

San Pellegrino 03/09/19 – Summer School “La Matematica oggi”

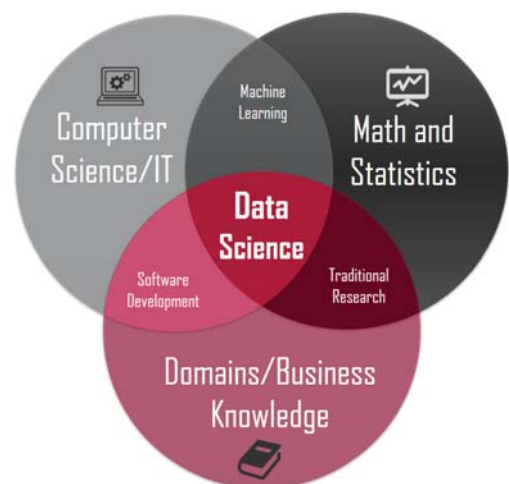
La statistica e la scienza dei dati

Mauro Gasparini Dipartimento di Scienze Matematiche “G.L. Lagrange”
Politecnico di Torino
mauro.gasparini@polito.it
<http://calvino.polito.it/~gasparini>

1

Statistica e Data Science

Evoluzione di un diagramma di Venn dovuto a Conway (2013).



<https://towardsdatascience.com/>

2

Un problema statistico: la classificazione

Decidere se un oggetto appartiene a una popolazione positiva oppure a una popolazione negativa tramite una serie di **rilevazioni** sull'oggetto. Esempi:

1. diagnosticare se un paziente sano o malato
2. assegnare a un segnale elettrico un valore binario (vero o falso)
3. giudicare un imputato è colpevole o innocente
4. decidere se un essere vivente è un animale o un vegetale
- ω 5. valutare se un debitore sarà insolvente o no

Un problema statistico è un problema di induzione

Cosa vuol dire **induzione**?

Induzione "scandalo della filosofia".

Aristotele, Bacon, Hume, Mill.

David Hume 1711–1776

Filosofo.
Il problema dell'induzione.
Il concetto di causa.



5

Cigni neri

“rara avis in terris, nigroque simillima cygno”
(Giovenale)

I cigni neri furono effettivamente scoperti in Australia alla fine del Seicento.

Recentemente, l'espressione “cigno nero” ha assunto un significato più specifico di evento imprevedibile e devastante, specialmente in economia e in politica (Taleb).



6

Il gioco della classificazione dei viventi

Torniamo al problema decidere se un essere vivente è un animale o un vegetale.

Poniamoci il problema di costruire delle domande a risposta binaria (per semplicità) sulla base delle quali prendere una decisione (a sua volta binaria, “dico vegetale” oppure “dico animale”).

Esempio: “E' verde”?

7



Idea: usare la probabilità condizionata

La domanda “E’ verde” non è perfetta, ma è “buona”. Perché?

Supponiamo di potere quantificare le seguenti **probabilità**

$$P(\textit{Animale}) = 0.45$$

$$P(\textit{Vegetale}) = \text{????}$$

$$P(\textit{Verde}|\textit{Animale}) = 0.05$$

$$P(\textit{Verde}|\textit{Vegetale}) = 0.70$$

Che tipo di probabilità sono queste? Cosa significa quella barra verticale?

∞

Riusciamo a vedere un motivo per cui la domanda è “buona”? Vediamo.

Aggiornamento bayesiano

Se osserviamo “Verde”, possiamo usare la probabilità condizionata inversa

$$\begin{aligned} P(\textit{Animale}|\textit{Verde}) &= \frac{P(\textit{Animale} \cap \textit{Verde})}{P(\textit{Verde})} \\ &= \frac{P(\textit{Animale}) P(\textit{Verde}|\textit{Animale})}{P(\textit{Animale}) P(\textit{Verde}|\textit{Animale}) + P(\textit{Vegetale}) P(\textit{Verde}|\textit{Vegetale})} \\ &= \frac{0.45 \times 0.05}{0.45 \times 0.05 + 0.55 \times 0.70} = 0.05521472 \end{aligned}$$

e conseguentemente

$$P(\textit{Vegetale}|\textit{Verde}) = 1 - 0.05521472 = 0.9447853$$

E’ un esempio di **aggiornamento bayesiano**.

∞

Thomas Bayes 1702–1761

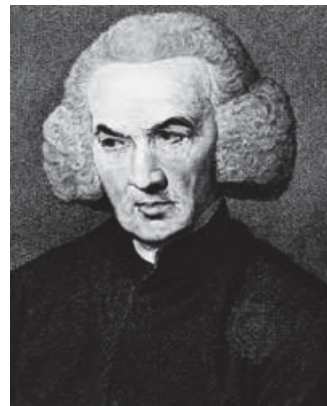
Teologo e matematico.
Teorema di Bayes.
Inconsapevole capostipite.
Pastore non conformista
(congregazionista).



10

Richard Price 1723 –1791

Amico e divulgatore di Bayes.
Contributi alle → Scienze Attuariali.
Pastore unitariano, radicale.
Amico della rivoluzione americana.



11

Rincariamo la dose.

Per essere sicuri, aggiungiamo un'altra domanda. Quale?

12

Conosciamo le probabilità?

Che possiamo conoscere o meno tutte le probabilità coinvolte è stato oggetto di lunghe discussioni negli ultimi due secoli tra statistici [bayesiani](#) e non [bayesiani](#).

13

Un bayesiano moderno: Bruno de Finetti 1906-1985

Fondatore della rinascita neobayesiana.
Scambiabilità.
Teorema di rappresentazione.
Schema della scommessa.
Additività finita.
Soggettivista estemo.
‘La probabilità non esiste’.
Fascista in gioventù, poi radicale.
Economia, Didattica della matematica.



14

Un non-bayesiano: Ronald A. Fisher 1890–1962

Padre della Statistica moderna.
Padre della → Genetica di Popolazione.
Moderna sintesi Darwin-Mendel.
 t di Student.
Regressione, ANOVA, DOE.
Randomizzazione.
Analisi multivariata (e.g. discriminante).



15

Quantificazione delle probabilità?

Per esempio, abbiamo quantificato $P(\textit{Animale}) = 0.45$ ma potrebbe essere molto minore.

Un alieno cosa userebbe?

Si accettano suggerimenti....

16

Pierre-Simon Laplace 1749-1827

Matematico e astronomo.

Ottenne indipendentemente il teorema di Bayes.

A priori uniforme.

Thorie analytique des probabilitis.

Regola di successione.

Minimi quadrati.

Determinismo.

Politicamente, un trasformista.



17

Quantificazione delle probabilità condizionate

E per quanto riguarda queste?

$$P(\text{Verde}|\text{Animale}) = 0.05$$

$$P(\text{Verde}|\text{Vegetale}) = 0.70$$

Le conosciamo?

Come farebbe un alieno a quantificarle?

Se l'alieno avesse a disposizione una serie di osservazioni analoghe, lo potrebbe fare. Come?

18

Esempio di training set

	verde?	altro?	vero		verde?	altro?	vero
1	1	1	vegetale	26	0	0	vegetale
2	0	1	vegetale	27	1	0	vegetale
3	0	1	vegetale	28	1	1	vegetale
4	1	0	animale	29	0	0	animale
5	0	0	animale	30	0	1	animale
6	1	1	animale	31	1	1	animale
7	1	1	vegetale	32	1	0	animale
8	0	0	vegetale	33	0	0	animale
9	1	1	animale	34	1	1	animale
10	0	1	animale	35	1	0	vegetale
11	1	1	animale	36	1	0	vegetale
12	1	1	vegetale	37	1	1	animale
13	0	1	animale	38	0	1	vegetale
14	0	1	animale	39	0	1	vegetale
15	1	1	animale	40	0	0	vegetale
16	1	0	vegetale	41	0	0	vegetale
17	1	0	vegetale	42	1	1	animale
18	1	0	animale	43	1	1	animale
19	1	1	vegetale	44	0	0	animale
20	1	0	vegetale	45	1	0	animale
21	1	1	animale	46	1	1	animale
22	0	1	animale	47	0	0	vegetale
23	0	1	animale	48	0	1	vegetale
24	1	1	animale	49	1	0	animale
25	1	1	vegetale	50	1	1	animale

19

Il training set come database statistico

Il *training set* è un semplice database statistico regolare (rettangolare):

- sulle righe ci sono le unità statistiche, o **oggetti** di interesse
- sulle colonne ci sono le **variabili**, o caratteristiche, rilevate sugli oggetti.

Ho generato quello nella pagina precedente con il seguente codice scritto nel linguaggio R:

```
data <- cbind(  
  as.data.frame(matrix(rbinom(100,1,.6),50,2)),  
  sample(c("animale","vegetale"),50,replace=T)  
)  
colnames(data) <- c("verde?", "altro?", "vero")
```

20

Apprendimento supervisionato e non

Questo è un esempio di **apprendimento supervisionato**, cioè arricchito da un insieme di addestramento, o **training set** (a maggior ragione, anche conoscere le probabilità si pu considerare una forma di apprendimento supervisionato).

Un problema piú difficile si presenta quando disponiamo di un database parziale, senza vero. Allora dobbiamo essere noi a decidere come partizionarlo in due gruppi, un esempio di **apprendimento non supervisionato**.

	verde?	altro?		verde?	altro?
1	1	1	11	0	0
2	0	1	12	1	0
3	0	1	13	1	1
4	1	0	14	0	0
5	0	0	15	0	1
6	1	1	16	1	1
7	1	1	17	1	0
8	0	0	18	0	0
9	1	1	19	1	1
10	0	1	20	1	0

21

La maledizione della dimensionalità

Due domande spesso non sono sufficienti.

Si presenta, crudele, un problema di **dimensionalità**.

Quante sono le possibili assegnazioni di questi $n = 20$ elementi a due popolazioni (partizioni di n)?

Quanti sono gli m possibili **profili** con d domande?

22

Quanti sono le possibili **assegnazioni** dei profili?

La maledizione della dimensionalità: esempio

Esercizio: scrivere tutti i possibili profili di risposta a tre domande binarie.

Calcolare tutte le possibili assegnazioni di questi profili.

23

Quale metodo usare per stimare le probabilità ?

Ammesso che le domande e quindi i profili siano un numero maneggevole, rimane il problema di come stimare le probabilità dei profili.

Una assunzione semplificatoria, che abbiamo fatto anche per il gioco della classificazione dei viventi, è la [indipendenza locale](#) delle risposte.

Una volta assunta questa indipendenza, si possono stimare le probabilità dei profili tramite il metodo della [massima verosimiglianza](#), dovuto a R.A.

Fisher. E' uno dei concetti più importanti della [Statistica](#) nonbayesiana moderna.

24

Metodi basati sulle distanze

Se invece la maledizione della dimensionalità è insuperabile, ci arrendiamo e ricorriamo a metodi più euristici fondati sulla [distanza](#) tra unità statistiche (o, specularmente, sulla loro [vicinanza](#), o [similarità](#)).

Una volta scelta una distanza tra due unità statistiche, cerchiamo di individuare due classi tali che:

1. entro le classi, le distanze tra le unità statistiche siano, in media, basse;
2. tra le classi, le distanze tra le unità statistiche siano, in media, alte.

Questa è una [impostazione algoritmica](#) più vicina al [Machine Learning](#).

25

Buone distanze

La distanza tra due unità con lo stesso profilo dovrebbe essere nulla.

Qual è una buona distanza tra due unità con diversi profili?

Che distanza conoscete dalla scuola superiore?

26

Distanza euclidea

La distanza euclidea tra due **vettori** $x = (x_1, \dots, x_n)$ e $y = (y_1, \dots, y_n)$ è

$$d(x, y) = \sqrt{\sum_{i=1}^n (x_i - y_i)^2}$$

Cos'è in una dimensione $n = 1$?

In due dimensioni $n = 2$?

27 Per x e y binari?

Una tabella per x e y binari

Raccogliamo in una **tabella di contingenza** il numero di concordanze e discordanze delle diverse possibili paia

$x \backslash y$	0	1	
0	n_{00}	n_{01}	n_{1+} zeri in x
1	n_{10}	n_{11}	n_{2+} uni in x
	n_{+1} zeri in y	n_{+2} uni in y	n oggetti

28

Varie distanze per x e y binari

La distanza euclidea è la radice del **numero** di concordanze.

Il **coefficiente** di concordanza è

$$\frac{n_{01} + n_{10}}{n} = \frac{n_{00} + n_{11}}{n} = \text{proporzione di concordanze.}$$

Il **coefficiente** di Jaccard è invece

$$\frac{n_{01} + n_{10}}{n_{01} + n_{10} + n_{11}}.$$

29

Ce ne sono altri.

Esempio ecologico

Due zone ecologiche x e y , di uguale ampiezza (*quadrats*) vengono confrontate per vedere se hanno o meno la specie 1, la specie 2, ..., la specie n .

Le distanze viste misurano il grado di diversità biologica delle due zone.

30

Riassumendo

Ho cercato di far capire, tramite il *fil rouge* del gioco della classificazione dei viventi, somiglianze e differenze tra Statistica e Machine Learning.

La Statistica ha dei fondamenti teorici più profondi, e meglio formulati dal punto di vista matematico.

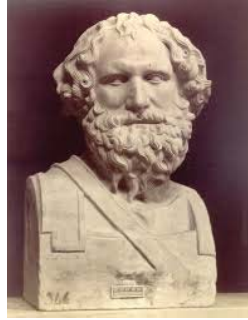
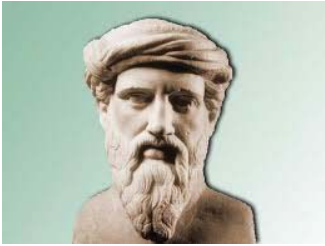
Il Machine Learning è più praticone ma può essere utile per quei casi in cui la Statistica fallisce (*Big Data*).

31 Entrambe costituiscono il cuore metodologico della *Data Science*.

Ne ho approfittato per menzionare diversi concetti matematici e non: probabilità, probabilità condizionata, induzione, database, classificazione supervisionata, classificazione non supervisionata, distanza.

Immagini di coda, se rimane tempo

Il più grande matematico dell'antichità.



... se la giocano.

33

“Il più grande matematico della modernità.”

E' abbastanza universalmente riconosciuto. Chi è?



34

Carl Friedrich Gauss 1777 –1855

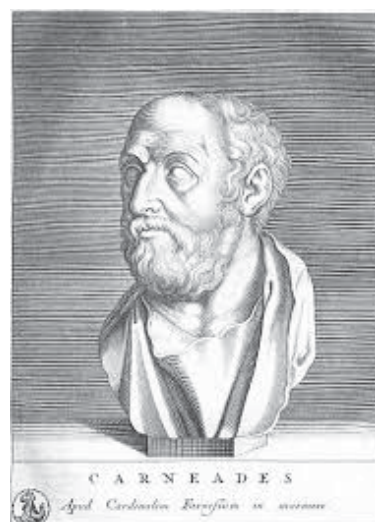
Algebra e Analisi.
Astronomia.
Probabilità.
Minimi quadrati.
→ Metrologia.



35

Carneade 214aC–129aC

Scettico greco.
Chi era costui?
Usò il termine “probabile”.



36

Cardano 1501–1576

Liber de ludo aleae, postumo.
Regole elementari della probabilità.
Applicazioni nel gioco dei dadi.
Coefficiente binomiale.
→ Combinatoria.



37

Pierre de Fermat 1601–1665

Corrispondenza con Pascal.
Matematico di primaria importanza.
Quello dell'Ultimo Teorema.
→ Teoria dei numeri.



38

Blaise Pascal 1623-1662

Corrispondenza con Fermat.
Geometra di primaria importanza.
Triangolo di Pascal (o di Tartaglia).
Apologeta cristiano.



39

Gottfried Achenwall 1719–1772

Giurista e filosofo.
Coni il termine *Statistica*.



40