



Laboratorio 4

A.A. 2025-2026

Enrico Longato





Dal lab 3: Risultati in presenza di NaN e inf

```
>> mean(elsa)
ans =
 Columns 1 through 13
  66.6290
 Columns 14 through 16
      NaN
>> nanmean(elsa)
ans =
 Columns 1 through 13
  66.6290
                               61.1979
                                                    3.9493 14.2644
 Columns 14 through 16
      Inf 135.8360 74.6381
>> nanstd(elsa)
  Columns 1 through 13
  10.4717
                      14.4517 15.3267
                                                                                           39.5524
                                                                                                    15.1565
                                                                                                             47.3685 125.1617
  Columns 14 through 16
      NaN 19.3198 11.3220
```

La media, in presenza di dati mancanti, viene restituita come dato mancante.

Razionale: "se non so uno dei numeri da sommare, non so neanche la media".

La media, in presenza di dati infiniti, viene restituita come infinito.

Razionale: "il numeratore dello stimatore è infinito perché uno dei numeri da sommare è infinito".

La SD, in presenza di dati infiniti, viene restituita come dato mancante.

Razionale: "poiché la media è infinita a sua volta, non so eseguire la differenza infinito (del dato) meno infinito (che è la media) che dovrei poi elevare al quadrato".



Dal lab 3: Indicizzazione booleana senza il find

1. Costruisco una matrice qualunque

2. Imponendo una condizione logica, creo una maschera booleana che ha 1 dove la condizione è vera e 0 altrove 3. Filtrando con la condizione logica direttamente (senza il find), trovo quello che mi aspetto, ovvero i valori che corrispondono alla condizione logica

4. NB: l'ordine sembra "strano", ma è quello del find lineare!

$$\begin{array}{c|c} 1 \\ 5 \end{array} \begin{array}{c|c} 2 \\ \hline 6 \end{array} \begin{array}{c|c} 3 \\ \hline 7 \end{array} \begin{array}{c|c} 4 \\ \hline 8 \end{array}$$





Dal lab 3: Negazione di una maschera

```
mask on v =
  2×4 logical array
>> ~mask on v
ans =
  2×4 logical array
```

Banalmente, la negazione elemento per elemento si fa con il simbolo tilde ~ (alt+126)





Dal lab 3: Funzioni con struttura f(qualcosa, dim)

Tantissime funzioni MATLAB di default agiscono colonna per colonna.

• Esempio: sum (A) dà un vettore riga contenente la somma di ciascuna colonna di A; se vogliamo un vettore colonna contenente la somma di ciascuna riga di A dobbiamo fare sum (A, 2)

Convenzionalmente, per una matrice 2D, "lavorare colonna per colonna" si dice "lavorare lungo la prima dimensione" (che, però, sarebbero le righe), mentre "lavorare riga per riga" si dice "lavorare lungo la seconda dimensione" (che, però, sarebbero le colonne).

• Da cui, sum (A) è come dire sum (A, 1); mentre sum (A, 2) si può dire solo così

L'ambiguità in termini fa, in effetti, confusione; dobbiamo inventarci uno stratagemma mnemonico. Eccone un paio (ma usate quello che volete):

- È come quando facciamo l'integrale: lo facciamo lungo x (la prima dimensione), ma il numero che esce dipende dal valore della funzione che leggiamo su y (la seconda dimensione).
- La dimensione da indicare a MATLAB è quella lungo cui un ideale ciclo for dovrebbe scorrere (per fare la somma di una colonna, il ciclo for scorrerebbe le righe, quindi la dimensione 1).





Dal lab 3: Eliminare i soggetti (= righe) con almeno un dato mancante

```
%% Calcolare la percentuale di nan per ogni variabile
% Calcoliamo il numero di nan per colonna
nan_counts = sum(isnan(elsa));

%% Rimuovere i soggetti che hanno, in almeno una delle variabili rimaste,
%% un valore mancante
% Sommo per righe con sum(..., 2) e guardo quando i nan sono più di 0

% Trasformiamo in percentuale
nan_percentages = 100 * nan_counts / size(elsa, 1);

%% Rimuovere le colonne con >20% dati mancanti
elsa_reduced = elsa(:, nan_percentages <= 20);

% Filtro
elsa_reduced = elsa(:, nan_percentages <= 20);

### Rimuovere i soggetti che hanno, in almeno una delle variabili rimaste,
% Sommo per righe con sum(..., 2) e guardo quando i nan sono più di 0

i_at_least_one_nan = sum(isnan(elsa_reduced), 2) > 0;
% Alternativamente, potevo fare any(isnan(elsa_reduced), 2)

#### Filtro
elsa_reduced = elsa(:, nan_percentages <= 20);
#### Piltro
elsa_reduced_filtered = elsa_reduced(~i_at_least_one_nan, :);</pre>
```

Esattamente speculare, con due ovvie differenze

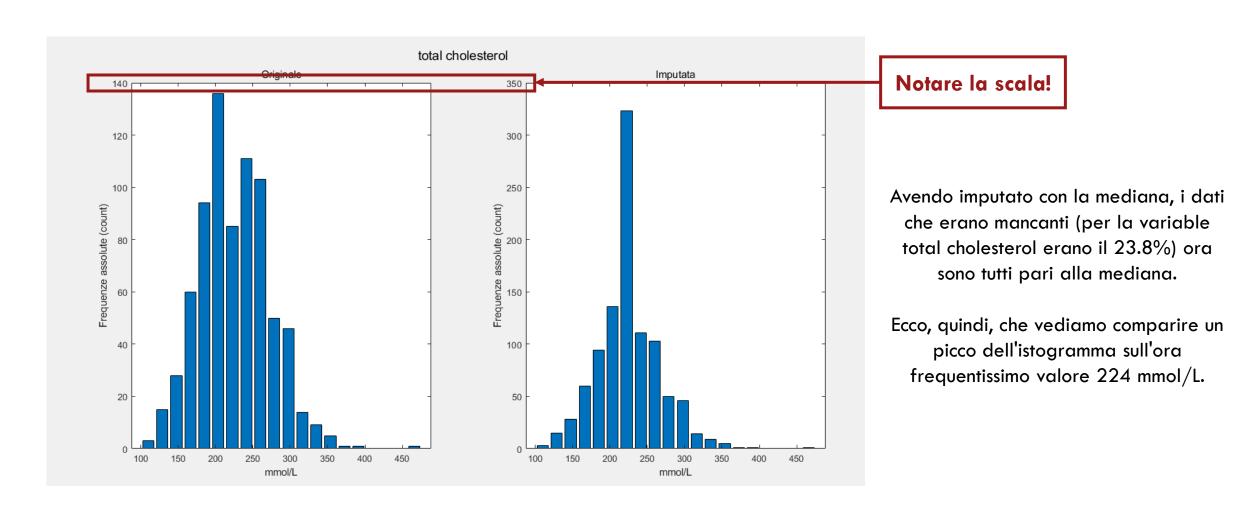
- 1. Per le colonne ho usato sum senza parametri, quindi come se fosse sum (isnan (elsa), 1); mentre per le righe ho usato sum con l'argomento dim = 2, cioè sum (isnan (elsa_reduced), 2)
- 2. Tengo tutte le righe (:, ...) quando filtro le colonne, mentre tengo tutte le colonne (..., :) quando filtro le righe.

Bonus: notare che per togliere soggetti, la logica di programmazione prevede di individuare quelli da tenere!





Dal lab 3: Esito dell'imputazione con la mediana







Laboratorio 4: Contenuti e obiettivi

- 1. Esercitazione "alla lavagna"
 - Caricamento ed esplorazione dati
 - z-test a un campione a due code
- 2. Esercizi da svolgere in autonomia (per superare la "paura del file bianco")
 - Statistiche descrittive e istogrammi
 - z-test a un campione a una coda
 - Variazione del livello di significatività
 - Riflessioni sull'ipotesi nulla
- 3. Ripasso di teoria (parte integrante del programma d'esame di teoria!)
 - Test statistici in pratica





Verifica di ipotesi

I test statici sono tentativi di confutare un'ipotesi (l'ipotesi nulla) che, di conseguenza, corrisponde al "contrario" di quello che vorremmo dimostrare.

- In altre parole, se ci interessa dimostrare che due fenomeni siano diversi, procediamo ipotizzando che siano uguali e cercando di confutare questo fatto.
- La confutazione passa dal ragionamento "i dati che ho a disposizione sono così incompatibili con l'ipotesi nulla che essa non può che essere falsa e, dunque, con un certo livello di confidenza, posso dire che l'ipotesi alternativa (quella che volevo dimostrare) è vera".





Statistica del test

Matematicamente, i test statistici si basano sulla costruzione di una statistica, cioè di una variabile aleatoria "intelligentemente congegnata" che ha le seguenti tre caratteristiche

- 1. Ha senso solo se prendiamo per vera l'ipotesi nulla.
- 2. Si può formulare e se ne possono calcolare i valori a partire da <u>informazioni facilmente reperibili</u> sui campioni statistici a disposizione (es.: media, varianza, numerosità, ...).
- 3. È <u>matematicamente "comoda"</u>, cioè, sappiamo "tutto" della sua distribuzione, che è completamente definita date l'ipotesi nulla presa per vera al punto 1 e le poche informazioni di cui al punto 2.
 - In pratica, conviene pensare alla **statistica del test** come un marchingegno matematico che noi **usiamo "così com'è"**, senza farci troppe domande.
 - Sempre in pratica, la **forma matematica complicata** della statistica del test è principalmente "di comodo": all'aumentare della complessità del test, essa **perde progressivamente riscontro intuitivo** e diventa progressivamente sempre più difficile (fino a praticamente impossibile) interpretarla in termini di fenomeni reali.
 - Più formalmente, vale che, in generale, <u>una statistica di test non è una buona statistica</u> <u>descrittiva</u> della popolazione (e le statistiche descrittive di una popolazione raramente sono buone statistiche di test).





Test statistici in pratica

- 1. Individuo la proprietà che voglio dimostrare e scelgo l'ipotesi nulla che, qualora dovesse essere rifiutata, mi darebbe la risposta che cerco.
- 2. Scelgo una regola di decisione, cioè un livello di significatività (lo posso scegliere "quando voglio", basta che sia **prima** di fare i calcoli).
- 3. Verifico alcune assunzioni di partenza sui dati (v. prossimo laboratorio; oggi ci fidiamo).
- 4. Assumo che l'ipotesi nulla sia vera (importante!).
- 5. Individuo la statistica corrispondente all'ipotesi nulla.
- 6. Calcolo i parametri necessari a partire dai dati a disposizione.
- 7. Calcolo il valore della statistica usando i numeri del punto sopra.
- 8. Confronto i p-value corrispondente con il livello di significatività e:
 - 1. Se il p-value è minore del livello di significatività, rifiuto l'ipotesi nulla.
 - 2. Se il p-value è maggiore del livello di significatività, mi chiudo in rispettoso silenzio, non potendo concludere niente sull'esperimento, (in questo caso, si dice che """"accetto"""" l'ipotesi nulla).





Test statistici in pratica (1 di 8)

• Individuo la proprietà che voglio dimostrare e scelgo l'ipotesi nulla che, qualora dovesse essere rifiutata, mi darebbe la risposta che cerco.

Voglio dimostrare che la media della popolazione da cui ho estratto il mio campione con varianza nota $33 \text{ (mg/dL)}^2 \text{ non } \text{ è } 100 \text{ mg/dL}$.

- Mi chiedo: "Esiste un'ipotesi nulla che mi aiuterebbe, se negata, in tal senso?"
- Risposta: "Sì, posso assumere H_0 : $\mu=100$ mg/dL"





Test statistici in pratica (2 di 8)

• Scelgo una regola di decisione, cioè un livello di significatività (lo posso scegliere "quando voglio", basta che sia **prima** di fare i calcoli).

Scelgo
$$\alpha = 0.05$$

Non c'è molto da dire: lo **scelgo "per protocollo"** sulla base di quanto severo voglio essere con il test:

- Se mi basta un ragionevole dubbio che l'ipotesi nulla sia falsa, sceglierò α più grande.
- Se mi serve essere molto più sicuro, sceglierò lpha più piccolo.

L'importante è sceglierlo prima (!!!) di fare i calcoli.





Test statistici in pratica (3 di 8)

• Verifico alcune assunzioni di partenza sui dati (v. prossimo laboratorio; oggi ci fidiamo).

La statistica del test si può costruire correttamente solo se sono valide contemporaneamente:

- 1. L'ipotesi nulla.
- 2. Alcune ipotesi su tutte le variabili aleatorie coinvolte (es. quelle di cui i dati sono realizzazione).

Sull'ipotesi nulla non si può sindacare.

Le altre ipotesi vanno verificate e bisogna essere consapevoli di quanto si stia deviando dal "caso ideale"; a volte può non essere un problema grave (v. stimatori robusti).

=> Prossimo laboratorio.





Test statistici in pratica (4 di 8)

Assumo che l'ipotesi nulla sia vera (importante!).

Devo interiorizzare e ricordami a ogni passaggio che l'ipotesi nulla è un'ipotesi e, dunque, è data per vera da ora in avanti.

Sembra controintuitivo, perché stiamo cercando di confutarla, ma dobbiamo darla per vera fino al momento in cui la confutiamo (*).

Altrimenti, la matematica perde di significato e, con essa, tutti i numeri che calcoliamo.

(*) Il ragionamento torna perché, tecnicamente, non confutiamo mai H_0 categoricamente; bensì solo con una certa "confidenza", corrispondente al livello di significatività α .





Test statistici in pratica (5 di 8)

• Individuo la statistica corrispondente all'ipotesi nulla.

All'atto pratico, non ci è mai richiesto di inventare la statistica del test.

• Semplicemente, cerchiamo in letteratura la statistica che meglio corrisponde all'ipotesi nulla e alla nostra situazione sperimentale.

Nel nostro caso, l'ipotesi nulla era H_0 : $\mu=100~{\rm mg/dL}$ e sapevamo che varianza del campione era 33 (mg/dL)² (e abbiamo glissato sulle altre ipotesi; v. prossimo lab).

- Dunque, usiamo uno ztest a un campione.
- Cerchiamo la funzione MATLAB corrispondente: ztest chiamata con un vettore di dati e due costanti (la media e la standard deviation nota).

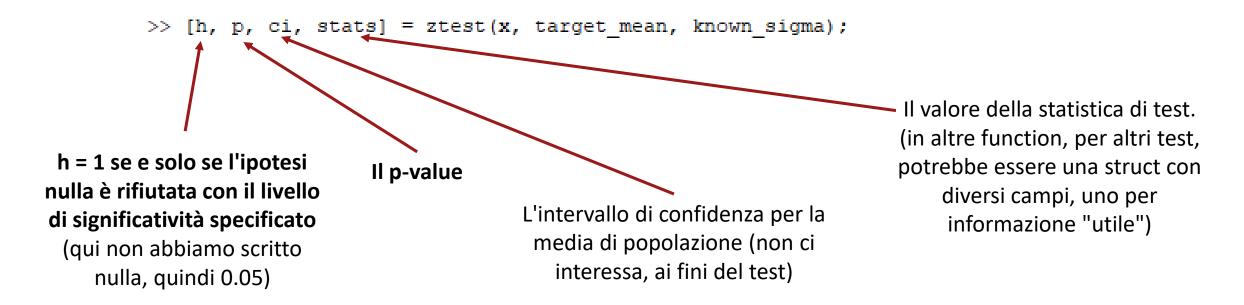




Test statistici in pratica (6 e 7 di 8)

- Calcolo i parametri necessari a partire dai dati a disposizione.
- Calcolo il valore della statistica usando i numeri del punto sopra.

MATLAB fa entrambe le cose per noi. In particolare, ci restituisce quanto segue







Test statistici in pratica (8 di 8)

- Confronto i p-value corrispondente con il livello di significatività e:
 - 1. Se il p-value è minore del livello di significatività, rifiuto l'ipotesi nulla.
 - 2. Se il p-value è maggiore del livello di significatività, mi chiudo in rispettoso silenzio, non potendo concludere niente sull'esperimento, (in questo caso, si dice che """"accetto"""" l'ipotesi nulla).

Si tratta di confrontare il valore di p restituito dalla function con $\alpha=0.05$ che avevamo deciso all'inizio.

• Possiamo addirittura leggere l'esito del confronto nel parametro di uscita h.

Il calcolo non è un problema, l'interpretazione potrebbe esserlo (v. slide successiva)







Test statistici in pratica (per la vita)

L'ipotesi nulla è un'ipotesi

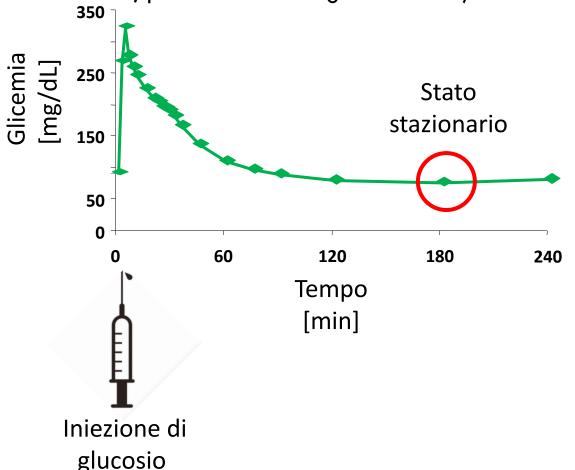
- Se il p-value è alto, ovvero i dati sembrano concordare con l'ipotesi nulla, questo NON MI DICE ASSOLUTAMENTE NULLA rispetto al fatto che l'ipotesi nulla sia vera oppure no.
- Se non <u>assumo (cioè do per assodato, cioè faccio finta!) che l'ipotesi</u> <u>nulla sia vera,</u> non posso neanche costruire la statistica del test e, quindi, non posso calcolare il p-value (che si basa sulla distribuzione della statistica del test).
 - Sarebbe come dire "Se questo (l'ipotesi nulla) è vero, allora è vero con una certa probabilità uguale al p-value".
 - Ma che senso ha? Se è vero, è vero e basta; cosa c'entra la probabilità?
 - Ecco, dunque, il cortocircuito logico che ci indica che l'ipotesi nulla non si può dimostrare, ma va presa per buona e, eventualmente, confutata.





CONTESTO DELL'ESERCITAZIONE

Utilizziamo una variante del dataset Intravenous Glucose Tollerance Test (IVGTT) dello laboratorio 2 (il nome file sarà diverso, perché non sono gli stessi dati).



Abbiamo a disposizione i dati di concentrazione di glucosio (glicemia) allo stato stazionario raccolti in individui appartenenti a due popolazioni:

- 1) Anziani (glucose_ss_elderly)
- 2) Giovani (glucose_ss_young)

I dati sono contenuti nel file IVGTT_SS.mat

- data: matrice double 2D (un gruppo di pazienti per colonna).
- labels: cell array con le etichette delle colonne
 (= nomi delle variabili / gruppi).
- units: cell array con le unità di misura.

NB: avere i dati nella matrice in questo modo è piuttosto inconsueto (perché dati sulla stessa riga non corrispondono allo stesso individuo), ma può capitare. Ci aiuta a fare esercizio di preprocessing.





Prima di svolgere l'esercitazione (oppure al bisogno), utilizzare il comando help di MATLAB seguito dal nome delle seguenti function, utili allo svolgimento degli esercizi.

Parte 1 (svolto)

--

Parte 2 (proposto, per casa, da fare alla fine solo se avanza tempo)

skewness, kurtosis, igr

Parte 3 (svolto)

• [h, p, ci, stats] = ztest(dati, media_target, standard_deviation_nota)

Parte 4 (proposto)

• [h, p, ci, stats] = ztest(dati, media_target, standard_deviation_nota, 'Tail', 'right')
oppure [h, p, ci, stats] = ztest(dati, media_target, standard_deviation_nota, 'Tail', 'left')

Parte 5 (proposto)

• __

Parte 6 (proposto)

• [h, p, ci, stats] = ztest(dati, media_target, standard_deviation_nota, 'Alpha', livello_di_significativita)





PARTE 1: CARICAMENTO E PULIZIA DATI (svolto)

- Caricare il file IVGTT_SS.mat
- Effettuare un'analisi naive sull'intera tabella, individuando e contando
 - I valori mancanti (NaN), con la funzione isnan
 - I valori infiniti (Inf), con la funzione isinf
 - I valori negativi, con la consueta indicizzazione logica
- Sostituire eventuali valori infiniti e valori negativi con NaN (v. laboratorio 3 per il razionale)





PARTE 2: CARICAMENTO E PULIZIA DATI (proposto, per casa, da fare alla fine solo se avanza tempo)

- Per ciascuna colonna della matrice data
 - Individuare e rimuovere i dati mancanti
 - Calcolare media e standard deviation
 - Calcolare mediana e IQR (funzione iqr)
 - Calcolare skewness e curtosi (skewness, kurtosis)
 - Mostrare a schermo tutte queste informazioni
- Rappresentare in una figura con 4 pannelli:
 - In alto a sinistra, il boxplot della variabile glucose_ss_elderly
 - In alto a destra, il boxplot della variabile glucose_ss_young
 - In basso a sinistra, l'istogramma delle frequenze assolute della variabile glucose_ss_elderly
 - In alto a destra, l'istogramma delle frequenze assolute della variabile glucose_ss_young
- Che idea ci possiamo fare sulla gaussianità delle due variabili?





PARTE 3: TEST STATISTICI (svolto)

Considerando solamente la variabile **glucose_ss_elderly**, privata dei suoi eventuali dati mancanti:

- Impostare un test statistico a un campione a due code per confrontare la media con il valore 100 mg/dL (livello di significatività $\alpha=0.05$). Assumere varianza nota = 33 (mg/dL)².
- Impostare un test statistico a un campione a due code per confrontare la media con il valore 82 mg/dL (livello di significatività $\alpha=0.05$). Assumere varianza nota = 33 (mg/dL)².
- Commentare entrambi i risultati ottenuti (con moltissima attenzione).

PARTE 4: TEST STATISTICI (proposto)

Considerando solamente la variabile **glucose_ss_elderly**, privata dei suoi eventuali dati mancanti:

- Impostare un test statistico a un campione e **a una coda** per confrontare la media con il valore 100 mg/dL (livello di significatività $\alpha=0.05$; ipotesi alternativa di interesse: $\mu>100$ mg/dL). Assumere varianza nota = 33 (mg/dL)².
 - Commentare il risultato (facile).
 - Darsene una spiegazione (difficile, lo rivedremo al prossimo laboratorio).





PARTE 5: L'IPOTESI NULLA NON SI PUO' DIMOSTRARE (proposto)

Considerando solamente la variabile **glucose_ss_elderly**, privata dei suoi eventuali dati mancanti:

- Impostare due test statistici a un campione e a due code per confrontare la media con i valori 81.5 e 82.5 mg/dL (livello di significatività $\alpha=0.05$). Assumere varianza nota = 33 (mg/dL)².
 - Commentare il risultato di ciascun test individualmente (facile).
 - Considerare cosa si può dire dopo averli fatti entrambi ("difficile", ma istruttivo).

PARTE 6: ESERCIZIO DI SINTASSI (proposto)

Considerando solamente la variabile **glucose_ss_elderly**, privata dei suoi eventuali dati mancanti:

- Ripetere l'esercizio PARTE 3 punto 1 (confronto con 100 mg/dL) usando un livello di significatività $\alpha = 0.01$ (consultare <u>autonomamente</u> l'help di MATLAB). Assumere varianza nota = 33 (mg/dL)².
- Commentare il risultato ottenuto

PARTE 7: <u>DOPO</u> IL PROSSIMO LABORATORIO (proposto per prepararsi all'esame)

• Svolgere le parti 3, 4, 5 e 6 utilizzando il test statistico appropriato in caso di varianza incognita (funzione **ttest**). Quindi <u>non</u> assumere più la varianza nota di prima!