, \mathbf{B}\boldsymbol{\Sigma}\mathbf{B}^T)$$

where $\mathbf{a} \in \mathbb{R}^m$ the matrix \mathbf{B} is $m \times n$.

In particular, if X is normal, all marginals are still normal being the marginal obtained via a linear transformation (therefore we get a normal) that merely extract the marginal/subset. Eg

$$\mathbf{Y} = \begin{bmatrix} X_1 \\ X_2 \\ X_4 \end{bmatrix} = \mathbf{0} + \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} X_1 \\ X_2 \\ X_3 \\ X_4 \end{bmatrix} = \mathbf{0} + \mathbf{B}\mathbf{X}$$

- if $X \sim \text{MVN}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ then independence of compenence amounts to null covariances: if $\boldsymbol{\Sigma}$ is diagonal, then X_1, \dots, X_n are independent;
- regarding $\det \boldsymbol{\Sigma}$:

- if $\det(\boldsymbol{\Sigma}) > 0$, then $\boldsymbol{\Sigma}$ is also positive-definite (> 0 , not only ≥ 0) and \mathbf{X} is absolutely continuous with density

$$f(\mathbf{x}) = (2\pi)^{-\frac{n}{2}} (\det \boldsymbol{\Sigma})^{-\frac{1}{2}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\}, \quad \forall \mathbf{x} \in \mathbb{R}^n$$

The univariate density we know is a special case where the matrix $\boldsymbol{\Sigma}$ is positive definite (otherwise matrix can be inverted).

For $n = 1$ the density $\boldsymbol{\Sigma}$ reduce to a scalar ($\boldsymbol{\Sigma} = \sigma^2$, variance of the variable) and $\boldsymbol{\mu}$ to a single number

$$f(x) = \frac{\exp \left(-\frac{(x-\mu)^2}{\sigma^2} \right)}{\sigma \sqrt{2\pi}}$$

- if otherwise $\det \Sigma = 0$ then X is still normal, but the distribution of X is no longer absolutely continuous.
For instance if $n = 1$ and $\sigma^2 = \Sigma = 0$ then

$$\varphi_X(t) = e^{-it\mu}$$

and $X = \mu$ is degenerate. In other terms if $n = 1$, the above definition implies that the degenerate random variable are considered normal.

Linear transformation proof. In order to prove that if $\mathbf{X} \sim \text{MVN}(\boldsymbol{\mu}, \Sigma)$ then $\mathbf{Y} = \mathbf{a} + \mathbf{B}\mathbf{X} \sim \text{MVN}(\mathbf{a} + \mathbf{B}\boldsymbol{\mu}, \mathbf{B}\Sigma\mathbf{B}^\top)$ we write the characteristic function of \mathbf{Y} according to the definition above. Let's evaluate it: NB: l'ha fatta il mio anno, non questo

$$\begin{aligned}\mathbb{E}\left[e^{it^\top \mathbf{Y}}\right] &= \mathbb{E}\left[e^{it^\top (\mathbf{a} + \mathbf{B}\mathbf{X})}\right] = \mathbb{E}\left[\underbrace{e^{it^\top \mathbf{a}}}_{\text{constant}} \cdot e^{it^\top \mathbf{B}\mathbf{X}}\right] = e^{it^\top \mathbf{a}} \underbrace{\mathbb{E}\left[e^{it^\top \mathbf{B}\mathbf{X}}\right]}_{\varphi_X(\mathbf{B}^\top \mathbf{t})} \\ &= e^{it^\top \mathbf{a}} \cdot e^{it^\top \mathbf{B}\boldsymbol{\mu} - \frac{1}{2} \mathbf{t}^\top \mathbf{B}\Sigma\mathbf{B}^\top \mathbf{t}} = \exp\left(it^\top (\mathbf{a} + \mathbf{B}\boldsymbol{\mu}) - \frac{1}{2} \mathbf{t}^\top (\mathbf{B}\Sigma\mathbf{B}^\top) \mathbf{t}\right) \\ &\iff Y \sim \text{N}(\mathbf{a} + \mathbf{B}\boldsymbol{\mu}, \mathbf{B}\Sigma\mathbf{B}^\top)\end{aligned}$$

□

Example 7.6.1 (Assignment 1 Viroli, Exercise 3). Suppose that \mathbf{X} is a bivariate Gaussian vector with components (X_1, X_2) which are marginally standard normally distributed and with correlations $1/2$:

1. What is the distribution of $Y_1 = 2X_1 - X_2$ and $Y_2 = X_1 - X_2/2$
2. find the linear transformation from \mathbf{X} to \mathbf{Y} and ask what is the distribution of \mathbf{Y}

Since $X_1, X_2 \sim \text{N}(0, 1)$ and considered that

$$\begin{aligned}\frac{1}{2} &= \text{Corr}(X_1, X_2) = \frac{\text{Cov}(X_1, X_2)}{\sqrt{\text{Var}[X_1]}\sqrt{\text{Var}[X_1]}} = \frac{\text{Cov}(X_1, X_2)}{1 \cdot 1} \\ &= \text{Cov}(X_1, X_2) = \text{Cov}(X_2, X_1)\end{aligned}$$

1. if $Y_1 = 2X_1 - X_2$ and $Y_2 = X_1 - X_2/2$, then Y_1, Y_2 will be linear combinations of correlated normals; the distributions of Y_1, Y_2 will be normals with mean the linear combinations of means:

$$\begin{aligned}\mathbb{E}[Y_1] &= \mathbb{E}[2X_1 - X_2] = 2\mathbb{E}[X_1] - \mathbb{E}[X_2] = 0 \\ \mathbb{E}[Y_2] &= \mathbb{E}\left[X_1 - \frac{1}{2}X_2\right] = \mathbb{E}[X_1] - \frac{1}{2}\mathbb{E}[X_2] = 0\end{aligned}$$

Applying $\text{Var}[aX + bY] = a^2 \text{Var}[X] + b^2 \text{Var}[Y] + 2ab \text{Cov}(X, Y)$ we have:

$$\begin{aligned}\text{Var}[Y_1] &= \text{Var}[2X_1 - X_2] = 4\text{Var}[X_1] + \text{Var}[X_2] + 2 \cdot 2(-1) \text{Cov}(X_1, X_2) \\ &= 4 + 1 - 4 \cdot \frac{1}{2} = 5 - 2 = 3\end{aligned}$$

$$\begin{aligned}\text{Var}[Y_2] &= \text{Var}\left[X_1 - \frac{1}{2}X_2\right] = \text{Var}[X_1] + \frac{1}{4} \text{Var}[X_2] + 2\left(-\frac{1}{2}\right) \text{Cov}(X_1, X_2) \\ &= 1 + \frac{1}{4} - \frac{1}{2} = \frac{3}{4}\end{aligned}$$

Therefore: $Y_1 \sim N(0, 3)$, $Y_2 \sim N(0, \frac{3}{4})$

2. in general, a linear trasformation of a multivariate normal is still normal; if $\mathbf{X} \sim \text{MVN}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ is a n -dimensional random vector and $\mathbf{Y} = \mathbf{A}\mathbf{X} + \mathbf{b}$, with \mathbf{A} an $m \times n$ matrix and $\mathbf{b} \in \mathbb{R}^m$, then \mathbf{Y} is a m -dimensional random vector and specifically $\mathbf{Y} \sim \text{MVN}(\mathbf{A}\boldsymbol{\mu} + \mathbf{b}, \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^\top)$.

In our case $m = n = 2$ and we have:

$$\mathbf{X} = \begin{bmatrix} X_1 \\ X_2 \end{bmatrix} \sim \text{MVN} \left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 & 1/2 \\ 1/2 & 1 \end{bmatrix} \right), \quad \mathbf{Y} = \begin{bmatrix} Y_1 \\ Y_2 \end{bmatrix} = \begin{bmatrix} 2X_1 - X_2 \\ X_1 - \frac{1}{2}X_2 \end{bmatrix}$$

so $\mathbf{Y} = \mathbf{A}\mathbf{X} + \mathbf{0}$, represent the linear transformation needed to obtain \mathbf{Y} , where

$$\mathbf{A} = \begin{bmatrix} 2 & -1 \\ 1 & -1/2 \end{bmatrix}$$

Therefore to evaluate the parameters of \mathbf{Y} :

$$\begin{aligned} \mathbf{A}\boldsymbol{\mu} + \mathbf{b} &= \begin{bmatrix} 2 & -1 \\ 1 & -1/2 \end{bmatrix} \begin{bmatrix} 0 \\ 0 \end{bmatrix} + \begin{bmatrix} 0 \\ 0 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \\ \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^\top &= \begin{bmatrix} 2 & -1 \\ 1 & -1/2 \end{bmatrix} \begin{bmatrix} 1 & 1/2 \\ 1/2 & 1 \end{bmatrix} \begin{bmatrix} 2 & 1 \\ -1 & -1/2 \end{bmatrix} = \begin{bmatrix} 3 & 3/2 \\ 3/2 & 3/4 \end{bmatrix} \end{aligned}$$

Finally:

$$\mathbf{Y} \sim \text{MVN} \left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 3 & 3/2 \\ 3/2 & 3/4 \end{bmatrix} \right)$$

Capitolo 8

Convergences and related topics

8.1 Convergence

Remark 276. Given a sequence X_1, \dots, X_n, \dots of real random variables and another real random variable X :

- we are interested in checking whether or not X_n converges to X as n goes to $+\infty$, written $X_n \rightarrow X$.
- there are 4 types/modes of convergence: in each case as n become larger, X_n get “closer” to X ; but *the way this happens is different* so one convergence does not necessarily imply others (we will see relationship between them in the following).

Remark 277. In the following all the standard calculus limits involved are meant for $n \rightarrow +\infty$.

Definition 8.1.1 (Almost sure convergence). X_n converge almost surely to X and we write $X_n \xrightarrow{a.s.} X$ if and only if

$$\begin{aligned}\mathbb{P}(\omega \in \Omega : X_n(\omega) \rightarrow X(\omega)) &= \mathbb{P}\left(\omega \in \Omega : \lim_{n \rightarrow +\infty} X_n(\omega) = X(\omega)\right) \\ &= 1\end{aligned}$$

Remark 278 (Interpretation). The idea is simple: for a fixed $\omega \in \Omega$, with varying n $X_n(\omega)$ is a sequence of real number (not random variables); this sequence can converge to the real number $X(\omega)$ or not.

If this is going to happen for all the elements of Ω then we met the condition.

Definition 8.1.2 (L_p convergence). Considered $p > 0$, X_n converges to X in L_p , written $X_n \xrightarrow{L_p} X$ if and only if:

1. all the X_n have moment of order p : $\mathbb{E}[|X_n|^p] < +\infty$;
2. X has moment of order p as well: $\mathbb{E}[|X|^p] < +\infty$;

3. finally $\mathbb{E}[|X_n - X|^p] \rightarrow 0$ that is

$$\lim_{n \rightarrow +\infty} \mathbb{E}[|X_n - X|^p] = 0$$

Definition 8.1.3 (Convergence in probability). X_n converges to X in probability, written $X_n \xrightarrow{p} X$, if and only if

$$\lim_{n \rightarrow +\infty} \mathbb{P}(|X_n - X| > \varepsilon) = 0, \quad \forall \varepsilon > 0$$

Definition 8.1.4 (Convergence in distribution). X_n converges to X in distribution, written $X_n \xrightarrow{d} X$, if and only if

$$\lim_{n \rightarrow +\infty} F_{X_n}(x) = F_X(x), \quad \forall x \in \mathbb{R} \text{ where } F \text{ is continuous}$$

where F_{X_n} and F_X are the distribution functions of X_n and X .

Remark 279 (Requirement of convergence in distribution). For this last case, intuitively, it would be more natural to require the convergence to hold on all the domain ($\forall x \in \mathbb{R}$, not only where F is continuous) that is

$$\lim_{n \rightarrow +\infty} F_{X_n}(x) = F_X(x), \quad \forall x \in \mathbb{R}$$

but this would be a too much severe requirement.

To understand why, suppose we have both degenerate $X_n = \frac{1}{n}$ and $X = 0$. Intuitively we would like X_n to converge to X . Here, for these degenerate, the distribution functions are:

$$F_{X_n}(x) = \begin{cases} 0 & \text{if } x < \frac{1}{n} \\ 1 & \text{if } x \geq \frac{1}{n} \end{cases}, \quad F_X(x) = \begin{cases} 0 & \text{if } x < 0 \\ 1 & \text{if } x \geq 0 \end{cases}$$

The unique discontinuity point of F is $x = 0$ where $F_X(0) = 1$ and $F_{X_n}(0) = 0$ but, if we exclude it

$$\begin{cases} \lim_{n \rightarrow +\infty} F_{X_n}(x) = 0 = F(x) & \text{if } x < 0 \\ \lim_{n \rightarrow +\infty} F_{X_n}(x) = 1 = F(x) & \text{if } x > 0 \end{cases}$$

so we can say $X_n \xrightarrow{d} X$.

However if we would require the stronger condition (convergence $\forall x \in \mathbb{R}$), it is no longer true that $\lim_{n \rightarrow \infty} F_{X_n}(x) = F(x)$ since for $x = 0$ we have

$$\lim_{n \rightarrow +\infty} F_{X_n}(0) = 0 \neq 1 = F_X(0)$$

and thus $X_n \not\xrightarrow{d} X$.

Thus if we would require $\lim_{n \rightarrow +\infty} F_{X_n}(x) = F_X(x), \forall x \in \mathbb{R}$ we would get the disturbing consequence that $X_n = \frac{1}{n}$ does not converge in distribution to $X = 0$ and this is a consequence we don't like.

Proposition 8.1.1 (Convergence in distribution of transformation). If $X_n \xrightarrow{d} X$ and $f : \mathbb{R} \rightarrow \mathbb{R}$ is a continuous function then $f(X_n) \xrightarrow{d} f(X)$

Important remark 64 (Connection among 4 modes of convergence). Summarized by the following schema (to be read as “if $X_n \xrightarrow{a.s.} X$, then $X_n \xrightarrow{p} X$ to the same X ”):

$$\begin{array}{ccccc} L_k \rightarrow & \xRightarrow{k>s} & L_s \rightarrow & & \\ & & \Downarrow & & \\ a.s. \rightarrow & \xRightarrow{} & p \rightarrow & \xRightarrow{} & d \rightarrow \end{array}$$

Finally, there’s only a special case of double implication between \xrightarrow{p} and \xrightarrow{d} in case of degenerate distribution. If X is degenerate ($X = a$ almost surely):

$$X_n \xrightarrow{p} X \iff X_n \xrightarrow{d} X$$

Important remark 65 (Some remarks on \xrightarrow{p}). We have that

1. if $X_n \xrightarrow{p} X$ and $X_n \xrightarrow{p} Y$, then $X = Y$ almost surely ($\mathbb{P}(X = Y) = 1$). In other terms the limit in probability is unique (provided it exists).

This fact has a useful consequence: suppose that we know $X_n \xrightarrow{p} X$ and we aim to prove $X_n \xrightarrow{a.s.} X$ or $X_n \xrightarrow{L_p} X$ (stronger types). In this case, the only possible limit can be X : suppose for example that $X_n \xrightarrow{a.s.} Z$ then, we have $X_n \xrightarrow{p} Z$ by the diagram, and by the previous result we have that $Z = X$ almost surely;

2. in general if $X_n \xrightarrow{p} X$ it may be that X_n fails to converge to X a.s.. However if $X_n \xrightarrow{p} X$, there is a subsequence $1 \leq n_1 < n_2 < n_3 < \dots$ such that $X_{n_j} \xrightarrow{a.s.} X$ as $j \rightarrow +\infty$.

In other terms, convergence in probability implies a.s. convergence along a suitable subsequence

Remark 280. Now some counterexamples to show that some double implications don’t work (as stated in the graph of convergence implications).

Example 8.1.1 ($X_n \xrightarrow{p} X \not\Rightarrow X_n \xrightarrow{L_p} X$ counterexample). Let X_n be such that $\mathbb{P}(X_n = 0) = \frac{n-1}{n}$ and $\mathbb{P}(X_n = n) = \frac{1}{n}$. Let $X = 0$. In this case $X_n \xrightarrow{p} X$ but $X_n \not\xrightarrow{L_p} X$, there’s convergence in probability but not in L_p (with $p = 1$):

- to prove convergence in probability we fix $\varepsilon > 0$. Then

$$\begin{aligned} \mathbb{P}(|X_n - X| > \varepsilon) &= \mathbb{P}(|X_n| > \varepsilon) \\ &\stackrel{(1)}{=} \mathbb{P}(|X_n| > \varepsilon, X_n = 0) + \mathbb{P}(|X_n| > \varepsilon, X_n = n) \\ &= 0 + \mathbb{P}(X_n = n) = \frac{1}{n} \end{aligned}$$

where in (1) $\mathbb{P}(A) = \sum_i \mathbb{P}(B_i) \cdot \mathbb{P}(A|B_i) = \sum_i \mathbb{P}(A \cap B_i)$ was applied (considered that X_n by assumption takes 2 values, 0 and n). Hence since $\frac{1}{n} \rightarrow 0$ we can state $X_n \xrightarrow{p} X$

- to show that $X_n \not\stackrel{L_1}{\rightarrow} X$ we note that

$$\begin{aligned}\mathbb{E}[|X_n - X|] &= \mathbb{E}[|X_n|] = |0| \mathbb{P}(X_n = 0) + |n| \mathbb{P}(X_n = n) \\ &= 0 + n \mathbb{P}(X_n = n) = n \frac{1}{n} = 1, \quad \forall n\end{aligned}$$

Therefore $\mathbb{E}[|X_n - X|] \not\rightarrow 0$ and thus $X_n \not\stackrel{L_1}{\rightarrow} X$

Dimostrazione. To prove the implication $X_n \xrightarrow{L_p} X \implies X_n \xrightarrow{p} X$ it suffices to use Tchebychev inequality. Suppose infact that $X_n \xrightarrow{L_p} X$: then given $\varepsilon > 0$, to have convergence in probability $\mathbb{P}(|X_n - X| > \varepsilon)$ must go to 0. Now we have that an upper bound for $\mathbb{P}(|X_n - X| > \varepsilon)$ is

$$\mathbb{P}(|X_n - X| > \varepsilon) \stackrel{(1)}{\leq} \frac{\mathbb{E}[|X_n - X|^p]}{\varepsilon^p} \stackrel{(2)}{\rightarrow} 0$$

where

- (1) due to Tchebychev
- (2) since by definition/assumption on $X_n \xrightarrow{L_p} X$

So given that the right part goes to 0, even the left part goes to 0 and saying that means that $X_n \xrightarrow{p} X$. \square

NB: new entry questo anno?

Example 8.1.2. Let $(X_n)_{n \in \mathbb{N}}$ be iid $\text{Unif}(0, 1)$ and define

$$Y_n = n \min(X_1, \dots, X_n)$$

Does Y_n converge in distribution and where?

To see it, fix x

- $\forall x > 0$ and consider distribution function

$$\mathbb{P}(Y_n \leq x) = \mathbb{P}\left(\min(X_1, \dots, X_n) \leq \frac{x}{n}\right) = 1 - \left(1 - F\left(\frac{x}{n}\right)\right)^n$$

where F is the distribution function of $X_1 \sim \text{Unif}(0, 1)$.

Hence $\mathbb{P}(Y_n \leq x) = 1 - \left(1 - \frac{x}{n}\right)^n$ converges to $1 - e^{-x}$ for $n \rightarrow \infty$.

- for $x \leq 0$ $\mathbb{P}(X_n \leq x) = 0, \forall n$

To summarize

$$\lim_{n \rightarrow +\infty} \mathbb{P}(X_n \leq x) = \begin{cases} 0 & \text{if } x \leq 0 \\ 1 - e^{-x} & \text{if } x > 0 \end{cases}$$

The latter is the distribution function of $\text{Exp}(1)$.

Hence $Y_n \xrightarrow{d} \text{Exp}(1)$

NB: Non fatto st'anno?

Example 8.1.3 (A counterexample for $X_n \xrightarrow{a.s.} X \not\implies X_n \xrightarrow{L_p} X$). As counterexample where a.s. convergence does not imply L_1 convergence considering the space:

$$(\Omega, \mathcal{A}, \mathbb{P}) = ([0, 1], \mathcal{B}([0, 1]), m)$$

with m the Lebesgue measure. In general the Lebesgue measure is *not* a probability measure, because on the real line it gives $+\infty$; but if defined on $[0, 1]$ its max is 1 so can be a probability measure.

We define also

$$X_n = n \cdot \mathbb{1}\left(\left[0, \frac{1}{n}\right]\right)\omega$$

$$X = 0$$

Here by construction we have that $\omega \in [0, 1]$; if

- $\omega \in (0, 1]$ then $\omega > \frac{1}{n}$ for large n . Therefore for large n we have that $X_n(\omega) = 0$.
- $\omega = 0$, we have that $X(0) = n\mathbb{1}([0, \frac{1}{n}])0 = n$ that goes to $+\infty$ as $n \rightarrow +\infty$.

Hence

$$\begin{aligned}\mathbb{P}(\omega \in \Omega : X_n(\omega) \rightarrow X(\omega)) &= \mathbb{P}(\omega \in \Omega : X(\omega) = 0) = \mathbb{P}(0, 1] \\ &= m(0, 1] = 1 - 0 = 1\end{aligned}$$

That is $X_n \xrightarrow{a.s.} X$.

However:

$$\begin{aligned}\mathbb{E}[|X_n - X|] &= \mathbb{E}[|X_n|] = \mathbb{E}[n \cdot \mathbb{1}([0, 1/n])\omega] = n \cdot \mathbb{E}[\mathbb{1}([0, 1/n])\omega] \\ &= n \cdot \mathbb{P}([0, 1/n]) = n \cdot m[0, 1/n] = n \cdot \frac{1}{n} \\ &= 1\end{aligned}$$

Hence here $X_n \xrightarrow{a.s.} X$ but $X_n \not\xrightarrow{L^1} X$.

Example 8.1.4. An example where convergence in distribution does not imply convergence in probability. Considering the same space:

NB: Non fatto st'anno?

$$(\Omega, \mathcal{A}, \mathbb{P}) = ([0, 1], \mathcal{B}([0, 1]), m)$$

now define $X_n = \mathbb{1}([0, 1/2])\omega$ and $X = \mathbb{1}((1/2, 1])\omega$. In this case we have that

$$|X_n - X| = 1, \quad \forall n$$

so X_n fails to converge to X in probability: $X_n \not\xrightarrow{P} X$. However the distribution functions are:

$$F(x) = \mathbb{P}(X \leq x) = \begin{cases} 0 & \text{if } x < 0 \\ \frac{1}{2} & \text{if } 0 \leq x < 1 \\ 1 & \text{if } x \geq 1 \end{cases}$$

and the other is the same:

$$F_n(x) = \mathbb{P}(X_n \leq x) = \begin{cases} 0 & \text{if } x < 0 \\ \frac{1}{2} & \text{if } 0 \leq x < 1 \\ 1 & \text{if } x \geq 1 \end{cases}$$

Hence since $F_n = F, \forall n$, we have that $X_n \xrightarrow{d} X$.

Example 8.1.5 (Esercizio dalla triennale). Si verifichi che $\frac{1}{n} \xrightarrow{d} 0$ con

- $\frac{1}{n}$ la variabile degenera X_n tale che $\mathbb{P}(X_n = \frac{1}{n}) = 1$
- 0 è la degenera X tale che $\mathbb{P}(X = 0) = 1$

Le funzioni di ripartizione sono

$$F_{X_n}(x) = \begin{cases} 1 & \text{se } x \geq \frac{1}{n} \\ 0 & \text{se } x < \frac{1}{n} \end{cases}, \quad F_X(x) = \begin{cases} 1 & \text{se } x \geq 0 \\ 0 & \text{se } x < 0 \end{cases}$$

Se

- $x < 0$ non ci sono problemi, infatti

$$F_{X_n}(x) = F(x) = 0, \quad \forall n \geq 1, \forall x < 0$$

- $x > 0$ si ha $\frac{1}{n} < x \forall n$ sufficientemente grande, quindi si ha

$$F_{X_n}(x) = 1 = F(x), \quad \forall n \text{ sufficientemente grande}, \forall x > 0$$

In definitiva $F(x) = \lim_{n \rightarrow +\infty} F_{X_n}(x), \forall x \neq 0$.

Questo basta per concludere che $1/n \xrightarrow{d} 0$. Infatti la funzione di ripartizione di X è discontinua in 0 , quindi, per avere convergenza in distribuzione, non occorre che $F_X(0) = \lim_n F_{X_n}(0)$. Ed infatti in questo esempio non è vero che $F_X(0) = \lim_n F_{X_n}(0)$ ($F_X(0) = 1$ ma $F_{X_n}(0) = 0, \forall n$)

8.2 Laws of large numbers

Remark 281. Laws (plural) because there are many of them (some of which are more famous/attractive).

Definition 8.2.1 (Sequence satisfying LLN). Let $(X_n)_{n \in \mathbb{N}} = X_1, X_2, \dots$ be a sequence of real rvs. We say it satisfies the law of large number if the sample mean

$$\bar{X}_n = \frac{1}{n} \sum_i X_i$$

converges to V for some random variable V .

Important remark 66 (Types of LLN). If

- convergence of sample mean is almost sure, $\bar{X}_n \xrightarrow{a.s.} V$, we speak of *strong law of large numbers*;
- convergence of sample mean is in probability, $\bar{X}_n \xrightarrow{p} V$, we speak of *weak law of large number*

Remark 282. Roughly speaking, any time we prove sample mean converges to a limit we have a law of large number. There are research papers that discover new large of large numbers frequently: they simply prove that a sample mean of certain sequences X_1, \dots, X_n converges to something.

In the following we'll see 3 strong laws of large numbers (SLLN) and one weak law of large numbers (WLLN).

Important remark 67. The limit V can be an arbitrary real rv; however the most important/famous (what most people think is LLN) is when the rvs are iid with existing mean, $\mathbb{E}[X_i], \forall i \in \mathbb{N}$, and V is degenerate $V = \delta_{\mathbb{E}[X_i]}$.

8.2.1 Strong laws

Theorem 8.2.1 (Kolmogorov strong law of large numbers (SLLN1)). *If $(X_n)_{n \in \mathbb{N}}$ is iid, then \bar{X}_n converges a.s. $\iff \mathbb{E}[|X_1|] < +\infty$.*

Moreover if $\mathbb{E}[|X_1|] < +\infty$, then $\bar{X}_n \xrightarrow{a.s.} \mathbb{E}[X_1]$ (that is $V = \mathbb{E}[X_1]$)

Anno mio, più sintetico: If $(X_n)_{n \in \mathbb{N}}$ is iid and $\mathbb{E}[|X_1|] < +\infty$, then $\bar{X}_n \xrightarrow{a.s.} \mathbb{E}[X_1]$.

Remark 283. Another strong law of large number follows, where we drop the iid hypothesis and replace it with some other condition.

Theorem 8.2.2 (A second example of strong LLN (SLLN2)). *Given a sequence of rvs $(X_n)_{n \in \mathbb{N}}$, if:*

- $\sup_n \mathbb{E}[X_n^2] < +\infty$ (higher second moment is still finite)
- random variables have common mean $\mathbb{E}[X_1] = \mathbb{E}[X_n], \forall n$
- $\text{Cov}(X_i, X_j) \leq 0, \forall i \neq j$

then again $\bar{X}_n \xrightarrow{a.s.} \mathbb{E}[X_1]$

Remark 284. Note that in this second version, the X_n are neither independent nor identically distributed; this version is not as popular as the first one, however it's practically very useful.

Dimostrazione. Let's prove SLLN2 (only convergence in probability, the almost sure is easier).

It suffices to apply Tchebychev inequality: given $\varepsilon > 0$ we have that

$$\begin{aligned}
 \mathbb{P}(|\bar{X}_n - \mathbb{E}[\bar{X}_n]| > \varepsilon) &\stackrel{(1)}{=} \mathbb{P}(|\bar{X}_n - \mathbb{E}[X_1]| > \varepsilon) \leq \frac{\text{Var}[\bar{X}_n]}{\varepsilon^2} = \frac{\text{Var}[\frac{1}{n} \sum_{i=1}^n X_i]}{\varepsilon^2} = \frac{\text{Var}[\sum_{i=1}^n X_i]}{n^2 \varepsilon^2} \\
 &= \frac{1}{n^2 \varepsilon^2} \left\{ \sum_{i=1}^n \text{Var}[X_i] + 2 \underbrace{\sum_{1 \leq i < j \leq n} \text{Cov}(X_i, X_j)}_{\leq 0} \right\} \\
 &\leq \frac{1}{n^2 \varepsilon^2} \sum_{i=1}^n \text{Var}[X_i] \leq \frac{1}{n^2 \varepsilon^2} \sum_{i=1}^n \mathbb{E}[X_i^2] \\
 &\stackrel{(2)}{\leq} \frac{nc}{n^2 \varepsilon^2} = \frac{c}{\varepsilon^2} \frac{1}{n} \rightarrow 0
 \end{aligned}$$

where

- in (1) before applying raw Tchebychev, we substituted applying equivalence of expected values

$$\mathbb{E}[\bar{X}_n] = \mathbb{E}\left[\frac{1}{n} \sum_{i=1}^n X_i\right] = \frac{1}{n} \sum_{i=1}^n \mathbb{E}[X_i] = \frac{n}{n} \mathbb{E}[X_1] = \mathbb{E}[X_1]$$

- c is such that $\mathbb{E}[X_n^2] \leq c, \forall n$.

This proves that $\bar{X}_n \xrightarrow{p} \mathbb{E}[X_1]$. Indeed, as claimed in the theorem, one also obtains $X_n \xrightarrow{a.s.} \mathbb{E}[X_1]$ but we will not prove almost sure convergence. \square

Remark 285. In the next example we have a strong law but the limit is not the mean. To state it we recall a definition.

Definition 8.2.2 (Stationary sequence of rv). A sequence $(X_n)_{n \in \mathbb{N}}$ is said to be *stationary* if the probability distribution of the sequence starting from two, is the same of the distribution of the unshifted sequence:

$$(X_2, X_3, X_4, \dots) \sim (X_1, X_2, X_3, \dots)$$

Hence the probability distribution of the sequence is invariant (doesn't change under shifts); in some framework this is the classical assumptions.

Theorem 8.2.3 (SLLN3). *If X_n is stationary and $\mathbb{E}[|X_1|] < +\infty$ (mean of X_1 exists), then $\bar{X}_n \xrightarrow{a.s.} V$ where V is a suitable rv (not necessarily degenerate).*

Remark 286. Two reasons why we mention the result above:

1. stationarity is an important assumption (like iid)
2. this is an example where we have a strong law (being the convergence as) but the limit is not the mean (this does not need to be the case).

8.2.2 Examples and consequences

Example 8.2.1 (A very classical example). We have an urn containing black and white balls from which we draw *with* replacement. The proportion p of white balls is not known, we want to make inference on it. Let:

$$X_i = \begin{cases} 1 & \text{if white ball drawn at trial } i \\ 0 & \text{if black ball drawn at trial } i \end{cases}$$

Since the drawing are with replacement sequence (X_i) are iid, and $\mathbb{E}[X_1] = p$, so by Kolmogorov's strong law we obtain that the sample mean converges to p , that is $\bar{X}_n \xrightarrow{a.s.} p$:

$$\frac{\text{n. palline bianche nelle prime } n \text{ prove}}{n} = \bar{X}_n \xrightarrow{a.s.} \mathbb{E}[X_1] = p$$

In other words, in this example the laws ensure us that the procedure based on intuition on making inference on p (look at extracted proportion) is going to converge to the desired quantity of interest and so if n is high, we can hope \bar{X}_n to be near p .

A question unsolved (for the moment?) is: how much big must be n in order \bar{X}_n be near to p .

Remark 287. A consequence of SLLN1 is the following

Theorem 8.2.4 (Glivenko-Cantelli thm). *If $(X_n)_{n \in \mathbb{N}}$ is iid, then*

$$\sup_t |F_n(t) - F_X(t)| \xrightarrow{a.s.} 0$$

where F_X is the distribution function common to the X_n and F_n is the so-called empirical distribution function

$$F_n(t) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}(X_i \leq t), \quad \forall t \in \mathbb{R}$$

Remark 288. This result just states that, the empirical distribution function F_n , regarded as an estimate based on the data X_1, \dots, X_n of the true/common distribution function F_X is consistent.

In fact F_n converges to F_X uniformly over t with probability 1.

The Glivenko-Cantelli thm is especially meaningful in nonparametric statistical inference

Example 8.2.2. The white black ball example 8.2.1 can be generalized as follows: let $(X_n)_{n \in \mathbb{N}}$ be iid but the distribution function F of X_1 is unknown. To make inference on F we fix a real number $x \in \mathbb{R}$ and we define the following random indicator variables

$$Y_i = \mathbb{1}(\cdot)X_i \leq x$$

Then $\{Y_i\}$ are still iid and

$$\mathbb{E}[Y_1] = \mathbb{E}[\mathbb{1}(\cdot)X_1 \leq x] = \mathbb{P}(X_1 \leq x) = F(x)$$

Hence

$$\frac{1}{n} \sum_{i=1}^n \mathbb{1}(\cdot)X_i \leq x \xrightarrow{a.s.} F(x)$$

In general:

$$F_n(x) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}(\cdot)X_i \leq x$$

is called the *empirical distribution function*, and can be regarded as an estimate of F . In fact, in statistical terms, the empirical distribution function is a *consistent* estimator of the true distribution function (that is, as the sample size goes $n \rightarrow \infty$, the procedure converges to the true value).

8.2.3 A weak law

Remark 289. Finally we state a weak law of large numbers.

Proposition 8.2.5. *If $(X_n)_{n \in \mathbb{N}}$ is iid then the following conditions are in \iff relation between them (if one is fulfilled all the others are):*

- $\bar{X}_n \xrightarrow{p} a$ for some constant a ;

NB: da mettere assieme a glivenko cantelli, è la dimostrazione

- $\varphi_{X_1}(t)$ is differentiable at 0 and $\varphi_{X_1}(0)' = ia$
- we have

$$\lim_{c \rightarrow +\infty} \mathbb{E}[X_1 \cdot \mathbb{1}(\cdot) | X_1| \leq c] = a$$

$$\lim_{c \rightarrow +\infty} c \cdot \mathbb{P}(|X_1| > c) = 0$$

Dimostrazione. Let's prove

$$\varphi_{X_1}(0)' = ia \implies \bar{X}_n \xrightarrow{p} a$$

Suppose infact $(X_n)_{n \in \mathbb{N}}$ is iid and exists the first derivative in point 0 of the characteristic function. Then the characteristic function of the sample mean is:

$$\begin{aligned} \varphi_{\bar{X}_n}(t) &= \mathbb{E}[e^{it\bar{X}_n}] = \mathbb{E}[e^{i\frac{t}{n} \sum_{i=1}^n X_i}] = \varphi_{\sum_{i=1}^n X_i}\left(\frac{t}{n}\right) \stackrel{(\text{II})}{=} \prod_{i=1}^n \varphi_{X_i}\left(\frac{t}{n}\right) \stackrel{(1)}{=} \left[\varphi_{X_1}\left(\frac{t}{n}\right)\right]^n \\ &\stackrel{(2)}{=} \left[\varphi_{X_1}(0) + \frac{t}{n}\varphi_{X_1}(0)' + o\left(\frac{t}{n}\right)\right]^n \stackrel{(3)}{=} \left[\varphi_{X_1}(0) + \frac{t}{n}ia + o\left(\frac{t}{n}\right)\right]^n \\ &= \left[1 + \frac{ita + n \cdot o\left(\frac{t}{n}\right)}{n}\right]^n \end{aligned}$$

where

- in (1) equally distributed
- in (2) we apply Taylor up to the first order
- in 3 we substituted $\varphi_{X_1}(0)' = ia$ by hypothesis

Now we use the general fact that if $z_n, z \in \mathbb{C}$ and $z_n \rightarrow z$ then

$$\left(1 + \frac{z_n}{n}\right)^n \rightarrow e^z$$

Especially considering the limit for final result

$$\lim_{n \rightarrow \infty} \left[1 + \frac{ita + n \cdot o\left(\frac{t}{n}\right)}{n}\right]^n = e^{ita}$$

In our case we used $z_n = ita + n \cdot o\left(\frac{t}{n}\right) \rightarrow iat$. We also recall that

$$o(x) \text{ as } x \rightarrow x_0 \iff \lim_{x \rightarrow x_0} \frac{o(x)}{x} = 0$$

Hence

$$no\left(\frac{t}{n}\right) = t \frac{o\left(\frac{t}{n}\right)}{\frac{t}{n}} \rightarrow 0 \text{ as } n \rightarrow +\infty$$

To summarize

$$\varphi_{\bar{X}_n}(t) \rightarrow e^{iat}, \quad \text{as } n \rightarrow +\infty$$

but e^{iat} is the characteristic function of the degenerate rv $X = a$ a.s.

By property 4 of characteristic functions we get that $\bar{X}_n \xrightarrow{d} a$, but since a degenerate rv, we also get that sample mean converges to a not only in distribution but also in probability $\bar{X}_n \xrightarrow{p} a$. \square

Example 8.2.3. This following example exhibits that in case $(X_n)_{n \in \mathbb{N}}$ is iid

- we could have $\mathbb{E}[|X_1|] = a < +\infty$ (X_1 has the mean) and convergence of sample mean to a is almost sure (not only in probability, by the strong law of Kolmogorov);
- if however $\mathbb{E}[|X_1|] = +\infty$, the strong law fails (and we don't have almost sure convergence). But it may be the case that characteristic function has first derivative at 0 (so $\exists \varphi_X(0)'$), or the other equivalent condition on limits holds, and thus weak law of large numbers holds given at least convergence in probability to a of sample mean.

Suppose that X_1 is absolutely continuous with density:

$$f(x) = \begin{cases} \frac{c}{x^2 \log|x|} & \text{if } |x| > 2 \\ 0 & \text{if } |x| \leq 2 \end{cases}$$

where c is the normalizing constant (to make integral=1). Then:

$$\begin{aligned} \mathbb{E}[|X|] &= \int_{-\infty}^{+\infty} |x| f(x) dx \stackrel{(1)}{=} 2 \int_0^{+\infty} |x| f(x) dx \stackrel{(2)}{=} 2c \int_2^{\infty} \frac{x}{x^2 \log x} dx \\ &= 2c \int_2^{+\infty} \frac{1}{x \log x} dx = +\infty \end{aligned}$$

where (1) because it's an even function and in (2) integral lower limit changes due to density

So this random variable does not have mean; hence, since (X_n) is iid the sequence \bar{X}_n does not converge almost surely and the SLLN1 does not hold.

However X_1 is symmetric (since f is an even function). The condition

$$\lim_{b \rightarrow +\infty} \mathbb{E}[X_1 \cdot \mathbb{1}(\cdot) | X_1| \leq b] = a$$

is certainly true with $a = 0$ if X_1 is symmetric (recalling that Y is symmetric if $Y \sim -Y$) in fact

$$X_1 \text{ symmetric} \implies X_1 \cdot \mathbb{1}(\cdot) | X_1| \leq b \text{ is symmetric}$$

and if a symmetric rv Y has the mean then $\mathbb{E}[Y] = 0$

$$Y \sim -Y \implies \mathbb{E}[Y] = \mathbb{E}[-Y] = -\mathbb{E}[Y] \implies \mathbb{E}[Y] = 0$$

Hence

$$\mathbb{E}[X_1 \cdot \mathbb{1}(\cdot) | X_1| \leq b] = 0, \quad \forall b > 0$$

Moreover $\forall b > 2$

$$\begin{aligned}
 b \cdot \mathbb{P}(|X_1| > b) &= b2 \mathbb{P}(X_1 > b) = 2b \int_b^{+\infty} f(x) \, dx \\
 &= 2bc \int_b^{+\infty} \frac{1}{x^2 \log x} \, dx \\
 &\leq \frac{2bc}{\log b} \int_b^{+\infty} \frac{1}{x^2} \, dx \\
 &= \frac{2bc}{\log b} \left[-\frac{1}{x} \right]_b^{+\infty} = \frac{2bc}{\log b} \frac{1}{b} \\
 &= \frac{2c}{\log b} \rightarrow 0 \quad \text{as } b \rightarrow +\infty
 \end{aligned}$$

Hence we concluded that $\bar{X}_n \xrightarrow{p} a = 0$.

8.3 Central limit theorem

8.3.1 CLT

Remark 290. Big topic of probability, one of the main findings together with law of large numbers.

Here as well, there are several CLTs (all fulfill the following general definition).

Definition 8.3.1 (Central limit theorem). A sequence $(X_n)_{n \in \mathbb{N}} = X_1, X_2, \dots$ of real random variable satisfies the CLT if there are two constants $a_n \in \mathbb{R}$ and $b_n > 0$ such that

$$\frac{\sum_{i=1}^n X_i - a_n}{b_n} \xrightarrow{d} N(0, 1)$$

Important remark 68 (CLT of sum). The sequence $(X_n)_{n \in \mathbb{N}}$ is arbitrary; one can think of sequence X_1, X_2, \dots as the sequence of observation.

To prove a CLT we need to find a_n and b_n for the ratio above to go in distribution to the standard normal.

The main/most important/natural special case is when:

$$\begin{aligned}
 a_n &= \mathbb{E} \left[\sum_{i=1}^n X_i \right] \quad \text{mean of the sum} \\
 b_n &= \sqrt{\text{Var} \left[\sum_{i=1}^n X_i \right]} \quad \text{sd of the sum}
 \end{aligned}$$

provided of course that such moment exist and $\text{Var} [\sum_{i=1}^n X_i] > 0$.

Under these choices we have that the standardization of the sum fulfill the CLT definition

Remark 291 (Natural/tipical application of CLT). Why CLT is so important in applications? Suppose we are interested in the distribution of $\sum_{i=1}^n X_i$ but we don't know how to evaluate it.

In this case, if CLT holds, we brutally replace such unknown distribution with $N(a_n, b_n^2)$. Obviously, making this replacement we make an error. However, thanks to the CLT, the error is expected to be small if n is large. In fact CLT implies that the distribution of standardized sample mean is close to standard normal

$$\frac{\sum_{i=1}^n X_i - a_n}{b_n} \sim N(0, 1)$$

Hence the distribution of the sum is close to

$$\sum_{i=1}^n X_i \sim a_n + b_n N(0, 1) = N(a_n, b_n^2)$$

If I adopt the normal for a fixed n we surely make an error, the distribution is not normal: but the distribution becomes normal as n gets larger, and the error smaller.

Remark 292. Now we start with some examples of CLT: in case LLN the most important is Kolmogorov one, similarly in CLT the main/most popular statement of this kind is the so-called CLT1.

Proposition 8.3.1 (CLT1). *If $(X_n)_{n \in \mathbb{N}}$ is sequence of iid rvs with $\mathbb{E}[X_i^2] < +\infty$ (finite second moments) and X_i is not degenerate, then the standardized sum converges in distribution to standard normal*

$$\frac{\sum_{i=1}^n X_i - \mathbb{E}[\sum_{i=1}^n X_i]}{\sqrt{\text{Var}[\sum_{i=1}^n X_i]}} = \frac{\sum_{i=1}^n X_i - n \mathbb{E}[X_i]}{\sqrt{n \text{Var}[X_i]}} \xrightarrow{d} N(0, 1)$$

Remark 293. Thus in CLT1, we have $a_n = n \mathbb{E}[X_i]$ and $b_n^2 = n \text{Var}[X_i]$. Not also that

$$\frac{\sum_{i=1}^n X_i - n \mathbb{E}[X_i]}{\sqrt{n \text{Var}[X_i]}} = \frac{\sum_{i=1}^n X_i - n \mathbb{E}[X_i]}{\sqrt{n} \sqrt{\text{Var}[X_i]}} = \frac{\sqrt{n}}{\sqrt{\text{Var}[X_i]}} (\bar{X}_n - \mu)$$

Dimostrazione. Let ϕ denote the characteristic function of the standardized single variable $\frac{X_1 - \mu}{\sigma}$ where $\mu = \mathbb{E}[X_1]$, $\sigma = \sqrt{\text{Var}[X_1]}$ (here I can divide for standard deviation cause looking at the assumption, the rv is not degenerate so the variance is positive). In the following we need the following facts where we use the fact that expected value and variance of standardized variables are 0 while second moment is 1:

$$\begin{aligned} \phi'(0) &= i \mathbb{E} \left[\frac{X_1 - \mathbb{E}[X_1]}{\sigma(X_1)} \right] = i \cdot 0 = 0 \\ \phi''(0) &= i^2 \mathbb{E} \left[\left(\frac{X_1 - \mathbb{E}[X_1]}{\sigma(X_1)} \right)^2 \right] = -1 \cdot 1 = -1 \end{aligned}$$

moment 1. The standardized sum is

$$\begin{aligned} Z_n &= \frac{\sum_{i=1}^n X_i - \mathbb{E}[\sum_{i=1}^n X_i]}{\sigma(\sum_{i=1}^n X_i)} \stackrel{(iid)}{=} \frac{\sum_{i=1}^n X_i - n \mathbb{E}[X_i]}{\sqrt{n \text{Var}[X_1]}} = \frac{\sum_{i=1}^n (X_i - \mathbb{E}[X_i])}{\sqrt{n \text{Var}[X_1]}} \\ &= \frac{1}{\sqrt{n}} \frac{\sum_{i=1}^n X_i - \mathbb{E}[X_i]}{\sigma(X_i)} \end{aligned}$$

Its characteristic function is

$$\begin{aligned}
 \varphi_{Z_n}(t) &= \varphi_{\frac{1}{\sqrt{n}} \frac{\sum_{i=1}^n X_i - \mathbb{E}[X_i]}{\sigma(X_i)}}(t) = \varphi_{\frac{\sum_{i=1}^n X_i - \mathbb{E}[X_i]}{\sigma(X_i)}}\left(\frac{t}{\sqrt{n}}\right) \\
 &\stackrel{(1)}{=} \left[\varphi_{\frac{X_i - \mathbb{E}[X_i]}{\sigma(X_i)}}\left(\frac{t}{\sqrt{n}}\right) \right]^n \\
 &\stackrel{(2)}{=} \left[\varphi(0) + \frac{t}{\sqrt{n}} \varphi'(0) + \frac{t^2}{n} \frac{1}{2} \varphi''(0) + o\left(\frac{t^2}{n}\right) \right]^n \\
 &= \left[1 + 0 - \frac{t^2/2}{n} + o\left(\frac{t^2}{n}\right) \right]^n \\
 &= \left[1 + \frac{-t^2/2 + n \cdot o\left(\frac{t^2}{n}\right)}{n} \right]^n
 \end{aligned}$$

where

- in (1) by iid: by independence the characteristic function of the sum is the product of the char function, and being identically distributed we have the power.
- in (2) as in the weak LLN proof, we use that rv by assumption have second moment finite; so we can say that the its characteristic function is C^2 and we can apply Taylor expansion (up to the the second order). So by Taylor (with Peano remainder)
- in (3) we substituted using previous result regarding single standardized variable characteristic function. since second moment exists, it exist the first as well and in the previous step we did the substitution

Now, for $n \rightarrow +\infty$ the last term developed for $\varphi_{Z_n}(t)$ converges¹ to $e^{-t^2/2}$ which is the characteristic function of a standard normal

$$\left[1 + \frac{-t^2/2 + n \cdot o\left(\frac{t^2}{n}\right)}{n} \right]^n \rightarrow e^{-t^2/2}$$

and this concludes the proof. \square

Remark 294. Let's see another version of CLT, among the several.

Proposition 8.3.2 (CLT2). *If $(X_n)_{n \in \mathbb{N}}$ are independent, with $\mathbb{E}[|X_i|^3] < +\infty$, $\mathbb{E}[X_i] = 0, \forall n$ and the following strange expression holds:*

$$\frac{\sum_{i=1}^n \mathbb{E}[|X_i - \mathbb{E}[X_i]|^3]}{(\sum_{i=1}^n \mathbb{E}[(X_i - \mathbb{E}[X_i])^2])^{\frac{3}{2}}} = \frac{\sum_{i=1}^n \mathbb{E}[|X_i|^3]}{(\sum_{i=1}^n \mathbb{E}[X_i^2])^{\frac{3}{2}}} \rightarrow 0$$

¹Again using the fact that in general if $a_n \rightarrow a$ then $(1 + \frac{a_n}{n})^n \rightarrow e^a$, and for us in it suffices to let $a_n = -\frac{t^2}{2} + n \cdot o\left(\frac{t^2}{n}\right) \rightarrow -\frac{t^2}{2}$

Then, as previously, the standardized sum converges to standard normal

$$\begin{aligned} \frac{\sum_{i=1}^n X_i - \mathbb{E}[\sum_{i=1}^n X_i]}{\sqrt{\text{Var}[\sum_{i=1}^n X_i]}} &= \frac{\sum_{i=1}^n X_i - n \mathbb{E}[X_i]}{\sqrt{n \text{Var}[X_i]}} = \frac{\sum_{i=1}^n (X_i - \mathbb{E}[X_i])}{\sigma(\sum_{i=1}^n X_i)} \\ &= \frac{\sum_{i=1}^n X_i}{\sigma(\sum_{i=1}^n X_i)} \xrightarrow{d} N(0, 1) \end{aligned}$$

Remark 295. Here the conclusion is the same as the CLT1 but the main difference is in the assumption where we are not forced to assume that rvs are identically distributed, they're just independent.

In order for the thm to still work, we need to replace that assumption with the new strange condition (don't try to attach a meaning to this condition: it's just a technical condition for the theorem to hold).

this second example is **useful because** it can be used when X_i are not identically distributed.

Important remark 69. Supponiamo che $(X_n)_{n \in \mathbb{N}}$ sia iid, con $\mathbb{E}[X_i^2] < +\infty$ e $\sigma = \text{Var}[X_i] > 0$. Se vogliamo conoscere la distribuzione di $\sum_{i=1}^n X_i$, grazie al CLT1 possiamo comunque dire che

$$\frac{\sum_{i=1}^n X_i - \mathbb{E}[\sum_{i=1}^n X_i]}{\sqrt{\text{Var}[\sum_{i=1}^n X_i]}} = \frac{\sum_{i=1}^n X_i - n\mu}{\sqrt{n\sigma^2}} \xrightarrow{d} N(0, 1)$$

con $\mu = \mathbb{E}[X_i]$ e $\sigma^2 = \text{Var}[X_i]$. Quindi per n grande la distribuzione di mio interesse di $\sum_{i=1}^n X_i$

$$\sum_{i=1}^n X_i \xrightarrow{d} N(n\mu, n\sigma^2)$$

Poi se ci interessa la media

$$\frac{\sum_{i=1}^n X_i}{n} \xrightarrow{d} \frac{1}{n} N(n\mu, n\sigma^2) = N\left(\mu, \frac{\sigma^2}{n}\right)$$

Per le proprietà di media e valore atteso

8.3.2 Examples

Example 8.3.1. Suppose $(X_n)_{n \in \mathbb{N}}$ are independent, all rvs with null mean ($\mathbb{E}[X_i] = 0, \forall i$), that variables are somehow bounded ($|X_i| \leq c, \forall i$) and $\text{Var}[\sum_{i=1}^n X_i] = \sum_{i=1}^n \text{Var}[X_i] \rightarrow +\infty$.

We are interested in convergence in distribution of the standardized sum

$$Z_n = \frac{\sum_{i=1}^n X_i}{\sigma(\sum_{i=1}^n X_i)}$$

The tool to study convergence in distribution for sum/mean is clt; in these cases we use the second version because we didn't say they are identically distributed. The random variable are independent, their mean is zero, thus in order to conclude that $Z_n \xrightarrow{d} N(0, 1)$ is enough to verify that "strange condition" holds, that is

$$\frac{\sum_{i=1}^n \mathbb{E}[|X_i|^3]}{(\sum_{i=1}^n \mathbb{E}[X_i^2])^{\frac{3}{2}}} \rightarrow 0$$

NB: Considerazione utile dalla triennale

NB: aggiuntina improvvisata mia

To answer we note that

$$|X_i|^3 = |X_i| X_i^2 \stackrel{(1)}{\leq} c X_i^2$$

with (1) by assumptions. Hence:

$$\begin{aligned} \frac{\sum_{i=1}^n \mathbb{E} [|X_i|^3]}{(\sum_{i=1}^n \mathbb{E} [X_i^2])^{\frac{3}{2}}} &\leq \frac{\sum_{i=1}^n \mathbb{E} [c \cdot X_i^2]}{(\sum_{i=1}^n \mathbb{E} [X_i^2])^{\frac{3}{2}}} = \frac{c(\sum_{i=1}^n \mathbb{E} [X_i^2])}{(\sum_{i=1}^n \mathbb{E} [X_i^2])^{\frac{3}{2}}} \\ &= \frac{c}{(\sum_{i=1}^n \mathbb{E} [X_i^2])^{\frac{1}{2}}} \rightarrow \frac{c}{\infty} = 0 \end{aligned}$$

where the denominator goes to $+\infty$ by assumption. So since the strange condition expression is upper bounded by 0, it goes to 0 as well. Hence $Z_n \xrightarrow{d} N(0, 1)$

Example 8.3.2. Suppose $(X_n)_{n \in \mathbb{N}}$ is iid with $\mathbb{E}[X_i] = 0$ and second moment $\mathbb{E}[X_i^2] = 1$, so variance $\text{Var}[X_i] = 1$. We're interested in convergence in distribution of this ratio:

$$Z_n = \frac{\sum_{i=1}^n X_i}{\sqrt{\sum_{i=1}^n X_i^2}}$$

We use CLT1 because of iid rvs. In fact Z_n can be written as (by dividing by \sqrt{n} both numerator and denominator):

$$Z_n = \frac{\frac{1}{\sqrt{n}} \sum_{i=1}^n X_i}{\sqrt{\frac{1}{n} \sum_{i=1}^n X_i^2}} = \frac{\sum_{i=1}^n X_i}{\sqrt{n}} \cdot \frac{1}{\sqrt{\frac{\sum_{i=1}^n X_i^2}{n}}}$$

Now

- by CLT1 $\frac{\sum_{i=1}^n X_i}{\sqrt{n}} \xrightarrow{d} N(0, 1)$ (think mean 0 and variance 1)
- since $(X_n)_{n \in \mathbb{N}}$ are iid $(X_n^2)_{n \in \mathbb{N}}$ are iid as well. Moreover $\mathbb{E}[X_1^2] = 1 < \infty$ by assumption. Thus Kolmogorov's SLLN applied to (X_i^2) we have that

$$\frac{\sum_{i=1}^n X_i^2}{n} \xrightarrow{a.s.} \mathbb{E}[X_i^2] = 1$$

Hence

$$Z_n = \frac{\sum_{i=1}^n X_i}{\sqrt{n}} \cdot \frac{1}{\sqrt{\frac{\sum_{i=1}^n X_i^2}{n}}} \xrightarrow{d} N(0, 1) \cdot \frac{1}{\sqrt{1}} = N(0, 1)$$

NB: sta considerazione (mio anno) per ora la lascio, si sa mai

Remark 296. In the above example as in the proof of CLT1, among other things, we used that if X_n is iid

$$\frac{\sum_{i=1}^n X_i - \mathbb{E}[\sum_{i=1}^n X_i]}{\sigma(\sum_{i=1}^n X_i)} = \frac{\sum_{i=1}^n (X_i - \mathbb{E}[X_i])}{\sqrt{n}\sigma(X_i)} = \frac{\sqrt{n}(\bar{X}_n - \mu)}{\sigma}$$

where $\sigma = \sigma(X_i)$ and $\mu = \mathbb{E}[X_i]$. In many theorem we write the quantity in that way.

Now $\sqrt{n} \rightarrow +\infty$ while $\bar{X}_n - \mu \xrightarrow{a.s.} 0$ if X_n is iid (and the moment exists).

$$\underbrace{\frac{\sqrt{n}}{\sigma}}_{\rightarrow +\infty} \cdot \underbrace{(\bar{X}_n - \mu)}_{\xrightarrow{a.s.} 0} \xrightarrow{d} N(0, 1)$$

Example 8.3.3. Suppose $(X_n)_{n \in \mathbb{N}}$ independent, $X_i \in \{-1, 0, 1\}$

$$\begin{aligned}\mathbb{P}(X_i = 1) &= \mathbb{P}(X_i = -1) = \frac{\alpha_i}{2} \\ \mathbb{P}(X_i = 0) &= 1 - \alpha_i\end{aligned}$$

$\forall i$. Let's find conditions on the constant α_i under which

$$Z_n = \frac{\sum_{i=1}^n X_i}{\sigma(\sum_{i=1}^n X_i)} \xrightarrow{d} N(0, 1)$$

These rvs can take only three values. We have that:

$$\begin{aligned}\mathbb{E}[X_i] &= 0 \cdot \mathbb{P}(X_i = 0) + 1 \cdot \mathbb{P}(X_i = 1) + (-1) \mathbb{P}(X_i = -1) = \frac{\alpha}{2} - \frac{\alpha}{2} = 0 \\ \mathbb{E}[X_i^2] &= 1 \cdot \mathbb{P}(X_i^2 = 1) + 0 \cdot \mathbb{P}(X_i^2 = 0) = \mathbb{P}(X_i = +1) + \mathbb{P}(X_i = -1) = \alpha_i \\ \text{Var}[X_i] &= \mathbb{E}[X_i^2] - (\mathbb{E}[X_i])^2 = \alpha_i\end{aligned}$$

Now since variables are independent and bounded ($|X_i| \leq c, \forall i$ if $c = 1$) by example 8.3.1 we can conclude that a sufficient condition for

$$Z_n = \frac{\sum_{i=1}^n X_i}{\sigma(\sum_{i=1}^n X_i)} \xrightarrow{d} N(0, 1)$$

provided that the sum $\sum_{i=1}^n \mathbb{E}[X_i^2] \rightarrow +\infty$. But since $\mathbb{E}[X_i^2] = \alpha_i$, we finally obtain

$$\sum_{i=1}^n \alpha_i \rightarrow +\infty \implies Z_n \xrightarrow{d} N(0, 1)$$

To prove that the converse holds, that is

$$Z_n \xrightarrow{d} N(0, 1) \implies \sum_{i=1}^n \alpha_i \rightarrow +\infty$$

Toward the contradiction suppose that $\sum_{i=1}^n \alpha_i \not\rightarrow +\infty$ that is

$$\lim_{n \rightarrow \infty} \sum_{i=1}^n \alpha_i = \alpha < +\infty$$

Note that we are summing non negative constants α_i which are the variances of the random variables; given independence we can rewrite the sum as

$$\sum_{i=1}^n \alpha_i = \sum_{i=1}^n \text{Var}[X_i] = \text{Var}\left[\sum_{i=1}^n X_i\right]$$

TODO: qui ho iniziato a improvvisare perché le note non sono esplicite

Now consider just the non-standardized sum of random variables

$$\begin{aligned}\sum_{i=1}^n X_i &= \sum_{i=1}^n X_i \cdot \frac{\sqrt{\sum_{i=1}^n \alpha_i}}{\sqrt{\sum_{i=1}^n \alpha_i}} = \sqrt{\sum_{i=1}^n \alpha_i} \cdot \frac{\sum_{i=1}^n X_i}{\sqrt{\sum_{i=1}^n \alpha_i}} \\ &= \underbrace{\sqrt{\sum_{i=1}^n \alpha_i}}_{\rightarrow \sqrt{\alpha}} \cdot \underbrace{\frac{\sum_{i=1}^n X_i}{\sqrt{\text{Var}[\sum_{i=1}^n X_i]}}}_{Z_n \xrightarrow{d} N(0, 1)} \xrightarrow{d} \sqrt{\alpha} N(0, 1) = N(0, \alpha)\end{aligned}$$

So under the assumptions above, $\sum_{i=1}^n X_i$ go in distribution to a normal. But this is a contradiction: X_i can assume only integer values (0, 1, -1) and so will be the sum $\sum_{i=1}^n X_i \in \mathbb{Z}$, while normal has domain on \mathbb{R} .

If $\sum_{i=1}^n X_i \xrightarrow{d} S$ the limit S must satisfy $\mathbb{P}(S \in \mathbb{Z}) = 1$. Instead $\mathbb{P}(N(0, \alpha) \in \mathbb{Z}) = 0$ and we have a contradiction.

Thus

$$Z_n \xrightarrow{d} N(0, 1) \implies \sum_{i=1}^n \alpha_i \rightarrow +\infty$$

Example 8.3.4. Find

$$\lim_{n \rightarrow +\infty} \mathbb{P}(\chi_n^2 > n + 7\sqrt{n})$$

where χ_n^2 denotes a chi-square with n degrees of freedom.

To evaluate such a limit, take an iid sequence $(Z_n)_{n \in \mathbb{N}}$ with $Z_i \sim N(0, 1)$. We know that $\chi_n^2 \sim \sum_{i=1}^n Z_i^2$; considered that $\mathbb{E}[Z_i] = 0$, $\text{Var}[Z_i] = 1$, moments of Z_i^2 can be obtained as:

$$\text{Var}[Z_i] = \mathbb{E}[Z_i^2] - (\mathbb{E}[Z_i])^2 = 1 \iff \mathbb{E}[Z_i^2] = 1$$

$$\text{Var}[\chi_n^2] = 2n = \text{Var}\left[\sum_{i=1}^n Z_i^2\right] \iff 2n = n \cdot \text{Var}[Z_i^2] \iff \text{Var}[Z_i^2] = 2$$

Thus $\mathbb{E}[Z_i^2] = 1$ and $\text{Var}[Z_i^2] = 2$. Now, back to the probability one obtains:

$$\begin{aligned} \mathbb{P}(\chi_n^2 > n + 7\sqrt{n}) &= \mathbb{P}\left(\sum_{i=1}^n Z_i^2 > n + 7\sqrt{n}\right) = \mathbb{P}\left(\frac{\sum_{i=1}^n Z_i^2 - n}{\sqrt{2}\sqrt{n}} > \frac{7}{\sqrt{2}}\right) \\ &= \mathbb{P}\left(\frac{\sum_{i=1}^n Z_i^2 - n \mathbb{E}[Z_i^2]}{\sqrt{n}\sqrt{\text{Var}[Z_i^2]}} > \frac{7}{\sqrt{2}}\right) = \mathbb{P}\left(\frac{\sum_{i=1}^n Z_i^2 - n \mathbb{E}[Z_i^2]}{\sqrt{\text{Var}[\sum_{i=1}^n Z_i^2]}} > \frac{7}{\sqrt{2}}\right) \end{aligned}$$

this converges by CLT1, as $n \rightarrow +\infty$ to $1 - \Phi(\frac{7}{\sqrt{2}})$ where Φ denotes the distribution function of $N(0, 1)$.

Example 8.3.5. Given any sequence $(X_n)_{n \in \mathbb{N}}$ of real rvs, the *empirical distribution function* is

$$F_n(t) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}(\cdot) X_i \leq t, \quad \forall t \in \mathbb{R}$$

If (X_n) is iid F_n can be regarded as an estimate of the distribution function common to the X_n based on the data (X_1, \dots, X_n) .

Suppose now that (X_n) is actually iid and denote by F the common distribution function of the X_n . Then

$$\sqrt{n}(F_n(t) - F(t)) \xrightarrow{d} N(0, F(t)(1 - F(t))), \forall t \in \mathbb{R} : 0 < F(t) < 1$$

Fix infatti uno such t and define $Y_n = \mathbb{1}(\cdot)X_n \leq t$. Since (X_n) is iid (Y_n) is still iid. Now some development we need after:

$$\begin{aligned}\mathbb{E}[Y_i] &= \mathbb{E}[\mathbb{1}(\cdot)X_i \leq t] = \mathbb{P}(X_i \leq t) = F(t) \\ \text{Var}[Y_i] &= \mathbb{E}[Y_i^2] - (\mathbb{E}[Y_i])^2 = \mathbb{E}[\mathbb{1}(\cdot)X_i \leq t^2] - F^2(t) \\ &\stackrel{(1)}{=} \mathbb{E}[\mathbb{1}(\cdot)X_i \leq t] - F^2(t) = F(t) - F^2(t) = F(t)(1 - F(t)) \\ \mathbb{E}\left[\sum_{i=1}^n Y_i\right] &= n \mathbb{E}[Y_i] = nF(t) \\ \text{Var}\left[\sum_{i=1}^n Y_i\right] &= n \text{Var}[Y_i] = nF(t)(1 - F(t))\end{aligned}$$

where in (1), $\mathbb{E}[\mathbb{1}(\cdot)X_i \leq t^2] = \mathbb{E}[\mathbb{1}(\cdot)X_i \leq t]$ since $\mathbb{1}(\cdot)X_i \leq t$ is an indicator. Finally, considering $\sqrt{n}[F_n(t) - F(t)]$, we can manipulate it a bit to see

$$\begin{aligned}\sqrt{n}[F_n(t) - F(t)] &= \frac{n[F_n(t) - F(t)]}{\sqrt{n}} = \frac{nF_n(t) - nF(t)}{\sqrt{n}} = \frac{n \cdot \frac{1}{n} \sum_{i=1}^n \mathbb{1}(\cdot)X_i \leq t - n \cdot \mathbb{E}[Y_i]}{\sqrt{n}} \\ &= \frac{\sum_{i=1}^n Y_i - \mathbb{E}[\sum_{i=1}^n Y_i]}{\sqrt{n}} \cdot \frac{\sqrt{\text{Var}[\sum_{i=1}^n Y_i]}}{\sqrt{\text{Var}[\sum_{i=1}^n Y_i]}} \\ &= \sqrt{n} \sqrt{F(t)(1 - F(t))} \cdot \frac{\sum_{i=1}^n Y_i - \mathbb{E}[\sum_{i=1}^n Y_i]}{\sqrt{n} \sqrt{n} \sqrt{F(t)(1 - F(t))}} \\ &= \sqrt{F(t)(1 - F(t))} \cdot \frac{\sum_{i=1}^n Y_i - \mathbb{E}[\sum_{i=1}^n Y_i]}{\sqrt{\text{Var}[\sum_{i=1}^n Y_i]}}\end{aligned}$$

by CLT1 thus this last goes to

$$\sqrt{F(t)(1 - F(t))} \cdot N(0, 1) = N(0, F(t)(1 - F(t)))$$

Example 8.3.6. Supponiamo che $Y_n \sim \text{Bin}(n, p)$; ora, poiché nel caso della binomiale si fanno estrazioni con reimmissione Y_n può essere scritta come somma di n variabili iid X_i NB: Esempio dalla triennale

$$Y_n = \sum_{i=1}^n X_i$$

dove X_i è l'indicatrice "bianca alla prova i " e

$$\begin{cases} \mathbb{P}(X_i = 1) = p \\ \mathbb{P}(X_i = 0) = 1 - p \end{cases}, \quad \forall i = 1, \dots, n$$

Quindi, ricordando che la media della binomiale è np e la varianza $np(1 - p)$, la seguente quantità standardizzata

$$\frac{Y_n - np}{\sqrt{np(1 - p)}} = \frac{\sum_{i=1}^n X_i - \mathbb{E}[\sum_{i=1}^n X_i]}{\sqrt{\text{Var}[\sum_{i=1}^n X_i]}} \xrightarrow{d} N(0, 1)$$

In sintesi se $Y_n \sim \text{Bin}(n, p)$ allora la quantità

$$\frac{Y_n - np}{\sqrt{np(1 - p)}} \xrightarrow{d} N(0, 1)$$

e quindi per n grande, la distribuzione di $\frac{Y_n - np}{\sqrt{np(1-p)}}$ può approssimarsi con una $N(0, 1)$

8.3.3 Berry-Esseen theorem

Remark 297. The most important reason for the CLT is so popular is: we are interested in the distribution of $\sum_{i=1}^n X_i$ but we are not able to evaluate it. Hence, we replace such unknown distribution with $N(a_n, b_n^2)$. If the CLT holds, namely if

$$\frac{\sum_{i=1}^n X_i - a_n}{b_n} \xrightarrow{d} N(0, 1)$$

the error we're making is small, for n large enough. Hence it is crucial to have a quantitative evaluation of the error.

Actually one of the reason of importance of CLT is the following thm which allows to evaluate the error we make in adopting the normal distribution for the sum of random variables.

Theorem 8.3.3 (Berry-Esseen Theorem). *Let's suppose condition of CLT1 plus existence of third moment holds, that is:*

- $(X_n)_{n \in \mathbb{N}}$ is iid
- X_i is non degenerate
- $\mathbb{E}[|X_i|^3] < +\infty$

Now consider the difference/error at point x :

$$\mathbb{P}\left(\frac{\sum_{i=1}^n X_i - \mathbb{E}[\sum_{i=1}^n X_i]}{\sqrt{\text{Var}[\sum_{i=1}^n X_i]}} \leq x\right) - \Phi(x)$$

where Φ is distribution function of $N(0, 1)$. By CLT1 the first term goes to standard normal $\Phi(x)$ so the above difference above goes to 0 as $n \rightarrow +\infty$. At finite n if we use standard normal we make an error, but this error is supped/bounded:

$$\sup_{x \in \mathbb{R}} \left| \mathbb{P}\left(\frac{\sum_{i=1}^n X_i - \mathbb{E}[\sum_{i=1}^n X_i]}{\sqrt{\text{Var}[\sum_{i=1}^n X_i]}} \leq x\right) - \Phi(x) \right| \leq \frac{c}{\sqrt{n}} \cdot \mathbb{E}\left[\left|\frac{X_i - \mathbb{E}[X_i]}{\sqrt{\text{Var}[X_i]}}\right|^3\right]$$

where $c \in (0, \frac{1}{2})$ (typically set it to $1/2$), and note that the error we make does not depend on considered x .

Example 8.3.7. For instance if the assumption by Berry holds and $n = 100$, we can say that the error made at any point x is

$$\leq \frac{1}{2} \frac{1}{10} \mathbb{E}\left[\left|\frac{X_i - \mathbb{E}[X_i]}{\sqrt{\text{Var}[X_i]}}\right|^3\right], \quad \forall x \in \mathbb{R}$$

Thus in practice to have a good estimate it's enough to know σ (or making assumption/educated guess).

Remark 298. One last remark on CLT: CLT1 allows to obtain some infos about speed of converge (also said convergence rate) in the Kolmogorov strong law of large numbers (the most important one). We see it below

Proposition 8.3.4. *Let's assume the condition of CLT1 and fix a sequence a_n of constants such that*

$$\frac{a_n}{\sqrt{n}} \rightarrow 0$$

Now by kolmogorov's strong law we can say that

$$\bar{X}_n - \mu \xrightarrow{a.s.} 0$$

where as before $\mu = \mathbb{E}[X_1]$. Moreover by CLT1 we have that

$$a_n(\bar{X}_n - \mu) = \frac{a_n}{\sqrt{n}} \sqrt{n}(\bar{X}_n - \mu)$$

and by assumption $\frac{a_n}{\sqrt{n}} \rightarrow 0$, while for CLT1 $\sqrt{n}(\bar{X}_n - \mu) \rightarrow N(0, \sigma^2)$ where $\sigma^2 = \text{Var}[X_i]$. Thus the product goes to 0

$$a_n(\bar{X}_n - \mu) \xrightarrow{p} 0$$

further it can be shown that one also obtains

$$a_n(\bar{X}_n - \mu) \xrightarrow{a.s.} 0$$

Remark 299. If we have only LLN we can say only $\bar{X}_n - \mu \xrightarrow{a.s.} 0$; using clt we can say much more $a_n(\bar{X}_n - \mu) \xrightarrow{a.s.} 0$.

Example 8.3.8. If I take $a_n = \sqrt{n}/\log n$ we have that

$$\frac{a_n}{\sqrt{n}} = \frac{1}{\log n} \rightarrow 0$$

and i get that

$$\frac{\sqrt{n}}{\log n}(\bar{X}_n - \mu) \xrightarrow{a.s.} 0$$

but $\sqrt{n}/\log n \rightarrow +\infty$ and

$$\frac{\sqrt{n}}{\log n}(\bar{X}_n - \mu) \rightarrow 0$$

even if $(\bar{X}_n - \mu)$ is multiplied by something that goes to $+\infty$.

8.4 Additional topics

8.4.1 Borel-Cantelli lemma

Remark 300. Let $\{A_i\}_{i \in \mathbb{N}}$ be a sequence of events, (A_i is any subset of sample space $A_i \subset \Omega$, $A_i \in \mathcal{A}$). Then we can define two new events.

Definition 8.4.1 (Limsup of the sequence). Defined as:

NB: remembering that intersection means \forall union means \exists)

$$\begin{aligned}\overline{\lim}_n A_i &= \bigcap_{n=1}^{\infty} \bigcup_{i=n}^{+\infty} A_i = \{\omega \in \Omega : \forall n \geq 1, \exists i \geq n \text{ such that } \omega \in A_i\} \\ &= \{\omega \in \Omega : \omega \in A_i \text{ for infinitely many } i\}\end{aligned}$$

Example 8.4.1. For instance if Bologna plays every sunday, A_i is Bologna wins at time i : limsup is event that Bologna wins infinite number of games.

Remark 301. Note that

- $\overline{\lim}_n A_i$ is still an event and $\overline{\lim}_n A_i \in \mathcal{A}$
- $\overline{\lim}_n A_i$ is true if and only if infinitely many of the A_i are true

NB: liminf non fatto quest'anno

Definition 8.4.2 (liminf of the sequence). Defined as

$$\underline{\lim}_n A_i = \bigcup_{n=1}^{\infty} \bigcap_{i=n}^{+\infty} A_i = \{\omega \in \Omega : \exists n \geq 1 \text{ such that } \omega \in A_i, \forall i \geq n\}$$

Example 8.4.2. Eg liminf is event there is an n such that from n on, Bologna wins every time.

Remark 302. By the Demorgan Law the complement of the limsup is the liminf of the complement, and the two events are connected by this equation

$$\left(\overline{\lim}_n A_i\right)^c = \left(\bigcap_{n=1}^{\infty} \bigcup_{i=n}^{+\infty} A_i\right)^c = \bigcup_{n=1}^{+\infty} \left(\bigcup_{i=n}^{+\infty} A_i\right)^c = \bigcup_{n=1}^{+\infty} \bigcap_{i=n}^{+\infty} A_i^c = \underline{\lim}_n A_i^c$$

Important remark 70. Borel-Cantelli lemma is a tool to evaluate the probability of the limsup $\mathbb{P}(\overline{\lim}_n A_n)$ under some assumptions.

Theorem 8.4.1 (Borel-Cantelli). *If*

- $\sum_i \mathbb{P}(A_i) < +\infty$ (that is converges) then the probability of the limsup is null: $\mathbb{P}(\overline{\lim}_n A_i) = 0$;
- $\sum_i \mathbb{P}(A_i) = +\infty$ (that is diverges) and the A_i are independent, then $\mathbb{P}(\overline{\lim}_n A_i) = 1$.

Important remark 71 (Two remarks). Regarding Borel-Cantelli:

1. Why the series of probability *necessarily* converges or diverges (can't be oscillating)? This is because it's the limit of a partial sum of positive or null numbers (probabilities).
2. if $\sum_{i=1}^n \mathbb{P}(A_n) = +\infty$ but the A_n are not independent, the Borel-Cantelli lemma does not apply (it does not cover any possible situation).

Remark 303. Proof is relatively easy but instead of it we make some examples to appreciate the use of the lemma.

Example 8.4.3. Suppose we have a coin and we throw it infinitely many times (or an urn with white and black balls from which we make drawing with replacement); we assume that the probability of tail is constant, $\mathbb{P}(T) = \alpha \in (0, 1)$, independently from the past.

Under these assumptions, for the second point of Borel-Cantelli, we observe *any* finite string of heads and tails infinitely many time with probability 1.

For example fix a finite sequence, say TTHHT; define the random variable X_i equal to indicator of the event

$$X_i = \mathbb{1}(\cdot)\text{tail at throw } i$$

We define also all non-overlapping sequences TTHHT below:

$$\begin{aligned} A_1 &= \{X_1 = 1, X_2 = 1, X_3 = 0, X_4 = 0, X_5 = 1\} \\ A_2 &= \{X_6 = 1, X_7 = 1, X_8 = 0, X_9 = 0, X_{10} = 1\} \\ A_3 &= \{X_{11} = 1, X_{12} = 1, X_{13} = 0, X_{14} = 0, X_{15} = 1\} \\ &\dots \end{aligned}$$

A_1 is the event where the string occurs at the first five trials; A_2 from trial 6 to 10 etc. Now we have that

- A_i are independent since defined them using different X_i (which are independent);
- for any A_i :

$$\mathbb{P}(A_i) = \alpha \cdot \alpha \cdot (1 - \alpha) \cdot (1 - \alpha) \cdot \alpha = \alpha^3(1 - \alpha)^2 > 0$$

Hence $\sum_{i=1}^n \mathbb{P}(A_n) = \sum_{i=1}^n \alpha^3(1 - \alpha)^2 = +\infty$ since is an infinite sum of positive constant.

Thus we met Borel-Cantelli (second point) requirements and one can conclude that

$$\mathbb{P}(\overline{\lim} A_i) = 1$$

and thus

$$\mathbb{P}(\text{observe TTHHT infinitely many times}) \geq \mathbb{P}(\overline{\lim} A_n) \stackrel{(1)}{=} 1$$

Il \geq presumo perché con gli A_i stiamo solo considerando eventi non overlappanti.

Example 8.4.4. Thanks to Borel-Cantelli it's simple to build an example where $X_n \xrightarrow{L_1} X$ but $X_n \not\xrightarrow{a.s.} X$ Take any sequence A_i of *independent* events such that $\mathbb{P}(A_i) = \frac{1}{i}$ and define

$$X_i = \mathbb{1}(\cdot)A_i = \begin{cases} 1 & \text{if } A_i \text{ is true} \\ 0 & \text{if } A_i \text{ is false} \end{cases}$$

Let $X = 0$ be degenerate. Now:

- we have that $X_i \xrightarrow{L_1} 0$ since:

$$\mathbb{E}[|X_i - X|] = \mathbb{E}[|X_i - 0|] = \mathbb{E}[|X_i|] = \mathbb{E}[\mathbb{1}(\cdot)A_i] = \mathbb{P}(A_i) = \frac{1}{i} \rightarrow 0$$

- let's see that it $X_i \xrightarrow{a.s.} X$ does not converge almost surely. Note that

- A_i are independent by assumption
- furthermore

$$\sum_i \mathbb{P}(A_i) = \sum_{i=1}^{\infty} \frac{1}{i} \stackrel{(1)}{=} +\infty$$

being (1) the harmonic series.

Hence by Borel-Cantelli we can say that $\mathbb{P}(\overline{\lim_n A_i}) = 1$.

On the other hand the complement events A_i^c

- are independent (if A_i are independent the complements are still independent)
- still

$$\sum_i \mathbb{P}(A_i^c) = \sum_i \frac{i-1}{i} = +\infty$$

Thus even here $\mathbb{P}(\overline{\lim_n A_i^c}) = 1$.

It follows that the intersection of two almost sure events is still almost sure, that is:

$$\mathbb{P}\left(\overline{\lim_n A_i} \cap \overline{\lim_n A_i^c}\right) = 1$$

Now

- fix an ω in this intersection $\omega \in (\overline{\lim_n A_i} \cap \overline{\lim_n A_i^c})$.
- since $\omega \in \overline{\lim_n A_i}$, $X_i(\omega) = \mathbb{1}(A_i)\omega = 1$ for infinitely many n
- similarly Since $\omega \in \overline{\lim_n A_i^c}$, $X_i(\omega) = \mathbb{1}(A_i)\omega = 0$ for infinitely many n

Hence $X_i(\omega)$ does not converge to any limit, so X_i does not converge almost surely.

Example 8.4.5. Let $(X_i)_{i \in \mathbb{N}}$ be iid rvs and suppose X_i is non degenerate. Under these assumption:

$$\mathbb{P}(X_i \text{ converges to a finite limit}) = 0$$

It's intuitive: if every student in a classrom choose a random number from the same distribution, the sequence will not converge to something; let's prove it formally.

Since X_i is non degenerate it can be shown (take this as given) that there are two numbers a, b with $a < b$ such that

$$\mathbb{P}(X_i \leq a) > 0 \vee \mathbb{P}(X_i \geq b) > 0$$

Now we define two events

$$\begin{aligned} A_i &= \{X_i \leq a\}, \\ B_i &= \{X_i \geq b\} \end{aligned}$$

What is the probability of limsup of A_i ? We have that

- A_i are independent (being X_i independent)
- being identically distributed we have that:

$$\sum_i \mathbb{P}(A_i) = \sum_i \mathbb{P}(X_i \leq a) \stackrel{(1)}{=} +\infty$$

with (1) because summing the same positive number infinite times

Hence $\mathbb{P}(\overline{\lim}_n A_i) = 1$.

By exactly the same arguments (B_i independent and with $\sum_i \mathbb{P}(B_i) = +\infty$) we conclude that $\mathbb{P}(\overline{\lim}_n B_i) = 1$.

Hence as before

$$\mathbb{P}\left(\overline{\lim}_n A_i \cap \overline{\lim}_n B_i\right) = 1$$

and then we fix ω in that intersection

$$\omega \in (\overline{\lim}_n A_i \cap \overline{\lim}_n B_i)$$

then $X_i(\omega)$ become a numerical sequence. Again this sequence does not converge:

- since $\omega \in \overline{\lim}_n A_i$, then $X_i(\omega) = X_i \leq a$ for infinitely many n
- otoh since $\omega \in \overline{\lim}_n B_i$, then $X_i(\omega) = X_i \geq b$ for infinitely many n

So having that $a < b$ the sequence can't converge to any limit and formally

$$\mathbb{P}(\omega \in \Omega : X_i(\omega) \text{ does not converge}) = 1$$

Remark 304. An incidentally (related to Borel-Cantelli) useful fact is the following

Important remark 72. We have

- recall that for any sequence $(A_i)_{i \in \mathbb{N}}$

$$\mathbb{P}(\cup_i A_i) \leq \sum_i \mathbb{P}(A_i)$$

- in particular, if $\mathbb{P}(A_i) = 0, \forall i$ then

$$\mathbb{P}(\cup_i A_i) \leq \sum_i \mathbb{P}(A_i) = 0 \tag{8.1}$$

- if $\mathbb{P}(A_i) = 1, \forall i$ then then the

$$\mathbb{P}(\cap_i A_i) = 1$$

In fact:

$$\mathbb{P}(\cap_i A_i) = 1 - \mathbb{P}((\cap_i A_i)^c) = 1 - \mathbb{P}(\cup_i A_i^c) \stackrel{(1)}{=} 1 - 0 = 1$$

where (1) by applying the previous one, eq 8.1, since $\mathbb{P}(A_i^c) = 0, \forall i$

8.4.2 Stable rvs

Remark 305. It's an important type of random variables, together with infinite divisible rvs

Definition 8.4.3 (Stable rv). A real rv X is said to be stable if exist an iid rvs sequence $\{X_i\}_{i \in \mathbb{N}}$ (with $X_i \sim Z$) and real constant a_n and $b_n > 0$ (actually real sequences i guess) such that

$$\frac{\sum_{i=1}^n X_i - a_n}{b_n} \xrightarrow{d} X$$

Example 8.4.6. $X \sim N(0, 1)$ is stable rv. In fact by CLT1, if Z is any non-degenerate rv with $\mathbb{E}[Z^2] < +\infty$ and if $(X_i)_{i \in \mathbb{N}}$ is iid with $X_i \sim Z$ then

$$\frac{\sum_{i=1}^n X_i - n \mathbb{E}[Z]}{\sqrt{n} \sqrt{\text{Var}[Z]}} \xrightarrow{d} N(0, 1)$$

Hence it suffices to let $a_n = n \mathbb{E}[Z]$ and $b_n = \sqrt{n} \sqrt{\text{Var}[Z]}$.

Remark 306. In general, the stable rvs are those rv which, like the normal can be obtained as the limit in distribution of the partial sums $\sum_{i=1}^n X_i$ (suitably normalized) of iid rvs.

NB: Qui la particolarità rispetto la definizione è $X_i \sim X$ e inoltre nello statement vi è direttamente \sim , non \xrightarrow{d}

Theorem 8.4.2 (Characterization). X is stable \iff considering the sequence $\{X_i\}_{i \in \mathbb{N}}$ iid with $X_i \sim X$, $\forall n \geq 1$ there are real constants $\alpha_n \in \mathbb{R}$ and $\beta_n > 0$ such that:

$$\frac{\sum_{i=1}^n X_i - \alpha_n}{\beta_n} \sim X$$

Remark 307. The idea of this theorem: given any rv X , take X_1, \dots, X_n iid with the same distribution as X , $X_i \sim X$. Then, in general, $\sum_{i=1}^n X_i \approx X$. However, if X is stable we can find constants α_n, β_n such that the normalized partial sum has the same distribution X of the summed variables:

$$\frac{\sum_{i=1}^n X_i - \alpha_n}{\beta_n} \sim X$$

Example 8.4.7. By applying CLT1 we found $N(0, 1)$ is stable. We can show even with characterization theorem: in fact, if $(X_i)_{i \in \mathbb{N}}$ are iid with $X_i \sim N(0, 1)$ then

$$\sum_{i=1}^n X_i \sim N(0, n)$$

so that

$$\frac{\sum_{i=1}^n X_i}{\sqrt{n}} \sim \frac{1}{\sqrt{n}} N(0, n) = N(0, 1)$$

Hence the previous characterization applies with $\alpha_n = 0$ and $\beta_n = \sqrt{n}$

Remark 308. Other example of stable rvs are the Cauchy and degenerate.

Example 8.4.8 (Cauchy). If X is Cauchy then the characteristic function of X is (take it as given)

$$\varphi_X(t) = e^{-|t|}, \quad \forall t \in \mathbb{R}$$

Now we have to verify the definition. Let's have $\{X_i\}_{i \in \mathbb{N}}$, with $X_i \sim \text{Ca}(\cdot)$, and verify that, if we set $a_n = 0$, $b_n = n$ (obtaining the sample mean):

$$\frac{\sum_{i=1}^n X_i - a_n}{b_n} = \frac{\sum_{i=1}^n X_i}{n} = \bar{X}_n$$

the sample mean is distributed as Cauchy as well; we do this by checking its characteristic function:

$$\begin{aligned} \varphi_{\bar{X}_n}(t) &= \varphi_{\frac{\sum_{i=1}^n X_i}{n}}(t) = \varphi_{\sum_{i=1}^n X_i}\left(\frac{t}{n}\right) \stackrel{(iid)}{=} \left[\varphi_{X_i}\left(\frac{t}{n}\right)\right]^n = \left[e^{-|\frac{t}{n}|}\right]^n \\ &= e^{-|t|} \end{aligned}$$

which is still Cauchy. Hence by taking $a_n = 0$, $b_n = n$ the quantity

$$\frac{\sum_{i=1}^n X_i - a_n}{b_n} = \bar{X}_n \sim Y_1$$

is Cauchy distributed and so by the characterization theorem Cauchy is example of stable rv.

Example 8.4.9 (Degenerate). Another stable rv; if $\mathbb{P}(X = a) = 1$ for some constant a then trivially

$$\frac{a \cdot n}{n} = a$$

thus

$$\frac{\sum_{i=1}^n X_i - \alpha_n}{\beta_n} \sim X$$

with $X_i \sim \delta_a$, $\alpha_n = 0$ and $\beta_n = n$.

Example 8.4.10. Let $\alpha \in (0, 2]$ and

$$\phi(t) = \exp\left(-\frac{|t|^\alpha}{2}\right), \quad \forall t \in \mathbb{R}$$

It can be shown that ϕ is the characteristic function of some rv X . Such X is actually stable.

In fact if X_1, \dots, X_n are iid with $X_i \sim X$ then

$$\begin{aligned} \varphi_{\frac{\sum_{i=1}^n X_i}{n^{1/\alpha}}}(t) &= \mathbb{E}\left[e^{i \frac{t}{n^{1/\alpha}} \sum_{i=1}^n X_i}\right] = \varphi_{\sum_{i=1}^n X_i}\left(\frac{t}{n^{1/\alpha}}\right) = \left[\varphi_{X_i}\left(\frac{t}{n^{1/\alpha}}\right)\right]^n \\ &= \left(\exp\left[-\frac{1}{2} \left|\frac{t}{n^{1/\alpha}}\right|^\alpha\right]\right)^n = \left(\exp\left[-\frac{1}{2} \frac{|t|^\alpha}{n}\right]\right)^n \\ &= \exp\left(-\frac{1}{2} |t|^\alpha\right) \end{aligned}$$

This proves that

$$\frac{\sum_{i=1}^n X_i}{n^{1/\alpha}} \sim X$$

so that once again the previous characterization applies with $\alpha_n = 0$ and $\beta_n = n^{1/\alpha}$

NB: bo secondo me qua ha sbagliato

8.4.3 Infinite divisible rvs

Remark 309. The infinite divisible rvs are a remarkable subclass of the stable rvs. Loosely speaking, a distribution is infinitely divisible if it can be expressed as the sum of an arbitrary number of independent and identically distributed (i.i.d.) random variables.

Definition 8.4.4 (Infinite divisible rv). A real rv X is infinite divisible if and only if $\forall n \geq 1$, there are Y_1, \dots, Y_n iid rvs such that $\sum_{i=1}^n Y_i \sim X$.

Remark 310. Important examples of infinitely divisible rvs are: the normal, the Cauchy, the degenerate, the Poisson, and the Gamma.

Proposition 8.4.3. *If X is stable then X is infinite divisible, but the viceversa does not hold. (so stable are a proper subset of infinite divisible)*

Dimostrazione. Infact

- if X is stable, we can write

$$X \sim \frac{\sum_{i=1}^n X_i - \alpha_n}{\beta_n}$$

where X_1, \dots, X_n are iid. Hence letting

$$Y_i = \frac{X_i - \alpha_n/n}{\beta_n}, \quad i = 1, \dots, n$$

one obtains Y_1, \dots, Y_n iid and

$$\sum_{i=1}^n Y_i = \frac{\sum_{i=1}^n X_i - \alpha_n}{\beta_n} \sim X$$

- viceversa does not hold: by counterexample we need a infinite divisible which is not stable.

It is sufficient to note that the only stable random variable X with finite variance/second moment ($\mathbb{E}[X^2] < \infty$) are the normal $N(\mu, \sigma^2)$ and the degenerate.

Hence if X is infinitely divisible, with finite variance, but neither degenerate nor normal, then X is an example of infinite divisible but not stable rv.

Example of infinite divisible but not stable are the exponential (or the poisson): the exponential has the second moment but it is neither normal nor degenerate thus it's not stable; however as we noted before is infinite divisible.

□

Remark 311. The following result characterizes the infinite divisible rv having finite second moment.

Theorem 8.4.4. *We have that*

1. X is infinite divisible and

2. $\mathbb{E}[X^2] < +\infty$ (has finite second moment)

$\iff X \sim X_1 + X_2 + X_3$ with X_1, X_2, X_3 independent, X_1 degenerate, $X_2 \sim N(0, \sigma^2)$ and X_3 generalized Poisson.

Remark 312. So this result describes the structure of infinite divisible random variables (with finite second moment).

In other words if X is infinitely divisible and $\mathbb{E}[X^2] < +\infty$ then X has the same distribution of the sum of 3 independent rvs, such that one is degenerate, the other is $N(0, \sigma^2)$ and the third is generalized Poisson. Let's see what is a generalized Poisson.

Definition 8.4.5 (Generalized poisson). X is generalized poisson if

$$X \sim \mathbb{1}(\cdot)N > 0 \cdot \sum_{i=1}^N Z_i$$

where:

- $N \sim \text{Pois}(\lambda)$
- (Z_i) is any iid sequence of rvs
- $N \perp\!\!\!\perp (Z_i)$

Important remark 73. We expect to find the poisson rv as a special case of this. Infact, if $Z_i = 1, \forall i$:

$$X \sim \mathbb{1}(\cdot)N > 0 \cdot N = N$$

but $N \sim \text{Pois}(\lambda)$. So as expected the poisson is just special case of the generalized Poisson.

Note also that if φ denotes the characteristic function common to the Z_i the characteristic function of X can be written as

$$\begin{aligned} \varphi_X(t) &= \mathbb{E}[e^{itX}] = \mathbb{E}\left[\mathbb{1}(\cdot)N = 0 \cdot 1 + \sum_{n=1}^{+\infty} \mathbb{1}(\cdot)N = ne^{it \sum_{i=1}^n Z_i}\right] \\ &= \mathbb{P}(N=0) + \sum_{n=1}^{+\infty} \mathbb{P}(N=n) \varphi_{\sum_{i=1}^n Z_i}(t) \\ &= \mathbb{P}(N=0) + \sum_{n=1}^{+\infty} \mathbb{P}(N=n) [\varphi_{Z_i}(t)]^n \\ &= \sum_{n=0}^{+\infty} \mathbb{P}(N=n) [\varphi_{Z_i}(t)]^n \\ &= \frac{e^{-\lambda} \lambda^n}{n!} [\varphi_{Z_i}(t)]^n \\ &= e^{-\lambda} \underbrace{\sum_{n=0}^{+\infty} \frac{(\lambda \varphi_{Z_i}(t))^n}{n!}}_{e^{\lambda \varphi_{Z_i}(t)}} \\ &= e^{-\lambda} e^{\lambda \varphi_{Z_i}(t)} = e^{\lambda(\varphi_{Z_i}(t)-1)} \end{aligned}$$

This latter is the characteristic function of a generalized poisson rv.

Theorem 8.4.5. *If X is infinite divisible and $\mathbb{P}(a \leq X \leq b) = 1$ for some a and b (X is bounded) then X is degenerate.*

Dimostrazione. Since X is infinite divisible $\forall n \geq 1$ we have $X \sim \sum_{i=1}^n X_{n_i}$, with X_{n_1}, \dots, X_{n_n} iid. Thus:

$$\text{Var}[X] = \text{Var}\left[\sum_{i=1}^n X_{n_i}\right] = \sum_{i=1}^n \text{Var}[X_{n_i}] = n \text{Var}[X_{n_i}] \leq n \mathbb{E}[X_{n_i}^2]$$

with last passage due to variance calculation formula.

Now since $\mathbb{P}(a \leq X \leq b) = 1$, we have that

- $\mathbb{P}(X_{n_i} > \frac{b}{n}) = 0$. Infact

$$0 = \mathbb{P}(X > b) = \mathbb{P}\left(\sum_{i=1}^n X_{n_i} > b\right) \geq \mathbb{P}\left(X_{n_i} > \frac{b}{n}, \forall i\right) \stackrel{(iid)}{=} \left[\mathbb{P}\left(X_{n_i} > \frac{b}{n}\right)\right]^n$$

and thus $\mathbb{P}(X_{n_i} > \frac{b}{n}) = 0$

- $\mathbb{P}(X_{n_i} < \frac{a}{n}) = 0$ by the same argument

Thus X_{n_i} stays between $\frac{a}{n}$ and $\frac{b}{n}$ almost surely

$$\mathbb{P}\left(\frac{a}{n} \leq X_{n_i} \leq \frac{b}{n}\right) = 1$$

and therefore

$$\mathbb{P}\left(|X_{n_i}| \leq \frac{\max(|a|, |b|)}{n}\right) = 1$$

Given this last equation, if we square both terms we can write that the expected value (of the squared rv) is less than the squared “domain superior limit”:

$$\mathbb{E}[X_{n_i}^2] \leq \frac{\max(|a|, |b|)^2}{n^2}$$

Hence.

$$\text{Var}[X] \leq n \mathbb{E}[X_{n_i}^2] \leq n \frac{\max(|a|, |b|)^2}{n^2} = \frac{\max(|a|, |b|)^2}{n}$$

and thus (being valid $\forall n \geq 1$ even an high value I guess)

$$\lim_{n \rightarrow +\infty} \frac{\max(|a|, |b|)^2}{n} = 0$$

thus X is degenerate. □

8.4.4 Examples

NB: non fatto quest'anno

Example 8.4.11 (Poisson). $X \sim \text{Pois}(\lambda)$ is infinite divisible. Infact, if Y_1, \dots, Y_n are independent and $Y_i \sim \text{Pois}(\lambda_i)$ then $\sum_{i=1}^n Y_i \sim \text{Pois}(\sum_{i=1}^n \lambda_i)$. Hence if $X \sim \text{Pois}(\lambda)$, it is sufficient to take X_{n_1}, \dots, X_{n_n} iid rvs with $X_{n_i} \sim \text{Pois}(\frac{\lambda}{n})$

non fatto quest'anno **Example 8.4.12** (Normal). $N(\mu, \sigma^2)$ is infinite divisible. Infact if X_1, \dots, X_n independent, $N(\mu_i, \sigma_i^2)$, then $\sum_{i=1}^n X_i \sim N(\sum_{i=1}^n \mu_i, \sum_{i=1}^n \sigma_i^2)$

non fatto quest'anno **Example 8.4.13** (Gamma). Another example is the gamma. In fact $X \sim \text{Gamma}(\alpha, \beta)$ iff X is absolutely continuous with density

$$f(x) = \begin{cases} \frac{\alpha^\beta}{\Gamma(\beta)} e^{-\alpha x} x^{\beta-1} & x > 0 \\ 0 & \text{elsewhere} \end{cases}$$

Note for $\beta = 1$ we get the $\text{Gamma}(\alpha, 1) = \text{Exp}(\alpha)$ so exponential is a special case of gamma.

Now if Y_1, \dots, Y_n indep and $Y_i \sim \text{Gamma}(\alpha, \beta_i)$ (with common α) then the sum of Y_i is still a gamma, that is $\sum_{i=1}^n Y_i \sim \text{Gamma}(\alpha, \sum_{i=1}^n \beta_i)$.

By the way, if Y_1, \dots, Y_n are iid $Y_i \sim \text{Exp}(\alpha)$, then the distribution of the sum $\sum_{i=1}^n Y_i = \text{Gamma}(\alpha, n)$ (n because $\beta = 1$ and the sum is n).

Using the above results it follows that Gamma is infinite divisible.

Capitolo 9

Convergence

Important remark 74 (Setup). Given a sequence of rvs, X_1, X_2, \dots , the aim is to study

$$(X_n)_{n \in \mathbb{N}} \xrightarrow[n \rightarrow \infty]{} X$$

We have four types of convergence:

1. convergence in probability (weak)
2. convergence in law/distribution (weak)
3. convergence in mean of order k (strong)
4. almost sure convergence (strong)

9.1 Convergence in probability

9.1.1 Definition

Remark 313. It's the first type of convergence: this is a weak type (it implies convergency in distribution but not stronger kinds of convergency)

Definition 9.1.1 (Convergence in probability). We say that a sequence $(X_n)_{n \in \mathbb{N}}$ converges in probability to the *limit distribution* X and we write:

$$(X_n)_{n \in \mathbb{N}} \xrightarrow[n \rightarrow \infty]{p} X$$

if alternatively (equivalent definitions), $\forall \varepsilon > 0$:

$$\mathbb{P}(|X_n - X| > \varepsilon) \xrightarrow[n \rightarrow \infty]{} 0 \quad (9.1)$$

$$\mathbb{P}(|X_n - X| < \varepsilon) \xrightarrow[n \rightarrow \infty]{} 1 \quad (9.2)$$

Remark 314. The limit distribution X can be any rv (gaussian etc) but as a special case it's when X_n converges to a δ_θ (the constant θ); it's peculiar since in inference the sequence can be an estimator collassing to a point (eg population mean) and can be a good property for an estimator.

9.1.2 Weak consistence

Definition 9.1.2 (Weak consistence). If $(X_n)_{n \in \mathbb{N}} \xrightarrow[n \rightarrow \infty]{p} \delta_\theta$ we say that X_n is (weakly) consistent for θ

Important remark 75. Weak consistency means converging probability (link between probability and inference)

Example 9.1.1. Considering a sequence of iid rvs $(X_n)_{n \in \mathbb{N}} \sim \text{Unif}(0, \theta)$, with $\theta > 0$, the transformation (max of the first n)

$$\max_{1 \leq i \leq n} X_i = X_{(n)}$$

Let's prove that $X_{(n)}$ is a consistent estimator for θ , that is:

$$X_{(n)} \xrightarrow{p} \delta_\theta$$

Remembering that $F_{(n)}(x) = [F_X(x)]^n$ we want to prove that

$$\mathbb{P}(|X_{(n)} - \theta| < \varepsilon) \rightarrow 1$$

Now we have:

$$\begin{aligned} \mathbb{P}(|X_{(n)} - \theta| < \varepsilon) &\stackrel{(1)}{=} \mathbb{P}(-X_{(n)} + \theta < \varepsilon) = \mathbb{P}(-X_{(n)} < \varepsilon - \theta) = \mathbb{P}(X_{(n)} > \theta - \varepsilon) \\ &= 1 - \mathbb{P}(X_{(n)} \leq \theta - \varepsilon) = 1 - F_{(n)}(\theta - \varepsilon) \\ &= 1 - [F_X(\theta - \varepsilon)]^n \end{aligned}$$

where in (1) since $X_{(n)} - \theta$ is negative or null (being θ the max of the uniform rvs) we can avoid the absolute value multiplying by -1 .

If $X \sim \text{Unif}(0, \theta)$, then $F_X(x) = \frac{x}{\theta}$, $0 \leq x \leq \theta$ so

$$\mathbb{P}(|X_{(n)} - \theta| < \varepsilon) = 1 - [F_X(\theta - \varepsilon)]^n = 1 - \left(\frac{\theta - \varepsilon}{\theta}\right)^n$$

and since $\frac{\theta - \varepsilon}{\theta} < 1$ with $0 < \varepsilon \leq \theta$

$$\lim_{n \rightarrow \infty} 1 - \left(\frac{\theta - \varepsilon}{\theta}\right)^n = 1$$

Proposition 9.1.1 (Sufficient conditions for weak consistence). *If*

$$\begin{cases} \lim_{n \rightarrow +\infty} \mathbb{E}[X_n] = \theta \\ \lim_{n \rightarrow +\infty} \text{Var}[X_n] = 0 \end{cases} \implies X_n \xrightarrow{p} \delta_\theta \quad (9.3)$$

Remark 315. The viceversa does not hold: eg X_n can converge in probability even if these conditions are not met.

Dimostrazione. Applying Tchebychev inequality

$$\mathbb{P}(|X_n - \mathbb{E}[X_n]| < \lambda \sigma(X_m)) \geq 1 - \frac{1}{\lambda^2}$$

Now we define/substitute $\varepsilon = \lambda\sigma(X_m)$ so that $\lambda^2 = \frac{\varepsilon^2}{\sigma^2(X_m)}$; therefore

$$\mathbb{P}(|X_n - \mathbb{E}[X_n]| < \varepsilon) \geq 1 - \frac{\sigma^2(X_n)}{\varepsilon^2}$$

if $n \rightarrow +\infty$ the last term go to zero so

$$\mathbb{P}(|X_n - \mathbb{E}[X_n]| < \varepsilon) \geq 1$$

and since this probability can't be larger than 1, it must be 1 so

$$\mathbb{P}(|X_n - \mathbb{E}[X_n]| < \varepsilon) = 1 \implies X_n \xrightarrow{P} \theta$$

□

Example 9.1.2. Let $X_n \sim \text{Geom}(p_n)$ with $p_n = 1 - \frac{1}{n}$, having pmf

$$\mathbb{P}(X_n = x) = p_n(1 - p_n)^{x-1}$$

with $\mathbb{E}[X_n] = \frac{1}{p_n}$, $\text{Var}[X_n] = \frac{1-p_n}{p_n^2}$. Let's prove that $X_n \xrightarrow{P} \delta_1$.

$$\begin{aligned} \lim_{n \rightarrow \infty} \mathbb{E}[X_n] &= \frac{1}{p_n} = \frac{1}{1 - \frac{1}{n}} \rightarrow 1 \\ \lim_{n \rightarrow \infty} \text{Var}[X_n] &= \frac{1 - p_n}{p_n^2} = \frac{1 - (1 - \frac{1}{n})}{(1 - \frac{1}{n})^2} = \frac{\frac{1}{n}}{(1 - \frac{1}{n})^2} \\ &= \frac{\frac{1}{n}}{(\frac{n-1}{n})^2} = \frac{n}{(n-1)^2} \rightarrow 0 \end{aligned}$$

Example 9.1.3 (Esame vecchio viroli). Let θ be the parameter of a population random variable X that follows a continuous uniform distribution on the interval $[\theta - 2, \theta + 1]$ and let $X = (X_1, \dots, X_n)$ be a simple random sample. Given the estimator $T_n(X) = \bar{X}_n + \frac{1}{2}$ decide if it is weakly consistent. We need

$$\begin{aligned} \mathbb{E}[T_n(X)] &= \mathbb{E}\left[\bar{X}_n + \frac{1}{2}\right] = \mathbb{E}[\bar{X}_n] + \frac{1}{2} = \mathbb{E}\left[\frac{X_1 + \dots + X_n}{n}\right] + \frac{1}{2} = \frac{1}{n} \mathbb{E}\left[\sum_{i=1}^n X_i\right] + \frac{1}{2} \\ &= \frac{1}{n} \cdot n \cdot \mathbb{E}[X_i] + \frac{1}{2} = \frac{\theta - 2 + \theta + 1}{2} + \frac{1}{2} = \frac{2\theta}{2} = \theta \\ \text{Var}[T_n(X)] &= \text{Var}\left[\bar{X}_n + \frac{1}{2}\right] = \text{Var}[\bar{X}_n] = \text{Var}\left[\frac{X_1 + \dots + X_n}{n}\right] = \frac{1}{n^2} \text{Var}\left[\sum_{i=1}^n X_i\right] \\ &= \frac{1}{n^2} \cdot n \cdot \text{Var}[X_i] = \frac{\text{Var}[X_i]}{n} = \frac{1}{12n}(\theta + 1 - \theta + 2)^2 = \frac{9}{12n} = \frac{3}{4n} \rightarrow 0 \end{aligned}$$

Therefore $T(X)$ is weakly consistent

Example 9.1.4 (Esame vecchio viroli). Let X_n be a sequence of iid exponential random variables with parameter 1. Study the convergence in probability of the minimum $X_{(1)}$.

The minimum of an exponential should converge to the minimum of the domain

so for the exponential is 0. We check the two sufficient condition but first let write the density function of the minimum

$$f_{X_1}(x) = n \cdot f_X(x) \cdot [1 - F_X(x)]^{n-1}$$

where for the Exp(1) we have

$$\begin{aligned} f(x) &= e^{-x} \\ F_X(x) &= 1 - e^{-x} \end{aligned}$$

and therefore

$$f_{X_{(1)}}(x) = n \cdot e^{-x} \cdot (1 - 1 + e^{-x})^{n-1} = n \cdot e^{-x(n-1)-x} = n \cdot e^{-nx}$$

We have

$$\mathbb{E}[X_{(1)}] = \int_0^{+\infty} x \cdot n \cdot e^{-nx} = - \int_0^{+\infty} x \cdot (-n) \cdot e^{-nx}$$

Sviluppiamo l'integrale indefinito e poi valutiamolo

$$\begin{aligned} - \int x(-n)e^{-nx} &= - \left[e^{-nx}x - \int e^{-nx} \right] = - \left[e^{-nx}x + \frac{1}{n} \int (-n)e^{-nx} \right] \\ &= - \left[e^{-nx}x + \frac{1}{n}e^{-nx} \right] = -e^{-nx} \left(x + \frac{1}{n} \right) \end{aligned}$$

Che valutato

$$\left[-e^{-nx} \left(x + \frac{1}{n} \right) \right]_0^{+\infty} = \frac{1}{n}$$

Per cui $\mathbb{E}[X_{(1)}] \rightarrow 0$.

Per la varianza calcoliamo il secondo momento

$$\mathbb{E}[X_{(1)}^2] = \int_0^{+\infty} x^2 \cdot n \cdot e^{-nx} = \dots = \frac{2}{n^2}$$

Per cui

$$\text{Var}[X_{(1)}] = \frac{2}{n^2} - \frac{1}{n} \rightarrow 0$$

Answer: $X_{(1)} \xrightarrow{p} 0$

Example 9.1.5 (Esame vecchio viroli). Let (X_1, \dots, X_n) a simple random sample from an exponential random variable

$$f_X(x) = \theta e^{-\theta x}$$

Study the convergence in probability of

$$T_n = 2 \frac{\sum_{i=1}^n X_i}{n} + 3$$

1. T_n converges in probability to a dirac at $(3 + \theta)/2$

2. T_n converges to a dirac at 2θ
3. T_n does not converge in probability
4. T_n converge to a dirac at $\frac{2}{\theta} + 3$ (dovrebbe essere questa)

For the exponential we have that $\mathbb{E}[X_i] = \frac{1}{\theta}$ and $\text{Var}[X_i] = \frac{1}{\theta^2}$. We check the two sufficient condition

$$\begin{aligned}\mathbb{E}[T_n] &= \mathbb{E}\left[2\frac{\sum_{i=1}^n X_i}{n} + 3\right] = 2\mathbb{E}\left[\frac{\sum_{i=1}^n X_i}{n}\right] + 3 = \frac{2}{n}n\mathbb{E}[X_i] + 3 \\ &= \frac{2}{\theta} + 3 \\ \text{Var}[T_n] &= \text{Var}\left[2\frac{\sum_{i=1}^n X_i}{n} + 3\right] = \frac{4}{n^2}\text{Var}\left[\sum_{i=1}^n X_i\right] = \frac{4}{n^2}n\text{Var}[X_i] \\ &= \frac{4}{n}\frac{1}{\lambda^2} \rightarrow 0\end{aligned}$$

So $T_n \xrightarrow{p} \delta_{\frac{2}{\theta}+3}$

9.1.3 Theorem: weak law of large numbers

Theorem 9.1.2 (Weak law of large numbers). *Let X_n be a sequence of iid rvs with $\mathbb{E}[X_n] = \theta$ and $\text{Var}[X_n] = \sigma^2 < +\infty$; if we define the partial mean as the mean of the first n rvs*

$$M_n = \frac{\sum_{i=1}^n X_i}{n} \quad (9.4)$$

then we have that

$$M_n \xrightarrow{p} \delta_\theta \quad (9.5)$$

Dimostrazione. We have that

$$\begin{aligned}\mathbb{E}[M_n] &= \frac{\sum_{i=1}^n \mathbb{E}[X_i]}{n} = \frac{n\theta}{n} = \theta \\ \text{Var}[M_n] &= \frac{1}{n^2} \sum_{i=1}^n \text{Var}[X_i] = \frac{n}{n^2}\sigma^2 = \frac{\sigma^2}{n}\end{aligned}$$

therefore since both

$$\begin{aligned}\lim_{n \rightarrow +\infty} \mathbb{E}[M_n] &= \theta \\ \lim_{n \rightarrow +\infty} \text{Var}[M_n] &= 0\end{aligned}$$

the sufficient conditions are met and $M_n \xrightarrow{p} \delta_\theta$ □

Example 9.1.6. Let X_1, \dots, X_n be independent rvs each distributed as Bernoulli with parameter p . Prove that $\frac{1}{n} \sum_{i=1}^n X_i^2 \xrightarrow{p} p$ as $n \rightarrow \infty$. According to the WLLN we have that

$$\frac{1}{n} \sum_{i=1}^n X_i^2 \xrightarrow{p} \mathbb{E}[X^2]$$

Now if $X \sim \text{Bern}(p)$, then $\mathbb{E}[X] = p$ and $\text{Var}[X] = p(1-p)$, so $\mathbb{E}[X^2] = \text{Var}[X] + \mathbb{E}[X]^2 = p(1-p) + p^2 = p$. Therefore

$$\frac{1}{n} \sum_{i=1}^n X_i^2 \xrightarrow{p} p$$

9.2 Convergence in law/distribution

Remark 316. We have two equivalent definition, by limit of distribution function or convergence in law/distribution of the moment generating function.

Definition 9.2.1 (Convergence in law (or distribution)). The sequence X_n converge in law (or distribution) to X , and we write $X_n \xrightarrow{d} X$, if and only if (\iff) ,

$$\lim_{n \rightarrow +\infty} F_{X_n}(x) = F_X(x)$$

$\forall x \in D_X$ in which $F_X(x)$ is continuous.

Definition 9.2.2 (Alternate definition).

$$X_n \xrightarrow{d} X \iff M_{X_n}(t) \rightarrow M_X(t), \forall t : |t| < \varepsilon \quad (9.6)$$

in a intorno di $t = 0$

Remark 317. Two theorem without proof before going on

Theorem 9.2.1. *Convergence in probability is stronger than convergence in distribution since $X_n \xrightarrow{p} X \implies X_n \xrightarrow{d} X$*

Theorem 9.2.2. *...but in the case of dirac we have both implication $X_n \xrightarrow{p} \delta_\theta \iff X_n \xrightarrow{d} \delta_\theta$*

Example 9.2.1. Let $(X_n)_{n \in \mathbb{N}}$ be iid standard normal, $X_n \sim N(0, 1)$. Defining the following variable

$$Y_n = \frac{X_1^2 + \dots + X_n^2}{n} = \frac{\chi_n^2}{n}$$

(at numerator we have a χ_n^2), prove that $Y_n \xrightarrow{d} \delta_1$.

We do it by moment generating function. Looking at the mgf of a chi square we have that:

$$M_{\chi_n^2}(t) = (1 - 2t)^{-n/2}$$

We have that $Y_n = \frac{\chi_n^2}{n}$ so its moment generating function (applying properties)

$$M_{Y_n}(t) = M_{\frac{\chi_n^2}{n}}(t) = M_{\chi_n^2}\left(\frac{t}{n}\right) = \left(1 - 2\frac{t}{n}\right)^{-\frac{n}{2}}$$

We then have

$$\lim_{n \rightarrow +\infty} M_{Y_n}(t) = \lim_{n \rightarrow +\infty} \left(1 - 2\frac{t}{n}\right)^{-\frac{n}{2}} = e^t$$

this remembering that

$$\lim_{n \rightarrow +\infty} \left(1 - \frac{a}{n}\right)^n = e^{-a}$$

$$\lim_{n \rightarrow +\infty} \left(1 + \frac{a}{n}\right)^n = e^a$$

So we have found that

$$\lim_{n \rightarrow +\infty} M_{Y_n}(t) = e^t$$

Now looking at δ_θ it has a simple moment generating function; if $X \sim \delta_\theta$

$$M_X(t) = \mathbb{E}[e^{tX}] = e^{t\theta}$$

therefore if $X \sim \delta_1$, its $M_X(t) = e^t$ and is the limit developed above.

Example 9.2.2. Let $X_n \sim \text{Bin}(n, \frac{\lambda}{n})$. Prove that $X_n \xrightarrow{d} \text{Pois}(\lambda)$. Here again there's a moving probability $X_1 \sim \text{Bin}(n, \lambda)$, $X_2 \sim \text{Bin}(n, \lambda/2)$, $\dots X_n \sim \text{Bin}(n, \lambda/n)$. The mgf of generic binomial rv

$$M_{X_n}(t) = \left(1 - \frac{\lambda}{n} + \frac{\lambda}{n}e^t\right)^n = \left(1 + \frac{\lambda}{n}(e^t - 1)\right)^n$$

Using $\lim(1 + a/n)^n = e^a$ we have that

$$\lim_{n \rightarrow +\infty} \left(1 + \frac{\lambda}{n}(e^t - 1)\right)^n = e^{\lambda(e^t - 1)}$$

But this is the mgf for $\text{Pois}(\lambda)$.

Example 9.2.3 (Viols S01E07). A rv X is said to have the two-parameter Pareto distribution with parameters α and β if its pdf is given by

$$f_X(x) = \frac{\alpha\beta^\alpha}{x^{\alpha+1}}, \quad x > \beta, \alpha, \beta > 0$$

1. show that the function just given is indeed a pdf
2. set $Y = \frac{X}{\beta}$ and show that its pdf is given by $f_Y(y) = \frac{\alpha}{y^{\alpha+1}}$, $y > 1$ and $\alpha > 0$, which is referred to as the one-parameter Pareto distribution
3. show that $\mathbb{E}[X] = \frac{\alpha\beta}{\alpha-1}$
4. show that $\text{Var}[X] = \frac{\alpha\beta^2}{(\alpha-1)^2(\alpha-2)}$ with $\alpha > 2, \beta > 0$
5. let $(X_n)_{n \in \mathbb{N}}$ with $X_n \sim \text{Bin}(n, \frac{\lambda}{n})$. Prove that $X_n \xrightarrow{d} \text{Pois}(\lambda)$.

We have:

1. in order to be a proper pdf

$$\int_{\beta}^{+\infty} \frac{x\beta}{x^{\alpha+1}} dx = 1$$

So

$$\begin{aligned} \int_{\beta}^{+\infty} \frac{x\beta}{x^{\alpha+1}} dx &= \alpha\beta^\alpha \cdot \int_{\beta}^{\infty} \frac{1}{x^{\alpha+1}} = \alpha\beta \cdot \left[-\frac{1}{\alpha} \cdot \frac{1}{x^\alpha}\right]_{\beta}^{\infty} \\ &= \alpha\beta^\alpha \frac{1}{\alpha} \frac{1}{\beta^\alpha} = 1 \end{aligned}$$

2. we apply

$$f_Y(y) = \left| \frac{\partial g^{-1}(y)}{\partial y} \right| f_X(g^{-1}(y))$$

having

$$g^{-1}(y) = \beta y$$

so

$$f_Y(y) = \beta \cdot \alpha \beta^\alpha \frac{1}{\beta^{\alpha+1} y^{\alpha+1}} \underbrace{\mathbb{1}(\beta, +\infty) \beta y}_{= \mathbb{1}(1, +\infty) y}$$

3. we have

$$\begin{aligned} \mathbb{E}[X] &= \int_{\beta}^{+\infty} \frac{x \cdot \alpha \beta^\alpha}{x^{\alpha+1}} dx = \alpha \beta^\alpha \cdot \underbrace{\int_{\beta}^{\infty} \frac{1}{x^\alpha} dx}_{(1)} \\ &= \alpha \beta^\alpha \frac{1}{\alpha-1} \frac{1}{\beta^{\alpha-1}} = \frac{\alpha}{\alpha-1} \beta \end{aligned}$$

where (1) is the kernel of a Pareto with parameters $\alpha = 1$ and $\beta = 0$, therefore $\alpha - 1 > 0$, $\alpha > 1$

4. we have that

$$\mathbb{E}[X^2] = \alpha \beta^\alpha \int_{\beta}^{\infty} x^2 \frac{1}{x^{\alpha+1}} dx = \alpha \beta^\alpha \int_{\beta}^{\infty} \frac{1}{x^{\alpha-1}} dx = \alpha \beta^\alpha \frac{1}{\alpha-2} \frac{1}{\beta^{\alpha-2}} = \frac{\beta^2 \alpha}{\alpha-2}$$

Therefore:

$$\begin{aligned} \text{Var}[X] &= \mathbb{E}[X^2] - \mathbb{E}[X]^2 = \frac{\beta^2 \alpha}{\alpha-2} - \frac{\alpha^2 \beta^2}{(\alpha-1)^2} \\ &= \frac{\beta^2 (\alpha-1)^2 - \alpha^2 \beta^2 (\alpha-2)}{(\alpha-2)(\alpha-1)^2} \\ &= \frac{\beta^2 \alpha^3 + \beta^2 \alpha - 2\beta^2 \alpha^2 - \alpha^3 \beta^2 + 2\alpha^2 \beta^2}{(\alpha-2)(\alpha-1)^2} \end{aligned}$$

5. take $M_{X_n}(t)$ of the binomial:

$$M_{X_n}(t) = \left(1 - \frac{\lambda}{n} + \frac{\lambda}{n} e^t\right)^n = \left(1 - \frac{\lambda}{n} (1 - e^t)\right)^n \stackrel{(1)}{=} \left(1 - \frac{a}{n}\right)^n$$

with (1) taking $a = \lambda(1 - e^t)$. Therefore

$$\lim_{n \rightarrow \infty} \left(1 - \frac{a}{n}\right)^n = e^{-a}$$

therefore $X_n \xrightarrow{d} X$ with $M_X(t) = e^{-\lambda(1-e^t)}$. But this happens $\iff X \sim \text{Pois}(\lambda)$.

Example 9.2.4. Let X be a continuous uniform random variable in $[0, 1]$. Let $Y = \frac{X}{1-X}$ and $Y_n = Y^{1/n}$:

1. Determine $f_Y(y)$ and $F_Y(y)$
2. Determine $F_{Y_n}(y)$
3. Study the convergence in law of Y_n

We have

1. that

$$Y = \frac{X}{1-X} = g(X) \implies g^{-1}(Y) = \frac{1}{1+Y} = X$$

Then

$$\begin{aligned} F_Y(y) &= \mathbb{P}(Y \leq y) = \mathbb{P}\left(\frac{X}{1-X} \leq y\right) = \mathbb{P}(X \leq y - yX) \\ &= \mathbb{P}(X + yX \leq y) = \mathbb{P}\left(X \leq \frac{y}{1+y}\right) = F_X\left(\frac{y}{1+y}\right) \\ &= \frac{y}{1+y} \end{aligned}$$

and so

$$\begin{aligned} f_Y(y) &= f_X(g^{-1}(y)) \cot \left| \frac{\partial g^{-1}(y)}{\partial y} \right| = \mathbb{1}([y/(1+y)])x \cdot |-y(1+y)^{-2} + (1+y)^{-1}| \\ &= \left| \frac{1+y-y}{(1+y)^2} \right| \mathbb{1}(x/1-x)y = \frac{1}{(1+y)^2} \mathbb{1}([0, +\infty])y \end{aligned}$$

2. for the second point

$$F_{Y_n}(y) = \mathbb{P}(Y_n \leq y) = \mathbb{P}(Y^{1/n} \leq y) = \mathbb{P}(Y \leq y^n) = F_Y(y^n) = \frac{y^n}{1+y^n}$$

3. for the third

$$\lim_{n \rightarrow \infty} F_Y(y^n) = \begin{cases} 0 & \text{if } y < 1 \\ 1/2 & \text{if } y = 1 \\ 1 & \text{if } y > 1 \end{cases}$$

so the F of a δ_1 and $F_{Y_n}(y)$ coincides except for $y = 1$, but it is a discontinuity point and we can ignore it.

Therefore $Y_n \xrightarrow{d} \delta_1$

9.2.1 Theorem: central limit theorem

Important remark 76. Fundamental theorem, basis for inference; this is why low number of patients does not permit to have a good approximation (it would be for $n \rightarrow +\infty$ but are needed at least 20/30 patients for the approximation start working)

Remark 318. This can be defined equivalently in terms of partial sum $\sum_{i=1}^n X_i$ or partial mean $\frac{\sum_{i=1}^n X_i}{n}$ of iid random variables with finite expected value and variance.

Proposition 9.2.3. *Let X_i be iid random variables, with mean $\mathbb{E}[X_i] = \mu$ and $\text{Var}[X_i] = \sigma^2$; let $S_n = \sum_{i=1}^n X_i$ be the partial sum and $M_n = \frac{\sum_{i=1}^n X_i}{n}$ the partial mean. If we define the standardized sum as*

$$Z_n = \frac{S_n - \mathbb{E}[S_n]}{\sqrt{\text{Var}[S_n]}} = \frac{S_n - n\mu}{\underbrace{\sqrt{n\sigma^2}}_{\text{no cov, } \perp\!\!\!\perp}} \stackrel{(1)}{=} \frac{M_n - \mu}{\frac{\sigma}{\sqrt{n}}}$$

where in (1) we divided everything by n .

Then $Z_n \xrightarrow{d} N(0, 1)$

Dimostrazione.

$$Z_n = \frac{S_n - n\mu}{\sigma \cdot \sqrt{n}} = \frac{\sum_{i=1}^n X_i - \sum_{i=1}^n \mu}{\sigma \cdot \sqrt{n}} = \sum_{i=1}^n \underbrace{\left(\frac{X_i - \mu}{\sigma} \right)}_{U_i} \cdot \frac{1}{\sqrt{n}} = \frac{\sum_{i=1}^n U_i}{\sqrt{n}}$$

with $\mathbb{E}[U_i] = 0$ and $\text{Var}[U_i] = 1$ (being standardized) and $\mathbb{E}[U_i^2] = 1$ as consequence of the first two using the variance formula $\text{Var}[U_i] = \mathbb{E}[U_i^2] - \mathbb{E}[U_i]^2$.

Now for the moment generating function of Z_n we have

$$M_{Z_n}(t) = M_{\frac{\sum U_i}{\sqrt{n}}}(t) \stackrel{(1)}{=} M_{\sum U_i}(t/\sqrt{n}) \stackrel{(2)}{=} \prod_{i=1}^n M_{U_i}(t/\sqrt{n}) \stackrel{(3)}{=} [M_U(t/\sqrt{n})]^n$$

with (1) by prop of mgf, (2) by independence and (3) since they are identically distributed. Since the mgf of standard normal is $e^{t^2/2}$, we want to prove that

$$\lim_{n \rightarrow +\infty} M_{Z_n}(t) = [M_U(t/\sqrt{n})]^n = e^{t^2/2}$$

We decompose $M_U(t/\sqrt{n})$ by Taylor (in point $t = 0$ so maclaurin) expansion. In general we have that

$$M_X(t) = 1 + t \mathbb{E}[X] + \frac{t^2}{2!} \mathbb{E}[X^2] + \frac{t^3}{3!} \mathbb{E}[X^3] + \dots$$

Applying this to $M_U(t/\sqrt{n})$ (two terms here are enough for what follows):

$$M_U(t/\sqrt{n}) = 1 + \frac{t}{\sqrt{n}} \underbrace{\mathbb{E}[U]}_{=0} + \frac{t^2}{n \cdot 2} \underbrace{\mathbb{E}[U^2]}_{=1} + \dots \simeq 1 + \frac{t^2}{2n}$$

therefore

$$M_{Z_n}(t) \simeq \left(1 + \frac{t^2}{2n}\right)^n$$

Finally

$$\lim_{n \rightarrow +\infty} M_{Z_n}(t) = \lim_{n \rightarrow +\infty} \left(1 + \frac{t^2}{2n}\right)^n = e^{t^2/2}$$

which is the mgf of $N(0, 1)$. □

9.3 Convergence in mean of order k

9.3.1 Definition

Definition 9.3.1. Let $k \in \mathbb{N}^+$. It's said that $X_n \xrightarrow{L_k} X$ if and only if

$$\lim_{n \rightarrow +\infty} \mathbb{E} [|X_n - X|^k] = 0$$

Important remark 77 (Convergence in quadratic mean). One of the most famous is for $n = 2$, $X_n \xrightarrow{L_2} X \iff \lim_{n \rightarrow \infty} \mathbb{E} [(X_n - X)^2] = 0$

9.3.2 Strong consistence

Important remark 78. In inference there are two types of consistency, *weak* consistency and *strong* consistency:

- weak type is convergence in probability
- strong is convergence in L_2 (quadratic mean)

In inference X_n is an estimator and θ is the parameter you want to estimate. Consistency is a good property for an estimator to have; *it's better to have strong because it implies weak.*

Definition 9.3.2 (Strong consistence). If $X_n \xrightarrow{L_2} \delta_\theta$ that is

$$\lim_{n \rightarrow +\infty} \mathbb{E} [(X_n - \theta)^2] = 0$$

we say that X_n is strongly consistent for θ .

Proposition 9.3.1. In this type of convergence we have this result

$$X_n \xrightarrow{L_2} \delta_\theta \iff \begin{cases} \lim_{n \rightarrow +\infty} \mathbb{E} [X_n] = \theta \\ \lim_{n \rightarrow +\infty} \text{Var} [X_n] = 0 \end{cases}$$

Example 9.3.1 (Esame vecchio viroli). Let Y_n be a sequence of independent poisson random variables with parameter $\lambda_n = 1/\sqrt{n}$. Study the convergence in quadratic mean of Y_n .

First we need to decide where it converges. Let's try $\mathbb{E} [Y_n]$

$$\begin{aligned} \mathbb{E} [Y_n] &= \lambda_n = \frac{1}{\sqrt{n}} \\ \lim_{n \rightarrow +\infty} \mathbb{E} [Y_n] &= \lim_{n \rightarrow +\infty} \frac{1}{\sqrt{n}} = 0 \end{aligned}$$

Does it converge to a δ_0 ? let's apply the definition

$$\lim_{n \rightarrow +\infty} \mathbb{E} [(Y_n - 0)^2] = \lim_{n \rightarrow +\infty} \mathbb{E} [Y_n^2]$$

To obtain the second moment of a poisson (considered that $\text{Var} [Y] = \lambda$):

$$\begin{aligned} \text{Var} [Y] &= \mathbb{E} [Y^2] - \mathbb{E} [Y]^2 \\ \lambda &= \mathbb{E} [Y^2] - \lambda^2 \implies \mathbb{E} [Y^2] = \lambda + \lambda^2 \end{aligned}$$

and so in our case $\mathbb{E}[Y_n^2] = \frac{1}{\sqrt{n}} + \frac{1}{n}$. Then

$$\lim_{n \rightarrow \infty} \frac{1}{\sqrt{n}} + \frac{1}{n} = 0$$

Therefore: $Y_n \xrightarrow{L_2} 0$

Example 9.3.2. Let $X_n \sim \text{Pois}(2/n)$; let's check that

1. $X_n \xrightarrow{d} \delta_0$
2. $X_n \xrightarrow{L_2} \delta_0$

We have that

1. for the Poisson distribution we have that $M_{X_n}(t) = e^{\frac{2}{n}(e^t - 1)}$. Taking the limit

$$\lim_{n \rightarrow +\infty} e^{\frac{2}{n}(e^t - 1)} = e^0 = 1$$

For δ_0 , the mgf is

$$M(t) = \mathbb{E}[e^{tX}] = \mathbb{E}[e^0] = 1$$

so same mgf we have proved the convergence

2. we have

$$\lim_{n \rightarrow +\infty} \mathbb{E}[(X_n - 0)^2] = \lim_{n \rightarrow +\infty} \mathbb{E}[X_n^2]$$

To obtain this we can use exploit formula; since X_n is a Poisson

$$\mathbb{E}[X_n] = \frac{2}{n}$$

$$\text{Var}[X_n] = \frac{2}{n}$$

$$\mathbb{E}[X_n^2] = \text{Var}[X_n] + \mathbb{E}[X_n]^2 = \frac{2}{n} + \frac{4}{n^2} = \frac{2n + 4}{n^2}$$

And finally

$$\lim_{n \rightarrow +\infty} \frac{2n + 4}{n^2} = 0$$

so it goes to δ_0

Example 9.3.3. Let $X_n \sim \text{Bern}\left(\frac{1}{n}\right) \cdot n$ so, its pmf be

$$X_n = \begin{cases} n & \text{with probability } 1/n \\ 0 & \text{with probability } 1 - 1/n \end{cases}$$

Study convergence in L_2 and probability.

We have that

$$\mathbb{E}[X_n] = n \cdot \frac{1}{n} = 1$$

$$\mathbb{E}[X_n^2] = n^2 \cdot \frac{1}{n} = n$$

$$\text{Var}[X_n] = n - 1^2 = n - 1$$

Now

- we can't conclude X_n converges in L_2 because of the (limit of the) variance

$$\begin{cases} \lim_{n \rightarrow \infty} \mathbb{E}[X_n] = 1 \\ \lim_{n \rightarrow \infty} \text{Var}[X_n] = +\infty \end{cases} \implies X_n \not\stackrel{L_2}{\longrightarrow} \delta_1$$

- if it converges in probability, where? to two possible distribution
 - what about δ_1 ? we have that

$$\mathbb{P}(|X_n - 1| < \varepsilon) \xrightarrow[n \rightarrow \infty]{} 1$$

convergence is not true because look at $X_n \sim \text{Bern}(1/n)$: 0 with larger and larger prob, 1 with lowering prob. Therefore $X_n - 1$ will be 1 with increasing prob and so $1 \not\leq \varepsilon, \forall \varepsilon \in \mathbb{R}$.

- what about δ_0 ? we have that

$$\mathbb{P}(|X_n| < \varepsilon) = \mathbb{P}(X_n < \varepsilon) \xrightarrow[n \rightarrow \infty]{} 1$$

this is true so $X_n \xrightarrow{p} \delta_0$.

So here we probed the convergence without the two sufficient condition (they're just sufficient, not needed; we can have convergence in prob even if we don't have the two sufficient conditions).

Example 9.3.4. Let X_1, X_2, \dots be a sequence of random variables such that

$$\mathbb{P}\left(X_n = \frac{1}{n}\right) = 1 - \frac{1}{n^2} \quad \mathbb{P}(X_n = n) = \frac{1}{n^2}$$

- Does X_n converge in quadratic mean?
- Does it converge in probability

Respectively

1. For L_2 we should prove $\lim_{n \rightarrow \infty} \mathbb{E}[(X_n - X)^2] = 0$ but who is X ? By reasoning we see that $X_n \rightarrow 0$ with probability $\rightarrow 1$ therefore we try with a δ_0

$$\lim_{n \rightarrow \infty} \mathbb{E}[(X_n - 0)^2] = \lim_{n \rightarrow \infty} \mathbb{E}[X_n^2]$$

The second moment is

$$\mathbb{E}[X_n^2] = \frac{1}{n^2} \cdot \left(1 - \frac{1}{n^2}\right) + n^2 \frac{1}{n^2} = \frac{n^2 - 1}{n^4} + 1$$

so

$$\lim_{n \rightarrow \infty} \mathbb{E}[X^2] = 1$$

so we conclude that $X_n \not\stackrel{L_2}{\longrightarrow} \delta_0$.

2. let's check the two sufficient assumptions for the convergence in probability:

(a) for the first we have

$$\mathbb{E}[X_n] = \frac{1}{n} \left(1 - \frac{1}{n^2} \right) = \frac{n^2 - 1}{n^3}$$

and $\lim_{n \rightarrow \infty} \mathbb{E}[X_n] = 0$

(b) for the second

$$\text{Var}[X_n] = \mathbb{E}[X_n^2] - \mathbb{E}[X_n]^2 = \frac{n^2 - 1}{n^2} + 1 - \frac{(n^2 - 1)^2}{n^6}$$

from which $\lim_{n \rightarrow \infty} \text{Var}[X_n] \rightarrow 1$

However by applying the definition

$$\lim_{n \rightarrow \infty} \mathbb{P}(|X_n| < \varepsilon) = 1$$

$= \lim_{n \rightarrow \infty} \mathbb{P}(X_n < \varepsilon) = 1$ and this is true since $X_n \rightarrow 0$ with probability $\rightarrow 1$ as $n \rightarrow \infty$

9.3.3 Theorem: strong law of large numbers

Remark 319. It's the most important theorem related to convergence in quadratic mean.

Theorem 9.3.2 (Strong law of large numbers). *Let $(X_n)_{n \in \mathbb{N}}$ be a sequence of independent random variables and assume $\mathbb{E}[X_n] = \mu$, $\text{Var}[X_n] = \sigma^2 < +\infty$. Then we say that the partial mean:*

$$M_n = \frac{\sum_{i=1}^n X_i}{n} \xrightarrow{L_2} \mu$$

Dimostrazione.

$$\begin{aligned} \mathbb{E}[(M_n - \mu)^2] &= \mathbb{E}\left[\left(\frac{\sum_{i=1}^n X_i}{n} - \frac{n\mu}{n}\right)\right] \stackrel{(1)}{=} \frac{1}{n^2} \mathbb{E}\left[\sum_{i=1}^n (X_i - \mu)^2\right] \\ &= \frac{1}{n^2} \sum_{i=1}^n \mathbb{E}[(X_i - \mu)^2] = \frac{1}{n^2} \sum_{i=1}^n \text{Var}[X_i] = \frac{1}{n^2} \cdot n \cdot \sigma^2 \\ &= \frac{\sigma^2}{n} \end{aligned}$$

where in (1) due to independence, the expectations of the cross products are all zeros, so the square of sums is the sum of squares. Finally

$$\lim_{n \rightarrow +\infty} \mathbb{E}[(M_n - \mu)^2] = \lim_{n \rightarrow +\infty} \frac{\sigma^2}{n} = 0$$

so $M_n \xrightarrow{L_2} \mu$

□

Example 9.3.5. Let $X_1, \dots, X_n \sim \text{Exp}(1)$. Find the distribution of $X_{(1)} = \min(X_1, \dots, X_n)$ and study its convergence.

Remembering that $F_{(1)} = 1 - [1 - F_X(x)]^n$ and being X exponential we have $F_X(x) = 1 - e^{-x}$, therefore:

$$F_{(1)}(x) = 1 - [1 - 1 + e^{-x}]^n = 1 - e^{-xn}$$

which is the pdf of $\text{Exp}(n)$. So even the minimum is distributed according to an exponential but of parameter n , which are the number of rvs we consider; that is $X_{(1)} \sim \text{Exp}(n)$.

Regarding the convergence to study,

- in this exercise, since we have the pdf of the minimum, it's convenient for us to try to study the limit of it, that is *in this case we study convergence in distribution* (using the cumulative distribution function, not the mgf or the characteristic function). If we find that the limit is a certain pdf we have the solution (finding which random variable gives that pdf). So let's study the limit of F :

$$\lim_{n \rightarrow \infty} F_{(1)}(x) = \lim_{n \rightarrow \infty} 1 - e^{-xn} = 1$$

At the same time 1 is equal to e^0 which is the cumulative distribution function of a δ_0 in 0: $e^0 = F_{\delta_0}(x)$.

Therefore the minimum converges in distribution to a Dirac in 0 but this also implies that it converge in probability:

$$X_{(1)} \xrightarrow{d} \delta_0 \implies X_{(1)} \xrightarrow{p} \delta_0$$

- now we could study a strong kind of convergence; in this case it's convenient to try studying the L_2 convergence, since we know the limiting distribution (the constant 0), so the expectation should be simpler. Furthermore the limit should be the same: if I know that it converges in distribution to a point, if it converges also in quadratic mean, then it should be at the same point (given the implication schema), it can't be another point.

$$\mathbb{E}[(X_{(1)} - 0)^2] = \underbrace{\mathbb{E}[X_{(1)}^2]}_{\text{second moment of Exp}(n)} = \underbrace{\frac{1}{n^2}}_{\text{variance}} + \underbrace{\frac{1}{n^2}}_{\text{second moment squared}} = \frac{2}{n^2}$$

Finally for the convergence in quadratic mean we should study the limit and check that it goes to 0. So:

$$\lim_{n \rightarrow +\infty} \frac{2}{n^2} = 0$$

Therefore we can conclude that

$$X_{(1)} \xrightarrow{L_2} \delta_0 \implies X_{(1)} \xrightarrow{L_1} \delta_0$$

9.4 Almost sure convergence

Remark 320. It's a strong convergence

TODO: non chiarissimo, la cumulata dovrebbe essere una step function non una costante, poi ok che da 0 in poi sia a 1.

Definition 9.4.1. A sequence converges almost surely to a limit distribution X , and we write $X_n \xrightarrow{a.s.} X \iff \mathbb{P}(\lim_{n \rightarrow \infty} |X_n - X| < \varepsilon) = 1$

Remark 321. Difficult to prove because it's not the limit of a probability but the probability of a limit.

Remark 322. The most important associated theorem with a.s. convergence is the following; somewhat similar to the strong/weak law large number.

Theorem 9.4.1 (Kolmogorov theorem). *Let $(X_n)_{n \in \mathbb{N}}$ be iid rvs such as $\mathbb{E}[X_n] = \mu$ is constant/fixed (no assumption on variance here); then it's possible to prove that the partial mean $M_n \xrightarrow{a.s.} \mu$*

Dimostrazione. No proof here, quite complicate. \square

Example 9.4.1. Let be $X_n \sim \text{Pois}(\lambda)$ a sequence of iid rvs; study the convergence of $Z_n = \frac{1}{n} \sum_{i=1}^n \frac{1}{1+X_i}$.

Let's define a continuous transformation of X_i that is $Y_i = \frac{1}{1+X_i}$ and so $Z_n = \frac{\sum_i Y_i}{n}$ is like a partial mean (we have many theorem associated to partial mean: weak/strong laws of large numbers and Kolmogorov theorem). Note that if X_1, \dots, X_n are iid then also Y_1, \dots, Y_n are iid as well (the transformation applied is the same and when we transform independent rv the independence is preserved, unless we combine different rvs).

If we can prove almost sure convergence then we have also the other one so it's convenient to start from the strongest, in case.

So according to Kolmogorov $M_n \xrightarrow{a.s.} \mu$ where in our case $\mu = \mathbb{E}[Y_i]$. Now let's see what is μ :

$$\begin{aligned} \mu &= \mathbb{E} \left[\frac{1}{1+X_i} \right] = \sum_{D_X} \frac{1}{1+x_i} \mathbb{P}(X_i = x_i) = \sum_{x=0}^{+\infty} \frac{1}{1+x} \frac{e^{-\lambda} \lambda^x}{x!} \\ &= \sum_{x=0}^{+\infty} \frac{1}{(x+1)!} e^{-\lambda} \lambda^x \cdot \frac{\lambda}{\lambda} = \frac{1}{\lambda} \sum_{x=0}^{+\infty} \frac{1}{(x+1)!} e^{-\lambda} \lambda^{x+1} \stackrel{(1)}{=} \frac{1}{\lambda} \sum_{t=1}^{+\infty} \underbrace{\frac{1}{t!} e^{-\lambda} \lambda^t}_{\text{Pois}(\lambda)} \\ &= \frac{1}{\lambda} (1 - e^{-\lambda}) \end{aligned}$$

where in (1) we made substitution $t = x + 1$ and considered that the sum is a Poisson without the probability for $t = 0$, starting the sum from 1). Therefore

$$Z_n \xrightarrow{a.s.} \mu = \frac{1}{\lambda} (1 - e^{-\lambda})$$

and then

$$Z_n \xrightarrow{a.s.} \delta_\mu \implies Z_n \xrightarrow{p} \delta_\mu \implies Z_n \xrightarrow{d} \delta_\mu$$

We can stop here since we proved all the convergences; if one can a strong type it's perfect.

Important remark 79. We don't need here to study L_k convergence since we already have a strong kind of convergence; it's enough to prove one of them. (We could try but it's not easy in the previous case).

Example 9.4.2. Study the convergence of $Y_n = (X_1 \cdot \dots \cdot X_n)^{1/n}$ where $X_i \sim \text{Unif}(0, 1)$ are iid rvs.

We need to think about a possible trick and it's given by the continuous mapping theorem (section below) which states that we can maintain convergence if we apply some continuous transformation (except for convergence in mean of order k , where g have to be both continuous and linear).

The transformation we should apply here is the logarithm because we have products and logarithm of a product is a sum.

Therefore consider the transformation $\log Y_n = \frac{1}{n} \sum_{i=1}^n \log X_i$; again we notice this is a partial mean and therefore could think of the strongest theorem we have, which is Kolmogorov; then we can say $M_n \xrightarrow{a.s.} \mu$, and as before we have to find $\mu = \mathbb{E}[\log X]$ where $X \sim \text{Unif}(0, 1)$. Therefore:

$$\mu = \mathbb{E}[\log X] = \int_0^1 \log x \cdot 1 \, dx \stackrel{(1)}{=} [x \log x - x]_0^1 = -1$$

where in (1) we did it by parts i guess. Therefore

$$\frac{1}{n} \sum_{i=1}^n \log X_i \xrightarrow{a.s.} -1$$

So by applying the continuous mapping theorem (we apply the inverse of the logarithm which is the exponential to both the sides of the convergence)

$$Y_n \xrightarrow{a.s.} e^{-1} = \frac{1}{e} \implies Y_n \xrightarrow{a.s., p, d} \delta_{\frac{1}{e}}$$

Example 9.4.3 (Assignment 1 Viroli, Exercise 4). Let X_1, \dots, X_n be a sequence of independent random variables with $X \sim \text{Exp}(\theta)$. Let $T_n = \frac{\sum_{i=1}^n e^{-X_i}}{n}$. Study the convergence of T_n as n goes to infinity.

By setting $Y_i = e^{-X_i}$ we have that $T_n = \frac{\sum_{i=1}^n Y_i}{n}$ so, being a partial mean of iid rvs with $\mathbb{E}[Y_i]$ constant (to be evaluated), we have that $T_n \xrightarrow{a.s.} \mathbb{E}[Y_i]$ by Kolmogorov theorem. Let's evaluate $\mathbb{E}[Y_i]$:

$$\begin{aligned} \mathbb{E}[Y_i] &= \mathbb{E}[e^{-X_i}] = \int_{D_X} e^{-x} \cdot \underbrace{f(x)}_{\text{Exp}(\theta)} \, dx = \int_0^{+\infty} e^{-x} \cdot \theta \cdot e^{-\theta x} \, dx \\ &= \theta \int_0^{+\infty} e^{-x-\theta x} \, dx = \theta \int_0^{+\infty} e^{-x-\theta x} \cdot \frac{(-1-\theta)}{(-1-\theta)} \, dx \\ &= \frac{\theta}{-1-\theta} \int_0^{+\infty} e^{-x-\theta x} \cdot (-1-\theta) \, dx = -\frac{\theta}{1+\theta} [e^{-x-\theta x}]_0^{+\infty} \\ &= -\frac{\theta}{1+\theta} [0 - 1] = \frac{\theta}{1+\theta} \end{aligned}$$

So we can conclude that

$$T_n \xrightarrow{a.s., p, d} \delta_{\frac{\theta}{1+\theta}}$$

Clearly as $\theta \rightarrow +\infty \implies T_n \rightarrow 1$; in figure ?? some heuristic checks for $\theta = 0.1, 1, 10$, (where if calculation above is ok, T_n should converge to $\frac{0.1}{1.1}, \frac{1}{2}, \frac{10}{11}$, horizontal dotted black lines).

9.5 Convergences properties

Proposition 9.5.1 (Properties). *Convergence implications are summarized in the following schema: to be read as “if $X_n \xrightarrow{a.s.} X$, then $X_n \xrightarrow{p}$ to the same X ”:*

$$\begin{array}{ccccc} L_k \rightarrow & \xRightarrow{k>s} & L_s \rightarrow & & \\ & & \Downarrow & & \\ a.s. \rightarrow & \xRightarrow{} & p \rightarrow & \xRightarrow{} & d \rightarrow \end{array}$$

Finally, there's only a special case of double implication between \xrightarrow{p} and \xrightarrow{d} :

$$\xrightarrow{p} \delta_\theta \iff \xrightarrow{d} \delta_\theta$$

Example 9.5.1 (Esame vecchio viroli). Indicate which of the following definitions is false: the convergence in mean of order 4 implies:

1. convergence in quadratic mean
2. the convergence in mean of order 3
3. the almost sure convergence
4. the convergence in distribution

We have that $\xrightarrow{L_4} \not\Rightarrow \xrightarrow{a.s.}$.

Theorem 9.5.2 (Continuous mapping theorem). *Let $(X_n)_{n \in \mathbb{N}}$ be rvs with some domain D_{X_n} . If g is a continuous function on the same domain D_{X_n} , the follow applies:*

$$X_n \xrightarrow{p} X \implies g(X_n) \xrightarrow{p} g(X) \quad (9.7)$$

$$X_n \xrightarrow{d} X \implies g(X_n) \xrightarrow{d} g(X) \quad (9.8)$$

$$\begin{cases} X_n \xrightarrow{L_k} X \\ g \text{ is linear} \end{cases} \implies g(X_n) \xrightarrow{L_k} g(X) \quad (9.9)$$

$$X_n \xrightarrow{a.s.} X \implies g(X_n) \xrightarrow{a.s.} g(X) \quad (9.10)$$

Remark 323. For the L_k case: if g is quadratic, log, exponential etc, being not a linear function, then the implication convergence doesn't hold.

Proposition 9.5.3 (Further properties). *We have that*

1. for convergence in probability

$$(X_n \xrightarrow{p} X \wedge Y_n \xrightarrow{p} Y) \implies aX_n + bY_n \xrightarrow{p} aX + bY \quad (9.11)$$

$$(X_n \xrightarrow{p} X \wedge Y_n \xrightarrow{p} Y) \implies X_n \cdot Y_n \xrightarrow{p} X \cdot Y \quad (9.12)$$

2. same as above applies for $\xrightarrow{a.s.}$

3. for $\xrightarrow{L_k}$ we only have

$$(X_n \xrightarrow{L_k} X \wedge Y_n \xrightarrow{L_k} Y) \implies aX_n + bY_n \xrightarrow{L_k} aX + bY \quad (9.13)$$

but the product does not hold

4. for \xrightarrow{d} we have Slutsky theorem:

$$(X_n \xrightarrow{d} X \wedge Y_n \xrightarrow{d} \delta_c) \implies \begin{cases} X_n + Y_n \xrightarrow{d} X + c \\ X_n \cdot Y_n \xrightarrow{d} cX \end{cases} \quad (9.14)$$

9.6 Delta method

Remark 324. This is a very useful tool for inference.

Important remark 80 (Motivation). From now on we think of this sequence X_n of random variable as an estimator for a parameter θ of interest; most of time n is the sample size. Imagine that you know that your estimator converges in distribution, as sample goes larger, to the constant θ

$$(X_n)_{n \in \mathbb{N}} \xrightarrow{d} \delta_\theta$$

So we can use our estimator to estimate θ .

Delta method is needed if we are interested not on θ but on a transformation on the parameter $g(\theta)$, with g continuous; this because using the continuous mapping theorem is not always optimal.

Example 9.6.1 (Motivating example: odd). Let $X_1, \dots, X_n \sim \text{Bern}(p)$ be independent, with $\mathbb{E}[X_i] = p$ and consider the partial mean $Y_n = \bar{X}_n = \frac{\sum_i X_i}{n}$. We know that, respectively by weak law of large number and by central limit theorem (it's a sum, not standardized) that:

$$\begin{aligned} Y_n &\xrightarrow{p} \delta_p \\ Y_n &\xrightarrow[\text{by CLT}]{d} N\left(p, \frac{p(1-p)}{n}\right) \end{aligned}$$

Some remarks:

1. the two limits above are not conflicting: by the clt we have a distribution but if $n \rightarrow \infty$ the variance of the gaussian goes to 0 and the distribution converges to a Dirac like the first one. In other terms these two results above are asymptotically equivalent (they are the same limit) since $\lim_{n \rightarrow \infty} \frac{p(1-p)}{n} = 0$.
2. the second result however is more useful to know: it's better for us to have a distribution rather than a point. According to gaussian distribution, we can construct intervals, we can test hypotheses, so we can use the idea that we have a distribution for this kind of things, very important from the inferential pov.

Now suppose we're interested not in p of event, but in its odd, that is:

$$g(p) = \frac{p}{1-p}$$

We know that (continuous mapping theorem), the transformation of the sequence converges to the transformation of the limit distribution:

$$Y_n \xrightarrow{p} p \implies g(Y_n) \xrightarrow{p} g(p) \iff \text{odd} \xrightarrow{p} \frac{p}{1-p} \iff \frac{\bar{x}}{1-\bar{x}} \xrightarrow{p} \frac{p}{1-p}$$

However this is a point results; we may be interested in constructing confidence intervals and hypothesis testing and for all that shit we need a proper distribution, not a point.

Therefore here comes the delta method.

Remark 325. To define the delta method first we need the generalized version of CLT.

Theorem 9.6.1 (Generalized version of the central limit theorem). *If we have that $\sqrt{n}(Y_n - \theta) \xrightarrow{d} Y$ converges to a limit distribution Y , then we also have the following equivalent facts (si riporta anche il primo) with $Z \sim N(0, 1)$*

$$\begin{cases} \sqrt{n}(Y_n - \theta) \xrightarrow{d} Y \\ \sqrt{n}(Y_n - \theta) \xrightarrow{d} \sigma Z \\ \frac{Y_n - \theta}{\sigma/\sqrt{n}} \xrightarrow{d} Z \sim N(0, 1) \\ Y_n \xrightarrow{d} Y \sim N(\theta, \sigma^2/n) \end{cases}$$

Important remark 81 (jargon/style). So we can say that a standardized random variable converges to Z , where $Z \sim N(0, 1)$, by writing it according to the first or the second expression. If one write according to first or second expression, one is using the so called generalized version of the central limit theorem.

Example 9.6.2 (Odd example continued). Coming back to our example we have that $Y_n \xrightarrow{d} N\left(p, \frac{p(1-p)}{n}\right)$; then we can rewrite using the generalized CLT

$$Y_n \xrightarrow{d} N\left(p, \frac{p(1-p)}{n}\right) \quad \text{centering ...}$$

$$Y_n - p \xrightarrow{d} N\left(0, \frac{p(1-p)}{n}\right) \quad \text{multiply both by } \sqrt{n} \dots$$

$$\sqrt{n}(Y_n - p) \xrightarrow{d} N(0, p(1-p)) \quad (1)$$

$$\sqrt{n}(Y_n - p) \xrightarrow{d} Z \cdot \sqrt{p(1-p)} \quad (2)$$

where in (1) and (2) remember that $cN(0, b) = N(0, bc^2)$ by the property of the standard gaussian and, again, $Z \sim N(0, 1)$. This is another example where starting from a gaussian I can rewrite it in a generalized form.

Last one is the generalized-CLT version-style; we need it for the delta method.

Proposition 9.6.2 (Delta method). *If the generalized CLT holds, that is:*

$$\sqrt{n}(Y_n - \theta) \xrightarrow{d} Y$$

we have that

$$\sqrt{n}(g(Y_n) - g(\theta)) \xrightarrow{d} g'(\theta) \cdot Y \quad (9.15)$$

Delta method proof. To answer consider Taylor expansion of the first order of $g(Y_n)$ at the point θ . It's sufficient to stop at first derivative:

$$g(Y_n) = g(\theta) + g'(\theta)(Y_n - \theta) + \dots$$

therefore

$$g(Y_n) - g(\theta) \simeq g'(\theta)(Y_n - \theta)$$

so multiplying by \sqrt{n}

$$\sqrt{n}(g(Y_n) - g(\theta)) \simeq g'(\theta) \underbrace{\sqrt{n}(Y_n - \theta)}_{\xrightarrow{d} Y}$$

Given the generalized version of the CLT the last part converges to Y so we have the final formula of the delta method which is

$$\sqrt{n}(g(Y_n) - g(\theta)) \xrightarrow{d} g'(\theta) \cdot Y$$

□

Important remark 82 (Motivation recap (general X)). Imagine we have a sequence which converges to a random variable X

$$(X_n)_{n \in \mathbb{N}} \xrightarrow{d} X$$

But are interested on $g(X_n)$ with g continuous (eg the odd). The question is what is the limit distribution of $\sqrt{n}(g(X_n) - g(\theta)) \xrightarrow{d} ?$

Delta method is a method to derive the limit distribution of a transformation starting from the limit distribution of the original variable.

The convergency is a convergency in distribution/law and it says that if the generalized clt holds, you have as result the same limit Y multiplied by the derivative of the transformation.

Example 9.6.3 (Odd example conclusion). The delta method is a tool that gives us a distribution for the odds. We can apply it since the generalized clt holds, as shown above:

$$Y_n \xrightarrow{d} N\left(p, \frac{p(1-p)}{n}\right) \implies \sqrt{n}(Y_n - p) \xrightarrow{d} \sqrt{p(1-p)} N(0, 1)$$

To apply the delta method formula we have to find the first derivative of the transformation

$$g(p) = \frac{p}{1-p}$$

$$g'(p) = \frac{1(1-p) - (-1)p}{(1-p)^2} = \frac{1-p+p}{(1-p)^2} = \frac{1}{1-p}$$

Now we can find the estimator for the odds and also its asymptotic distribution. Now with \bar{x} as our estimator for p we can say that

$$\sqrt{n}(\bar{x} - p) \xrightarrow{d} N(0, p(1-p))$$

and according to the delta method we can say that

$$\begin{aligned}\sqrt{n}(g(Y_n) - g(\theta)) &\xrightarrow{d} g'(\theta) \cdot Y \\ \sqrt{n}\left(\frac{\bar{x}}{1-\bar{x}} - \frac{p}{1-p}\right) &\xrightarrow{d} \frac{1}{(1-p)^2} \cdot N(0, p(1-p)) \\ &\xrightarrow{d} N\left(0, \frac{p(1-p)}{(1-p)^4}\right) \\ &\xrightarrow{d} N\left(0, \frac{p}{(1-p)^3}\right)\end{aligned}$$

Example 9.6.4 (Logarithm of the mean). Having X_1, \dots, X_n are iid with dist $f(x)$ (whatever distribution), $\mathbb{E}[X] = \mu$, $\text{Var}[X] = \sigma^2$ if we take the average $Y_n = \bar{X}_n = \sum_{i=1}^n X_i/n$ as our estimator, with the clt we have the

$$\sqrt{n}(\bar{x} - \mu) \xrightarrow{d} \sigma N(0, 1)$$

Now what is the distribution of the estimator for the logarithm of μ $g(\mu) = \log(\mu)$? Applying the delta method we have:

$$\begin{aligned}g(\mu) &= \log(\mu) \\ g'(\mu) &= \frac{1}{\mu}\end{aligned}$$

So:

$$\begin{aligned}\sqrt{n}(g(\bar{x}) - g(\mu)) &\xrightarrow{d} g'(\mu) \cdot \sigma \cdot N(0, 1) \\ \sqrt{n}(\log(\bar{x}) - \log \mu) &\xrightarrow{d} \frac{1}{\mu} \cdot \sigma \cdot N(0, 1) \\ &\xrightarrow{d} N\left(0, \frac{\sigma^2}{\mu^2}\right)\end{aligned}$$

OR better, in explicit way:

$$\log \bar{x} \xrightarrow{d} N\left(\log \mu, \frac{\sigma^2}{n\mu^2}\right)$$

Example 9.6.5. Let X_1, \dots, X_n iid, with $X_i \sim f_X(x)$, $\mathbb{E}[X] = \mu$, $\text{Var}[X] = \sigma^2$. Find the asymptotic distribution of the second moment \bar{X}_n^2 .

We know that by CLT

$$\sqrt{n}(\bar{x} - \mu) \xrightarrow{d} \sigma N(0, 1)$$

According to Delta method

$$\sqrt{n}(\bar{x}^2 - \mu^2) \xrightarrow{d} g'(\mu)\sigma N(0, 1)$$

with

$$\begin{aligned}g(\mu) &= \mu^2 \\ g'(\mu) &= 2\mu\end{aligned}$$

then we conclude that

$$\begin{aligned}\sqrt{n}(\bar{x}^2 - \mu^2) &\xrightarrow{d} 2\mu\sigma N(0, 1) \\ &\xrightarrow{d} N(0, 4\mu^2\sigma^2)\end{aligned}$$

Example 9.6.6 (Esame vecchio viroli). Let $\hat{\theta}_n$ be an estimator for θ with the following asymptotic distribution

$$\sqrt{n}(\hat{\theta}_n - \theta) \xrightarrow{d} \sqrt{\theta} N(0, 1)$$

Use the delta method to derive the asymptotic distribution of $g(\hat{\theta}_n) = \log \hat{\theta}_n$:

1. $\sqrt{n}(\log \hat{\theta}_n - \log \theta) \xrightarrow{d} N(0, \frac{1}{4})$
2. $\sqrt{n}(\log \hat{\theta}_n - \log \theta) \xrightarrow{d} N(0, \frac{1}{\theta})$
3. $\sqrt{n}(\log \hat{\theta}_n - \log \theta) \xrightarrow{d} N(0, \frac{1}{\theta^2})$
4. $\sqrt{n}(\log \hat{\theta}_n - \log \theta) \xrightarrow{d} N(0, \frac{2}{\theta^2})$

We have that $g(\theta) = \log(\theta)$ and $g'(\theta) = \frac{1}{\theta}$ so, by the delta method

$$\begin{aligned}\sqrt{n}(\log(\hat{\theta}_n) - \log(\theta)) &\xrightarrow{d} g'(\theta) \cdot \sqrt{\theta} N(0, 1) \\ &\xrightarrow{d} \frac{1}{\theta} \cdot \sqrt{\theta} N(0, 1) \\ &\xrightarrow{d} N\left(0, \frac{1}{\theta}\right)\end{aligned}$$

as reported by Bigo as well

Example 9.6.7 (Esame vecchio viroli). Let $\hat{\theta}_n$ be an estimator for θ with the following asymptotic distribution

$$\sqrt{n}(\hat{\theta}_n - \theta) \xrightarrow{d} \sqrt{\theta} N(0, 1)$$

Use the delta method to derive the asymptotic distribution of $g(\hat{\theta}_n) = \frac{\hat{\theta}_n^2}{2} + 2$:

1. $\sqrt{n}(\frac{\hat{\theta}_n^2}{2} - \frac{\theta^2}{2}) \xrightarrow{d} N(0, \theta^3) + 2$
2. $\sqrt{n}(\frac{\hat{\theta}_n^2}{2} - \frac{\theta^2}{2}) \xrightarrow{d} N(0, \theta^3)$
3. $\sqrt{n}(\frac{\hat{\theta}_n^2}{2} - \frac{\theta^2}{2}) \xrightarrow{d} N(0, \frac{\theta^2}{2})$
4. $\sqrt{n}(\frac{\hat{\theta}_n^2}{2} - \frac{\theta^2}{2}) \xrightarrow{d} N(0, \frac{\theta^4}{4})$

qui si ha che $g(x) = \frac{x^2}{2} + 2$ da cui $g'(x) = x$ e $g'(\theta) = \theta$. Per cui

$$\begin{aligned}\sqrt{n}(g(\hat{\theta}_n) - g(\theta)) &\xrightarrow{d} g'(\theta) \sqrt{\theta} N(0, 1) \\ \sqrt{n}\left(\frac{\hat{\theta}_n^2}{2} + 2 - \frac{\theta^2}{2} - 2\right) &\xrightarrow{d} \theta^{\frac{3}{2}} N(0, 1) \\ \sqrt{n}\left(\frac{\hat{\theta}_n^2}{2} - \frac{\theta^2}{2}\right) &\xrightarrow{d} N(0, \theta^3)\end{aligned}$$

Example 9.6.8 (Esame vecchio viroli). Let $\hat{\theta}_n$ be an estimator for θ with the following asymptotic distribution

$$\sqrt{n}(\hat{\theta}_n - \theta) \xrightarrow{d} \frac{2}{\theta} N(0, 1)$$

Use the delta method to derive the asymptotic distribution of $g(\hat{\theta}_n) = \sqrt{\hat{\theta}_n}$:

- $\sqrt{n}(\sqrt{\hat{\theta}_n} - \sqrt{\theta}) \xrightarrow{d} N(0, \frac{4}{boh})$
- $\sqrt{n}(\sqrt{\hat{\theta}_n} - \sqrt{\theta}) \xrightarrow{d} N(0, \frac{1}{\theta_3})$
- $\sqrt{n}(\sqrt{\hat{\theta}_n} - \sqrt{\theta}) \xrightarrow{d} N(0, \frac{2}{boh})$
- $\sqrt{n}(\sqrt{\hat{\theta}_n} - \sqrt{\theta}) \xrightarrow{d} N(0, \frac{1}{\theta^2})$

qui si ha $g(x) = \sqrt{x}$, $g'(x) = \frac{1}{2\sqrt{x}}$ e $g'(\theta) = \frac{1}{2\sqrt{\theta}}$. Da cui

$$\begin{aligned} \sqrt{n}(\sqrt{\hat{\theta}_n} - \sqrt{\theta}) &\xrightarrow{d} \frac{1}{2\sqrt{\theta}} \frac{2}{\theta} N(0, 1) \\ &\xrightarrow{d} N(0, \theta^{-3}) \end{aligned}$$

Example 9.6.9. Let X_1, \dots, X_n be independent $\text{Geom}(p)$

1. Does $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$ converge in probability?
2. what is its limiting distribution?
3. and what is the distribution of $\frac{1}{\bar{X}_n}$

We have that

1. According to the WLLN $\bar{X}_n \xrightarrow{p} \mathbb{E}[X] = \frac{1}{p}$.
2. The limiting distribution can be derived by the CLT

$$\sqrt{n}\left(\bar{X}_n - \frac{1}{p}\right) \xrightarrow{d} N\left(0, \frac{1-p}{p^2}\right)$$

with $\frac{1-p}{p^2}$ as variance.

3. The limiting distribution of $\frac{1}{\bar{X}_n}$ can be found by the Delta method. We have that

$$\begin{aligned} g(x) &= \frac{1}{x} \\ g'(x) &= -\frac{1}{x^2} \end{aligned}$$

So considering $\theta = \frac{1}{p}$ we have that

$$\begin{aligned}\sqrt{n}\left(\bar{X}_n - \frac{1}{p}\right) &\xrightarrow{d} N\left(0, \frac{1-p}{p^2}\right) \\ \Rightarrow \\ \sqrt{n}\left(\frac{1}{\bar{X}_n} - p\right) &\xrightarrow{d} g'(\theta) N\left(0, \frac{1-p}{p^2}\right) \\ \sqrt{n}\left(\frac{1}{\bar{X}_n} - p\right) &\xrightarrow{d} -\frac{1}{(1/p)^2} N\left(0, \frac{1-p}{p^2}\right) \\ \sqrt{n}\left(\frac{1}{\bar{X}_n} - p\right) &\xrightarrow{d} -p^2 N\left(0, \frac{1-p}{p^2}\right) \\ \sqrt{n}\left(\frac{1}{\bar{X}_n - p}\right) &\xrightarrow{d} N\left(0, \frac{p^4(1-p)}{p^2}\right)\end{aligned}$$

Example 9.6.10. Let X_1, \dots, X_n a sequence of independent rvs with $X \sim \text{Exp}(\theta)$. Let $T_n = \sum_{i=1}^n \frac{X_i}{2n}$

1. Does T_n converge in probability?
2. Find the limiting distribution of T_n by CLT
3. find the limiting distribution of $\log(T_n)$

For

1. the convergence in probability we have that

$$\begin{aligned}\mathbb{E}[T_n] &= \frac{\sum_{i=1}^n \mathbb{E}[X_i]}{2n} = \frac{\sum_{i=1}^n \frac{1}{\theta}}{2n} = \frac{1}{2\theta} \\ \text{Var}[T_n] &= \frac{1}{4n^2} \sum_{i=1}^n \text{Var}[X_i] = \frac{n}{4n^2\theta^2} = \frac{1}{4n\theta^2}\end{aligned}$$

therefore $T_n \xrightarrow{p} \delta_{1/2\theta}$

2. for the convergence in distribution by CLT let's first study $T_n^* = 2T_n = \frac{\sum_{i=1}^n X_i}{n} = \bar{X}_n$. By CLT

$$\sqrt{n}(\bar{X}_n - 1/\theta) \xrightarrow{d} N\left(0, \frac{1}{\theta^2}\right)$$

since

$$\frac{\bar{X}_n - 1/\theta}{\frac{1}{\sqrt{n}\theta}} \xrightarrow{d} N(0, 1)$$

with $\mathbb{E}[\bar{X}_n] = \frac{1}{\theta}$, $\text{Var}[\bar{X}_n] = \frac{1}{n\theta^2}$. So by the continuous mapping theorem

$$\frac{T_n - 1/2\theta}{\frac{1}{2\sqrt{n}\theta}} \xrightarrow{d} N(0, 1)$$

and the generalized form is

$$\sqrt{n}\left(T_n \frac{1}{2\theta}\right) \xrightarrow{d} \frac{1}{2\theta} \cdot N(0, 1)$$

3. for the convergence of $\log T_n$, by the delta method

$$\begin{aligned}\sqrt{n}\left(\log T_n - \log \frac{1}{2\theta}\right) &\xrightarrow{d} g'(\theta) \cdot \frac{1}{2\theta} \text{N}(0, 1) \\ &\xrightarrow{d} 2\theta \frac{1}{2\theta} \text{N}(0, 1)\end{aligned}$$

Capitolo 10

Simulation

10.1 Sampling values from rvs

Important remark 83. This is important for practical/inferential reasons: some methods in statistics need sampling from rvs to get estimates, and not always distribution are available/easy to use (eg complicate, not popular, not well known, or because we don't know completely the analytical stuff eg we know the kernel not the normalization constant).

Important remark 84. So it's important in difficult situation to have a method for sampling from the distribution (to have some values), because if we draw infinite time we can obtain the distribution.

Important remark 85. The methods available are summarized in table 10.1. Vi-
roli will do univariate methods. The MCMC stuff (Gibbs sampling, Metropolis-
Hasting) will be done in Bayesian statistics.

10.1.1 Inversion method

Remark 326. This is the simpler method

Important remark 86. If $X \sim f_X(x)$ whatever f , then its $F_X(x) = U$ can be thought as a new random variable $U \sim \text{Unif}(0, 1)$ (this result is called *probability integral transform*).

Definition 10.1.1 (Inversion method). If our aim is to draw values from X , a solution with a two step procedure is as follows:

1. draw different values u_1, \dots, u_n from $\text{Unif}(0, 1)$;
2. compute $F_X^{-1}(u_1), \dots, F_X^{-1}(u_n)$ obtaining $x_1, \dots, x_n \sim f_X(x)$

Univariate	Multivariate
Inversion	Gibbs sampling
Accept-reject	Metropolis-Hasting
Sampling and resampling	...

Tabella 10.1: Sampling methods

Therefore this method requires we know F (and obtain its inverse).

Example 10.1.1. Let $X \sim \text{Exp}(\lambda)$, with known $F_X(x) = 1 - e^{-\lambda x} = u$; but imagine we are not able to draw from the exponential distribution. Knowing F we can obtain its inverse and apply inversion method. For the inverse:

$$\begin{aligned} 1 - u &= e^{-\lambda x} \\ \log(1 - u) &= -\lambda x \\ -\frac{1}{\lambda} \log(1 - u) &= x \end{aligned}$$

Following the algorithm:

1. we generate u_1, \dots, u_n from $\text{Unif}(0, 1)$
2. we calculate $x_1 = -\frac{1}{\lambda} \log(1 - u_1), \dots, x_n = -\frac{1}{\lambda} \log(1 - u_n)$

Remark 327. From a practical point of view it's not very useful:

1. it's already implemented for common distribution in statistical software: eg `rexp` uses this method;
2. there are very few rvs for which the pdf is known and is invertible.

10.1.2 Accept-reject method

Important remark 87 (Setup). We

- are interested in generating values from $\pi(x)$ which is the target distribution (not a known one eg exp, normal etc). It's known in part, analytically (eg we know at least the kernel concerning x , not necessarily integral-normalizing-to-1 constants), but we are not able to draw values from it.
- we choose $p(x)$, a perfectly-known distribution from which we can draw values.

Remark 328. One could invent a distribution by specifying the kernel (a function of x), setup the domain, and deriving the normalization constant by integration like this

$$1 = \int_{D_X} c \cdot (\text{kernel in } x) \, dx = c \int_{D_X} (\text{kernel in } x) \, dx = \dots$$

Immagininig we're not able to solve the integral and don't know or we don't know to generate values from this distribution. Then we can use the accept-reject method.

Definition 10.1.2 (Accept-reject method). The algorithm is the following:

1. draw a value x from the proposal $p(x)$
2. draw a value u from $\text{Unif}(0, 1)$

3. check if

$$u < \frac{\pi(x)}{M \cdot p(x)} \quad (10.1)$$

where for $\pi(x)$ and $p(x)$ we mean densities and M is a positive constant (so overall the right hand ratio is positive). M has to be fixed in advance, such that:

$$\pi(x) < M \cdot p(x), \quad \forall x \in D_X$$

One should check this condition

4. if 10.1 is true then *accept* x , if false then *reject* x
5. you repeat from 1, until you have enough elements for the application's need

Important remark 88. Some remarks:

- accept and reject because the rule specify when keeping our simulation as suitable value or not
- first of all we should choose the proposal p . How we should choose p ? First we should know something about the target:
 1. know the domain space D_X of the target (eg if one wants positive values or values between $-\infty$ and $+\infty$). *The proposal should respect the domain space.*
 2. if we know that the target is symmetric (or asymmetric), the proposal should be symmetric (respectively asymmetric) as well.
- regarding M : its said that M should guarantee that the ratio

$$\frac{\pi(x)}{M \cdot p(x)}$$

is a value between 0 and 1, since it's compared with a draw from $\text{Unif}(0, 1)$; so we should choose M high enough. So from here the condition to be checked

$$\pi(x) < M \cdot p(x), \quad \forall x \in D_X$$

If this condition is not satisfied, the method doesn't work.

Since we don't know what is the target so it's difficult to practically choose M :

- an *option in practice* is to choose M very large (eg 1000). in this case i'm quite sure the inequality will be respected.
- but if it's too large we will have a method where acceptance is very rare. So the method could be slow (method has a tradeoff).

Proof of accept-reject. We aim is to prove that what we generate is a realization of the distribution of interest, and in math terms that the density f of the value we accept x (conditioned to being accepted) is equal to the target

$$f\left(x \middle| u < \frac{\pi(x)}{Mp(x)}\right) = \pi(x)$$

In order to prove that we start by writing/expanding the conditional density, which is the ratio between joint density and at denominator the probability of conditioning. Thus by definition we have:

$$f\left(x \middle| u < \frac{\pi(x)}{Mp(x)}\right) = \frac{\mathbb{P}\left(x \cap u < \frac{\pi(x)}{Mp(x)}\right)}{\mathbb{P}\left(u < \frac{\pi(x)}{Mp(x)}\right)}$$

Now, given that the intersection can be written twofold (conditioning on the first or the second event)

$$\mathbb{P}(A \cap B) = \mathbb{P}(A|B) \mathbb{P}(B) = \mathbb{P}(B|A) \mathbb{P}(A)$$

we can rewrite the numerator, which is a joint density, this way:

$$\begin{aligned} \frac{\mathbb{P}\left(x \cap u < \frac{\pi(x)}{Mp(x)}\right)}{\mathbb{P}\left(u < \frac{\pi(x)}{Mp(x)}\right)} &\stackrel{(1)}{=} \frac{p(x) \cdot \mathbb{P}\left(u < \frac{\pi(x)}{Mp(x)} \middle| x\right)}{\int_{D_x} p(x) \cdot \mathbb{P}\left(u < \frac{\pi(x)}{Mp(x)} \middle| x\right) dx} \stackrel{(2)}{=} \frac{p(x) \cdot \frac{\pi(x)}{Mp(x)}}{\int_{D_x} p(x) \frac{\pi(x)}{Mp(x)} dx} \\ &\stackrel{(3)}{=} \frac{\pi(x)}{\underbrace{\int_{D_x} \pi(x) dx}_{=1}} \stackrel{(4)}{=} \pi(x) \end{aligned}$$

where:

- (1) because we write the denominator as well as (integral of) joint density, given the fact that it's a marginal density so the way to do it is $\int_{D_Y} f(x, y) dy = f(x)$.
Furthermore informally/put another way (Luca's view) it seems basically the theorem of total probability for $\mathbb{P}\left(u < \frac{\pi(x)}{Mp(x)}\right)$;
- (2) remembering that u is coming from a $\text{Unif}(0, 1)$ distribution; therefore the probability that a $\text{Unif}(0, 1)$ is lower than a constant $c \in [0 - 1]$, is the constant c itself (here our constant is $\frac{\pi(x)}{Mp(x)}$);
- (3) we moved constant M and simplified;
- (4) the denominator is the integral of the target distribution over the domain so it must be 1.

□

Important remark 89. As said it works iff M is carefully chosen to make the ratio $\frac{\pi(x)}{Mp(x)}$ between 0 and 1: if it doesn't, the property of the uniform used at (2) in proof above doesn't work any more.

Acceptance probability A things important for the algorithm to be computed: we have said that M the probability of acceptance of drawn values from the proposal distribution.

Idea is that if you take M too large you will accept few values (we will see in lab), so there is an inverse correlation between these two quantities. But again

it's important that M is large enough to make the ratio is lower than 1.

This is the tradeoff: now we view this tradeoff in math terms. We want to compute the acceptance probability.

Let's call acceptance probability of the algorithm *alpha*; it's defined as

$$\begin{aligned}\alpha &= \mathbb{P}(\text{accepted}) \stackrel{(1)}{=} \mathbb{P}\left(U < \frac{\pi(x)}{Mp(x)}\right) \\ &\stackrel{(2)}{=} \int_{D_x} p(x) \mathbb{P}\left(U < \frac{\pi(x)}{Mp(x)} \mid X = x\right) dx \\ &= \int_{D_x} p(x) \frac{\pi(x)}{Mp(x)} dx = \frac{1}{M} \underbrace{\int_{D_x} \pi(x) dx}_{=1} \\ &= \frac{1}{M}\end{aligned}$$

where in:

- (1) we wrote capital U (meaning $\text{Unif}(0, 1)$) since it's not a single extraction but a random variable that originate a probability
- (2) we rewrite as integral of joint probability (as made for the accept reject method proof), or more explicitly here

$$\int_{D_y} f(x, y) dy = \int_{D_y} f(y) f(x|y) dy$$

TODO: to be reported above maybe

So given that $\alpha = \frac{1}{M}$ we have a very precise relation regarding the trade off we talked above.

Important remark 90. Observe: it works even if we don't know fully the target, but we know the target unless a normalization constant.

What about a situation in which the target can be decomposed in a kernel part $k(x)$ times a constant, that is in situations such as $\pi(x) = k(x) \cdot c$ (where we know $k(x)$ but not c)? Eg we want to generate a rv from a distribution with kernel: $\exp(-\log(x)) \cdot c$ (we don't know c).

This doesn't matter since the method works in any case: we repeat the proof with a different perspective.

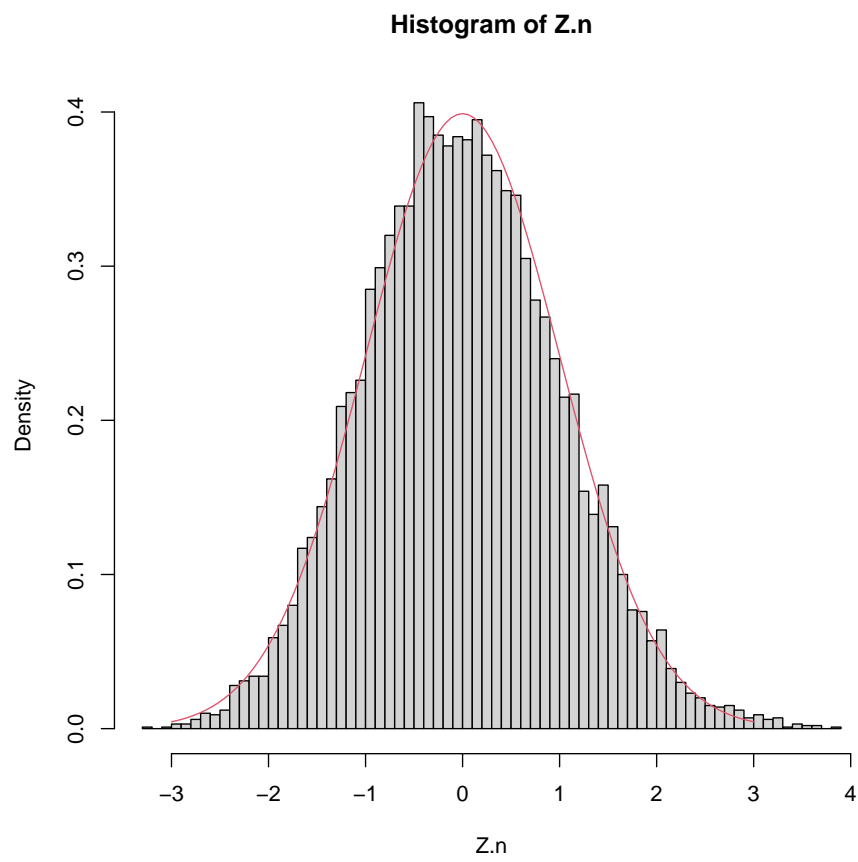
Dimostrazione. In proof the difference is at numerator of where we substituted $k(x) \cdot c$ instead of $\pi(x)$. So we aim is to prove that $f\left(x|u < \frac{k(x) \cdot c}{Mp(x)}\right) = \pi(x)$:

$$\begin{aligned}f\left(x|u < \frac{k(x) \cdot c}{Mp(x)}\right) &= \frac{\mathbb{P}\left(x \cap u < \frac{k(x) \cdot c}{Mp(x)}\right)}{\mathbb{P}\left(u < \frac{k(x) \cdot c}{Mp(x)}\right)} = \frac{p(x) \mathbb{P}\left(u < \frac{k(x) \cdot c}{Mp(x)} \mid x\right)}{\int_{D_x} p(x) \mathbb{P}\left(u < \frac{k(x) \cdot c}{Mp(x)} \mid x\right) dx} \\ &= \frac{p(x) \frac{k(x) \cdot c}{Mp(x)}}{\int_{D_x} p(x) \frac{k(x) \cdot c}{Mp(x)} dx} = \frac{k(x) \cdot c}{\underbrace{\int_D k(x) \cdot c dx}_{=1}} = k(x) \cdot c \\ &= \pi(x)\end{aligned}$$

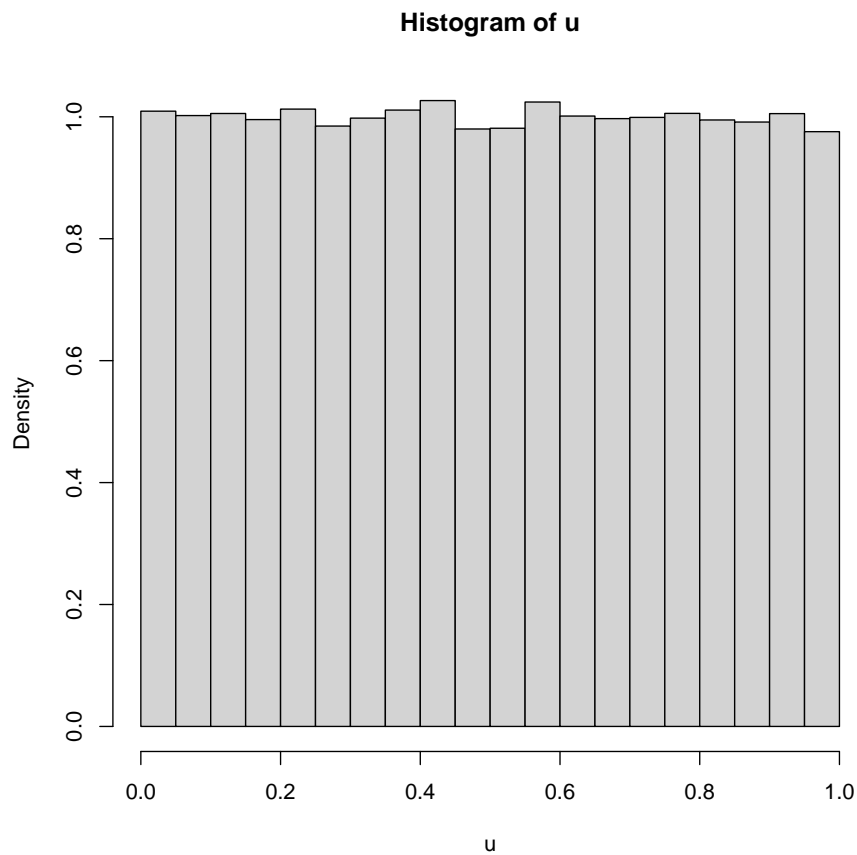
So we proved that we can use the algorithm even without knowing the normalization constant. \square

10.2 R exercises

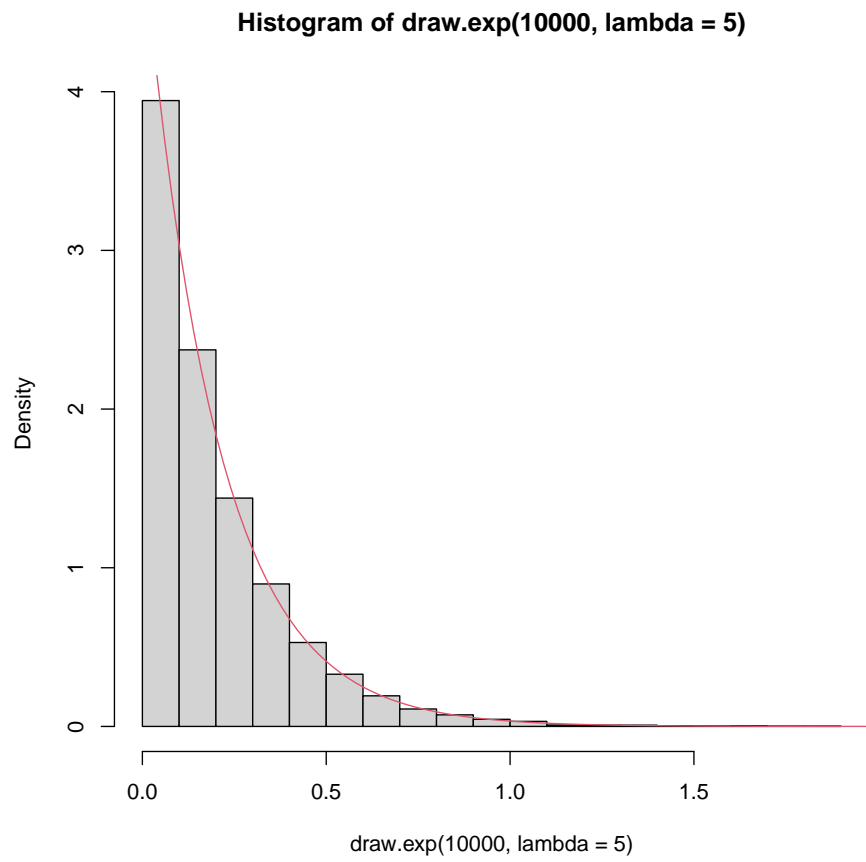
10.2.1 CLT



10.2.2 Inversion method

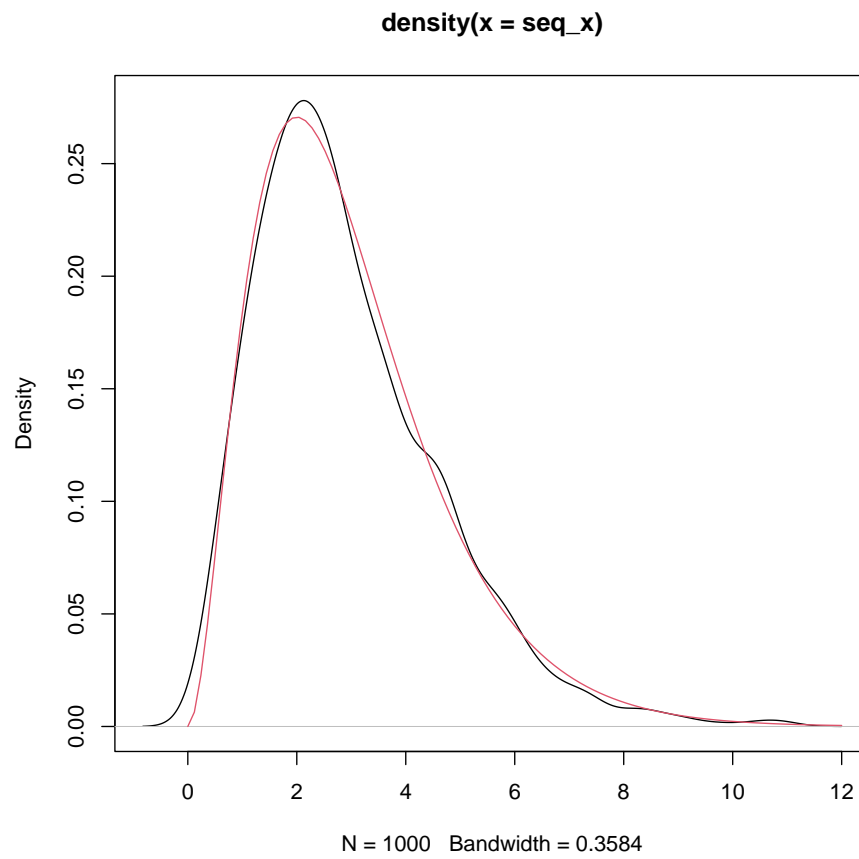


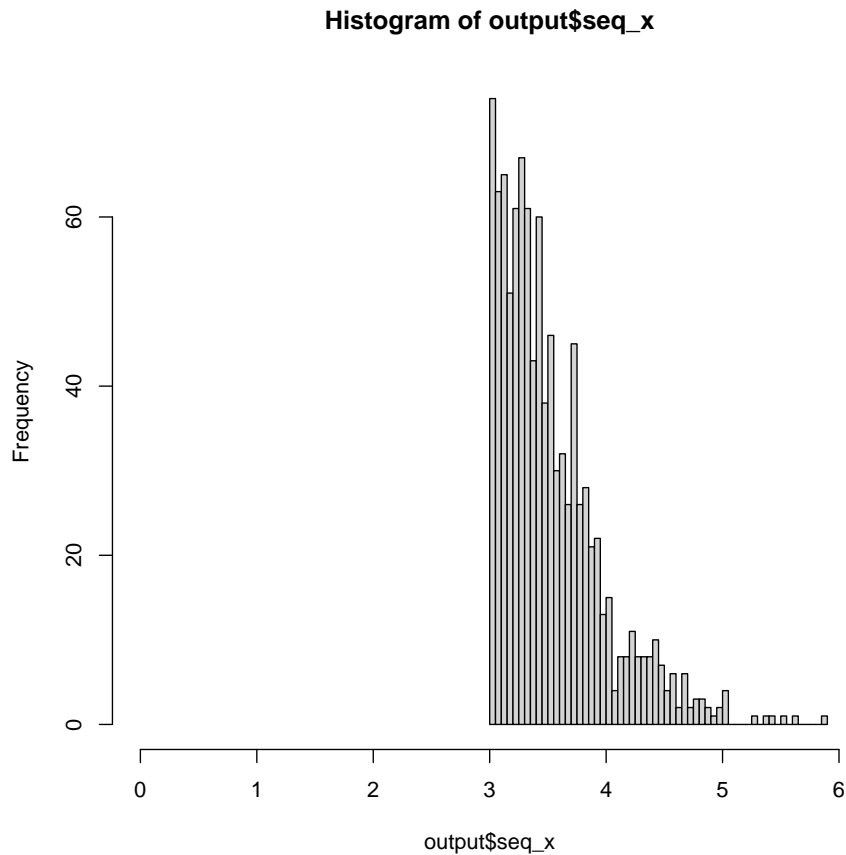
```
## [1] 0.00769923 0.05389461 0.37726227 0.64083592 0.48585141 0.40095990
## [7] 0.80671933 0.64703530 0.52924708 0.70472953 0.93310669 0.53174683
## [13] 0.72222778 0.05559444 0.38916108 0.72412759 0.06889311 0.48225177
## [19] 0.37576242 0.63033697 0.41235876 0.88651135 0.20557944 0.43905609
## [25] 0.07339266 0.51374863 0.59624038 0.17368263 0.21577842 0.51044896
## [31] 0.57314269 0.01199880 0.08399160 0.58794121 0.11558844 0.80911909
## [37] 0.66383362 0.64683532 0.52784722 0.69493051 0.86451355 0.05159484
## [43] 0.36116388 0.52814719 0.69703030 0.87921208 0.15448455 0.08139186
## [49] 0.56974303 0.98820118 0.91740826 0.42185781 0.95300470 0.67103290
## [55] 0.69723028 0.88061194 0.16428357 0.14998500 0.04989501 0.34926507
## [61] 0.44485551 0.11398860 0.79792021 0.58544146 0.09809019 0.68663134
## [67] 0.80641936 0.64493551 0.51454855 0.60183982 0.21287871 0.49015098
## [73] 0.43105689 0.01739826 0.12178782 0.85251475 0.96760324 0.77322268
## [79] 0.41255874 0.88791121 0.21537846 0.50764924 0.55354465 0.87481252
## [85] 0.12368763 0.86581342 0.06069393 0.42485751 0.97400260 0.81801820
## [91] 0.72612739 0.08289171 0.58024198 0.06169383 0.43185681 0.02299770
## [97] 0.16098390 0.12688731 0.88821118 0.21747825
```



10.2.3 Accept-reject

```
## [1] 101852
```



10.3 Other simulation based stuff

10.3.1 Definizioni

Definition 10.3.1 (Simulazione di Monte Carlo). Un algoritmo che genera molteplici campioni di dati utilizzando un processo generativo controllato (e composto matematicamente da una parte fissa ed una aleatoria) per simulare una popolazione/sue caratteristiche; questo al fine di esaminare i pattern rinvenibili nei campioni generati

Definition 10.3.2 (Resampling methods). Metodi che estraggono molteplici campioni da dati osservati, piuttosto che da un meccanismo teorico. Vi appartengono, tra l'altro test di permutazione, jackknife, bootstrap.

Definition 10.3.3 (Test di permutazione). I campioni sono generati riassegnando casualmente i pazienti a differenti gruppi (da porre a confronto); si crea di fatto una situazione da ipotesi nulla (no differenze tra gruppi) e una distribuzione dello stimatore di interesse sotto tale ipotesi. Si confronta poi la stima ottenuta nel campione con tale distribuzione sotto ipotesi nulla per inferire eventualmente rifiutando la stessa.

Definition 10.3.4 (Jackknife). I campioni sono generati iterativamente rimuovendo una osservazione o un gruppo di osservazioni dai dati disponibili; si creano di fatto subset del dataset di partenza. Non approfondiremo, focalizzandoci sul bootstrap.

Definition 10.3.5 (Bootstrap). I campioni sono generati estraendo con ripetizione osservazioni o gruppi di osservazioni dai dati disponibili.

10.3.2 Metodo di monte carlo

10.3.2.1 Calcolo di campioni

Mediante metodo di monte carlo, occorre:

- individuare il test e il criterio di rifiuto dell'ipotesi nulla;
- impostare la sequenza ampiezze campionarie testate;
- per ciascuna di essa creare un numero sufficiente di campioni (1000-10000), generati sotto ipotesi alternativa;
- verificare a che numerosità la probabilità con cui il test rifiuta l'ipotesi nulla supera la potenza pre-fissata (es 80%)

10.3.2.2 Validazione stimatori

Definition 10.3.6 (Correttezza). Per un dato sample size, il valore atteso dello stimatore deve essere il più possibile vicino al valore vero della popolazione.

Definition 10.3.7 (Efficienza). Per un dato sample size, la variabilità dello stimatore deve essere il minore possibile.

Definition 10.3.8 (Consistenza). In campioni via via crescenti, il valore atteso si avvicina al valore vero per effetto della diminuzione della variabilità della stima.

Si ha che:

- **mean1** è sia unbiased che consistente
- **mean2** è biased a meno che non sia applicato a vettori di 10 unità (e mano a mano che il campione aumenta il bias aumenta verso il basso); è anche inconsistente
- **mean3** è unbiased ma non consistente
- **mean4** è biased ma consistente

Validazione di correttezza Nel seguente adoperiamo due metriche per la validazione della correttezza.

Vogliamo minimizzare entrambe: sceglieremo lo stimatore con il valore minimo di ABS o MSE. In entrambi i casi è **mean1**

```
##          m1          m2          m3          m4
## 0.1110486 3.9988007 2.4970018 0.1398323
##          m1          m2          m3          m4
## 0.01900102 15.99419687 6.25870320 0.02968860
```

Validazione di consistenza Dobbiamo calcolare la stessa metrica (ABS o MSE) aumentando via via il campione (passiamo da 50 a 500) e verificando che lo stimatore tende al valore effettivo (attraverso una riduzione di ABS o MSE). Si nota appunto che gli stimatori consistenti sono il primo e il quarto

```
##          ss sim          m1          m2          m3          m4
## [1,] 50    1 4.854466 0.9922330 2.480582 4.953536
## [2,] 50    2 4.745560 0.9627959 2.406990 4.842408
## [3,] 50    3 4.967219 0.9424115 2.356029 5.068591
## [4,] 50    4 4.951787 0.9873412 2.468353 5.052844
## [5,] 50    5 4.905349 0.9348728 2.337182 5.005458
## [6,] 50    6 4.890084 0.9916527 2.479132 4.989881
##          ss sim          m1          m2          m3          m4
## [9995,] 500 995 5.093528 0.09791140 2.447785 5.103736
## [9996,] 500 996 5.124492 0.10559905 2.639976 5.134762
## [9997,] 500 997 4.881977 0.09614737 2.403684 4.891761
## [9998,] 500 998 4.970079 0.09718721 2.429680 4.980039
## [9999,] 500 999 4.974599 0.10641193 2.660298 4.984568
## [10000,] 500 1000 4.966426 0.09617232 2.404308 4.976378
##          ss          m1          m2          m3          m4
## 50    50 0.11433367 3.997841 2.494604 0.14063697
## 100 100 0.07934855 4.498721 2.493603 0.09099866
## 150 150 0.06818283 4.667379 2.505341 0.07328250
## 200 200 0.05687781 4.749895 2.498951 0.06104778
## 250 250 0.04948347 4.800259 2.503243 0.05142773
## 300 300 0.04461717 4.833104 2.496567 0.04742551
## 350 350 0.04278968 4.857122 2.499632 0.04445638
## 400 400 0.03968947 4.875146 2.502921 0.04105650
## 450 450 0.03855951 4.888710 2.495975 0.03962197
## 500 500 0.03696599 4.900024 2.500612 0.03769749
```

Validazione stimatori: variabilità e coverage

Definition 10.3.9 (Coverage probability). La proporzione di campioni simulati in cui l'intervallo di confidenza di un parametro stimato sui dati include il vero parametro della popolazione (utilizzato nella generazione dei campioni stessi).

Remark 329. Una volta individuato uno stimatore e intervallo di confidenza di interesse (o più se li si vuole confrontare):

- simulare un numero sufficiente di campioni (es 1000) sulla base del parametro scelto (es proporzione pari a 0.1) e l'ampiezza campionaria scelta (es 100)

- valutare se l'intervallo di confidenza nominale al 95% include il parametro oggetto di stima in che percentuale:
 - se molto inferiore a 95 si tratta di uno stimatore eccessivamente liberale e pericoloso (variabilità della stima è una sottostima);
 - se molto superiore a 95 si tratta di uno stimatore conservativo.

Remark 330. In una seconda fase per valutare il coverage in varie situazioni si può:

- far variare l'ampiezza campionaria di ciascun per valutare il coverage dello stimatore a diversi sample size;
- far variare il parametro di interesse per valutare il coverage in situazioni differenti (es percentuali tra 0 e 1).

Example 10.3.1 (Coverage dell'intervallo di Clopper Pearson). Facciamo una simulazione con $\pi = 0.1$ basata su 1000 campioni da 100 pazienti ciascuno. Si conclude che il coverage è buono (volendo si può aggiungere l'intervallo di confidenza della stima di coverage stessa, la cui ampiezza dipende comunque dal numero di ripetizioni adottate).

```
## [1] 0.957
```

10.3.3 Test di permutazione/randomizzazione

In questi test invece di assumere una forma particolare della distribuzione nulla al fine di effettuare il test, si utilizzano i dati osservati per crearne una. Se il nostro obiettivo è verificare la differenza tra due trattamenti si riallocano casualmente l'outcome (e rompendo il legame bivariato) e calcolando la stima di interesse. La distribuzione delle stime in campioni ripetuti costituirà la nostra distribuzione dello stimatore sotto ipotesi nulla. Abbiamo due varianti.

Definition 10.3.10 (Test di permutazione). Costruiscono la distribuzione in maniera esatta effettuando tutte le permutazioni possibili

Example 10.3.2. Se ad esempio abbiamo due gruppi di 5 pazienti ciascuno, il numero di campioni da generare è $\binom{10}{5}$

Remark 331. Il problema è che in campioni grandi questa soluzione diviene intensiva e si può optare per ...

Definition 10.3.11 (Randomization test). Si costruisce la distribuzione nulla effettuando un numero alto di permutazioni (non tutte).

10.3.3.1 Test di permutazione

```
## treatment outcome
## 1      1 5.2889932
## 2      1 5.5227244
## 3      1 5.7360698
```

```
## 4      0 -0.6683198
## 5      0  1.9418637
## 6      0  0.9191380
## [1] 20
##      [,1] [,2] [,3] [,4] [,5] [,6]
## [1,]    1    1    1    0    0    0
## [2,]    1    1    0    1    0    0
## [3,]    1    0    1    1    0    0
## [4,]    0    1    1    1    0    0
## [5,]    1    1    0    0    1    0
## [6,]    1    0    1    0    1    0
##      1      1      1      1      1      1
## 4.7850352 0.5154421 0.6576724 0.8134931 2.2555645 2.3977947 2.5536155
##      1      1      1      1      1      1
## -1.7159776 -1.8717983 -1.5737473 -2.2555645 -2.3977947 1.8717983 1.7159776
##      1      1      1      1      1      1
## -2.5536155 1.5737473 -0.8134931 -0.6576724 -0.5154421 -4.7850352
## [1] 0.05
```

Si mostra come se l'ipotesi nulla fosse vera 4.78 sarebbe il più grande dei valori e si osserverebbe in un campione su venti, corrispondendo ad un p-value esatto pari a 0.05.

10.3.3.2 Test di randomizzazione

Cambiano solo due cose, la numerosità del campione di partenza e il meccanismo di generazione delle ripetizioni. Supponendo di avere un campione di 100 pazienti, 50 trattati e 50 controlli, il numero di permutazioni sarebbe 1.0089134×10^{29} , il che può essere poco praticabile. Optiamo invece per costruire la nulla sulla base di 1000 ripetizioni.

```
## [1] 0
```

10.3.4 Bootstrap

L'idea è avendo un singolo campione di ipotizzare che sia rappresentativo della popolazione dal quale è estratto e che possa essere utilizzato per generare altri campioni. Pertanto

- la stima puntuale di un parametro di interesse si ottiene dal nostro campione (la stima può riguardare anche più parametri contemporaneamente, come i coefficienti di un modello);
- per ottenere la distribuzione dello stimatore (e una stima di variabilità es intervallo di confidenza) si procede ad effettuare estrazioni *con reinserimento* al fine di generare altri campioni *aventi la stessa numerosità* del campione di partenza (se non fosse con reinserimento si otterrebbe lo stesso campione di partenza);

- si applica poi lo stimatore per ottenere la stima su quel campione generato e si ricostruisce la distribuzione dello stimatore ripetendo il processo più volte; il beneficio è che lo stimatore può essere complesso a piacere ma la stima della sua distribuzione campionaria è sempre fattibile.

Example 10.3.3 (Esempio di una media). Il metodo non parametrico è solitamente più generale, anche per distribuzioni dello stimatore che non siano normali (come in questo caso); affinché funzioni bene costruire molti campioni bootstrap.

```
##      2.5%    97.5%
## 3.647167 4.483127
## [1] 0.2065029
## [1] 0.2142785
## [1] 3.641902 4.481858
```

Example 10.3.4 (Stimatore generico e utilizzo di R). Per uno stimatore generico facciamo uso delle facilities di R, ossia delle funzioni `boot` e `boot.ci` del pacchetto `boot`. Facciamo un esempio con la correlazione tra educazione e fertilità utilizzando il dataset `swiss`

```
## [1] -0.6637889
##
## ORDINARY NONPARAMETRIC BOOTSTRAP
##
##
## Call:
## boot::boot(data = swiss, statistic = boot_f, R = 10000)
##
##
## Bootstrap Statistics :
##      original    bias    std. error
## t1*  -0.6637889  0.0191199   0.1086864
##           [,1]
## [1,] -0.6601276
## [2,] -0.4821279
## [3,] -0.5415851
## [4,] -0.5500972
## [5,] -0.7330102
## [6,] -0.5844682

## Warning in boot::boot.ci(boot_res):  varianze bootstrap necessarie
## per intervalli studentizzati

## BOOTSTRAP CONFIDENCE INTERVAL CALCULATIONS
## Based on 10000 bootstrap replicates
##
## CALL :
## boot::boot.ci(boot.out = boot_res)
##
```

```
## Intervals :
## Level      Normal      Basic
## 95%   (-0.8959, -0.4699 )   (-0.9496, -0.5258 )
##
## Level      Percentile      BCa
## 95%   (-0.8018, -0.3779 )   (-0.8186, -0.4246 )
## Calculations and Intervals on Original Scale
```

Example 10.3.5 (Più stimatori contemporaneamente). La cosa che cambia è che la funzione utilizzata restituisce un vettore di numeri

```
##
## ORDINARY NONPARAMETRIC BOOTSTRAP
##
## Call:
## boot::boot(data = db, statistic = boot_f, R = 1000)
##
##
## Bootstrap Statistics :
##      original      bias      std. error
## t1* 1.729245  3.692618e-04 0.201888150
## t2* 3.005978 -3.685578e-06 0.003410122
##      [,1]      [,2]
## [1,] 1.958497 3.003515
## [2,] 1.983657 3.002827
## [3,] 1.570135 3.007322
## [4,] 1.474370 3.009072
## [5,] 1.874145 3.002240
## [6,] 1.898579 3.001049

## Warning in boot::boot.ci(boot_res, index = 1):  varianze bootstrap
## necessarie per intervalli studentizzati

## BOOTSTRAP CONFIDENCE INTERVAL CALCULATIONS
## Based on 1000 bootstrap replicates
##
## CALL :
## boot::boot.ci(boot.out = boot_res, index = 1)
##
## Intervals :
## Level      Normal      Basic
## 95%   ( 1.333,  2.125 )   ( 1.339,  2.119 )
##
## Level      Percentile      BCa
## 95%   ( 1.340,  2.119 )   ( 1.376,  2.155 )
## Calculations and Intervals on Original Scale

## Warning in boot::boot.ci(boot_res, index = 2):  varianze bootstrap
## necessarie per intervalli studentizzati
```



```
## BOOTSTRAP CONFIDENCE INTERVAL CALCULATIONS
## Based on 1000 bootstrap replicates
##
## CALL :
## boot::boot.ci(boot.out = boot_res, index = 2)
##
## Intervals :
## Level      Normal          Basic
## 95%   ( 2.999,  3.013 )   ( 3.000,  3.013 )
##
## Level      Percentile      BCa
## 95%   ( 2.999,  3.012 )   ( 2.999,  3.012 )
## Calculations and Intervals on Original Scale
```