

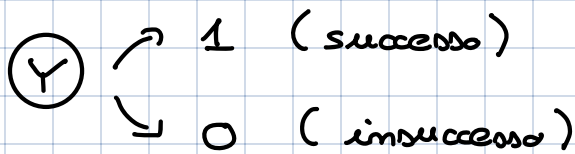
## DATI CATEGORICI

- ex) • TIPO DI SANGUE  $\rightarrow \{O^+, O^-, A^+, A^-, B^+, B^-, AB^+, AB^-\}$  (8 LIVELLI) NO ORDINE
- STATO DI APPARTENENZA  $\rightarrow \{ITALIA, GIAPPONE \dots\}$  ( $\sim 200$  LIVELLI) NO ORDINE
- RISPOSTA A UN FARMACO  $\rightarrow \{+, 0, -\}$  (3 LIVELLI) SÌ ORDINE  
(ha senso dire che '+' > '-')

MA

LA DISTANZA FRA I LIVELLI NON È BEN DEFINITA

## VARIABILE BINARIA



$$P_Y(y) = \begin{cases} p & \text{se } y = 1 \\ 1-p & \text{se } y = 0 \end{cases}$$

$Y_i \stackrel{iid}{\sim} \text{Bernoulli}(\pi)$

$$N = \sum_{i=1}^n Y_i \sim \text{Binomiale}(n, \pi)$$

$$P_N(n) = \binom{n}{N} p^N (1-p)^{n-N}$$

$n = n^*$  esperimenti  
 $N = n^*$  successi

## TEST D'IPOTESI

$$\begin{aligned} H_0 &: \pi = \pi^0 \\ H_A &: \pi > \pi^0 \end{aligned}$$

al livello  $\alpha$ ,  
unilaterale

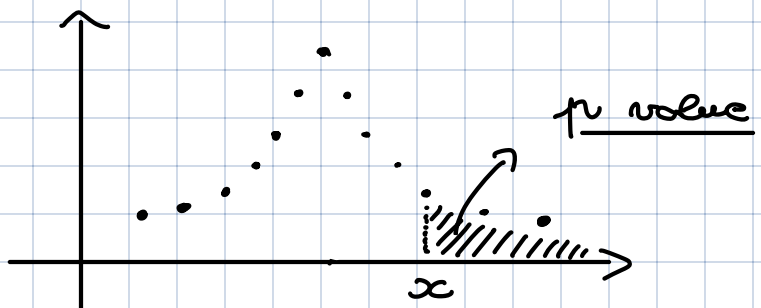
- esperimento  $\rightarrow x$  successi

$\rightarrow$  RIFIUTA  $H_0$  SE

$\searrow$  ACCETTO SE

$$\begin{aligned} \sum_{x=m+1}^n P_N(x) &< \alpha \\ \sum_{x=m+1}^n P_N(x) &\geq \alpha \end{aligned}$$

**OSS:** Accettare significa non avere abbastanza evidenza per rifiutare



# INTERVALLI DI CONFIDENZA (di FIDUCIA)

- CLIPPER PEARSON ①
- WALD ②
- WILSON ③

$\alpha$ : livello di fiducia

$$\textcircled{1} \quad L^{CP} : \min \{ \pi \mid IP(\text{Bim}(n, \pi) \leq x) > \frac{\alpha}{2} \}$$

$$U^{CP} : \max \{ \pi \mid IP(\text{Bim}(n, \pi) > x) > \frac{\alpha}{2} \}$$

OSS: È un intervallo "esatto" perché utilizza la densità binomiale ma ha copertura  $> 1-\alpha$  (conservativo)

$$\textcircled{2} \quad \sqrt{n} \left( \frac{N}{n} - \pi \right) \xrightarrow[\text{TLC}]{n \rightarrow \infty} N(0, \pi(1-\pi))$$

$$\frac{N}{n} = \hat{\pi} \quad IP \left( \left| \frac{\sqrt{n}(\hat{\pi} - \pi)}{\sqrt{\pi(1-\pi)}} \right| < z_{\frac{1-\alpha}{2}} \right) = 1-\alpha \quad (*)$$

$$IP \left( -z_{\frac{1-\alpha}{2}} \leq \frac{\sqrt{n}(\hat{\pi} - \pi)}{\sqrt{\pi(1-\pi)}} \leq z_{\frac{1-\alpha}{2}} \right) = 1-\alpha$$

SOSTITUIRE  $\pi$  CON  $\hat{\pi}$  !!

$$IP \left( \underbrace{\hat{\pi} - z_{\frac{1-\alpha}{2}} \sqrt{\frac{\hat{\pi}(1-\hat{\pi})}{n}}}_{L^{WALD}} \leq \pi \leq \underbrace{\hat{\pi} + z_{\frac{1-\alpha}{2}} \sqrt{\frac{\hat{\pi}(1-\hat{\pi})}{n}}}_{U^{WALD}} \right) = 1-\alpha$$

$$I.C.^{WALD} = [L^{WALD}, U^{WALD}]$$

$$\textcircled{3} \quad (*)^2 = \left( \frac{\sqrt{n}(\hat{\pi} - \pi)}{\sqrt{\pi(1-\pi)}} \right)^2 = z_{\frac{1-\alpha}{2}}^2$$

$$(\hat{\pi} - \pi)^2 = \frac{\pi(1-\pi) z_{\frac{1-\alpha}{2}}^2}{n}$$

$$\hat{\pi}^2 - 2\pi\hat{\pi} + \pi^2 = \frac{\pi}{n} z_{\frac{1-\alpha}{2}}^2 - \pi^2 \frac{z_{\frac{1-\alpha}{2}}^2}{n}$$

$$\pi^2 \left( 1 + \frac{z_{\frac{1-\alpha}{2}}^2}{n} \right) - \pi \left( 2\hat{\pi} + \frac{z_{\frac{1-\alpha}{2}}^2}{n} \right) \pi + \hat{\pi}^2 = 0$$

EQ II GRADO IN  $\pi$

Wilson  
I.C. :  $\frac{1}{1 + \frac{z_{\alpha/2}^2}{n}} \left( \hat{\pi} + \frac{z_{\alpha/2}^2}{2n} \pm \frac{z_{\alpha/2}}{2n} \sqrt{4n(1-\hat{\pi})\hat{\pi} + z_{\alpha/2}^2} \right)$

## VARIABILE CATEGORICA MULTILIVELLO

Variabile categorica  $C$  che può assumere  $D$  livelli.

↓  
CODIFICA  
STANDARD

$$C = i \quad i = 1 \dots D$$

↓  
CODIFICA  
ONE-HOT

$$\vec{B} = (\overbrace{0, 0}^{x_1, x_2}, \dots, \overbrace{0, 1, 0 \dots 0}^{x_i})$$

↓  
 $i$ -esima  
posizione

OSS: Sono sufficienti  $d = D - 1$   
bit per codificare  $\vec{B}$

I parametri che governano la distribuzione  
di  $\vec{B}$  sono

$$\vec{\pi} = (\pi_1, \dots, \pi_d)$$

$$\pi_i = P(\vec{B} = (0, \dots, \underbrace{0, 1}_i, \dots, 0))$$

$i$ -esima  
posizione  $\forall i = 1 \dots d$

Risultati che ci serviranno dopo

- $E[X_k] = \pi_k$
- $Cov(X_k, X_l) = E[X_k X_l] - E[X_k] E[X_l]$ 

$$= \begin{cases} -\pi_k \pi_l & k \neq l \\ \pi_k - \pi_k^2 & k = l \end{cases}$$

↓  
 $\pi_k(1 - \pi_k)$

2 LIVELLI

$$Y \sim \text{Bernoulli}(\pi)$$

D LIVELLI

$$\vec{B} \sim \text{Multinomial}(\vec{\pi})$$

$$\text{can}_d \vec{\pi} = (\pi_1 \dots \pi_d)$$

OSS:  $\sum_{i=1}^D \pi_i \neq 1$   
 $\pi_D = 1 - \sum_{i=1}^{D-1} \pi_i$

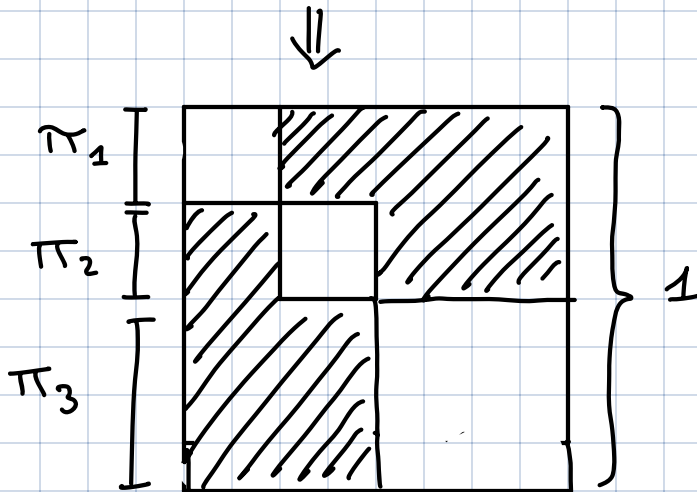
## INDICI DI VARIABILITA' PER MULTINOMI

- $E = \sum_{i=1}^D \pi_i \log(\pi_i)$

ENTROPIA DI SHANNON

- $G = 1 - \sum_{i=1}^D \pi_i^2$

DIVERSITA' DI GINI



— • — • — • — •

$$Y_i \sim \text{Bernoulli}(\pi)$$

$$N = \sum_{i=1}^n Y_i \sim \text{Bin}(n, \pi)$$

$$\vec{B}_i \sim \text{Multinomial}(\vec{\pi})$$

$$\vec{N} = \sum_{i=1}^n \vec{B}_i \sim \text{Multinomial}(n, \vec{\pi})$$

↙

$$f_N(m_1, \dots, m_D) = \frac{n!}{m_1! m_2! \dots m_D!} \prod_{i=1}^D \pi_i^{m_i}$$

## RISULTATO ASIMPTOTICO (TCL MULTIVARIATO)

$$\sqrt{n} \left( \frac{N_1}{n} - \pi_1, \dots, \frac{N_D}{n} - \pi_D \right) \xrightarrow{n \rightarrow \infty} \mathcal{N}_D(\vec{0}, \Sigma)$$

$$\text{con } \Sigma = [\sigma_{ij}] = \begin{cases} \pi_i(1-\pi_i) & i=j \\ -\pi_i\pi_j & i \neq j \end{cases}$$

TEST PER  $\vec{\pi}$

$$\begin{aligned} H_0: & \vec{\pi} = \vec{\pi}_0 \\ H_A: & \vec{\pi} \neq \vec{\pi}_0 \end{aligned}$$

### ① Caso Binomiale

$$\text{RIFIUTO } H_0 \text{ SE } \left| \frac{N - n\pi^0}{\sqrt{n\pi^0(1-\pi^0)}} \right| > Z_{\frac{\alpha}{2}}$$

$$\text{O EQUIVALENTEMENTE } \left( \frac{N - n\pi^0}{\sqrt{n\pi^0(1-\pi^0)}} \right)^2 > \underbrace{Z_{\frac{\alpha}{2}}^2}_{= \chi^2_{\alpha}(1)}$$

$$\begin{aligned} \frac{N^2 - 2Nn\pi^0 + n^2\pi^{0^2}}{n\pi^0(1-\pi^0)} &= \dots = \frac{(N - \pi^0 n)^2}{\pi^0 n} + \frac{(n - N - n(1-\pi^0))^2}{n(1-\pi^0)} \\ &= \sum_{i=1}^2 \frac{(N_i - n\pi_i^0)^2}{n\pi_i^0} = \chi^2 \end{aligned}$$

CHI-QUADRO DI PEARSON

dove  $\pi_1^0$  = probabilità successo sotto  $H^0$   
 $\pi_2^0$  = probabilità insuccesso sotto  $H^0$

$$\chi^2 > \chi^2_{\alpha}$$

RIFIUTO  $H_0$

$$\chi^2 \leq \chi^2_{\alpha}$$

ACCETTO  $H_0$

TEST  
CHI-QUADRO

## ② Caso multinomiale

**TEO** 
$$\sum_{i=1}^D \frac{(N_i - n\pi_i^0)^2}{n\pi_i^0} \longrightarrow \chi^2(d)$$

**Dim)** ① 
$$\sum_{i=1}^D \frac{(N_i - n\pi_i^0)^2}{n\pi_i^0} = \left( \frac{N_1}{n} - \pi_1, \dots \right) \Sigma^{-1} \left( \frac{N_1}{n}, \dots \right)^T$$

↙  
x CASA

HINT:  $\Sigma^{-1} = [\chi_{ij}] = \begin{cases} \frac{1}{\pi_i} + \frac{1}{\pi_D} & \text{se } i=j \\ \frac{1}{\pi_D} & \text{se } i \neq j \end{cases}$

② 
$$\left( \frac{N_1}{n} - \pi_1, \dots \right) \Sigma^{-1} \left( \frac{N_1}{n} - \pi_1, \dots \right)^T \longrightarrow \chi^2(d)$$

$$X^2 = \sum_{i=1}^D \frac{(N_i - n\pi_i^0)^2}{n\pi_i^0}$$

**TEST**  
**CHI-QUADRO**

↗  $X^2 > \chi^2(d)$  RIFUTATO  $H_0$

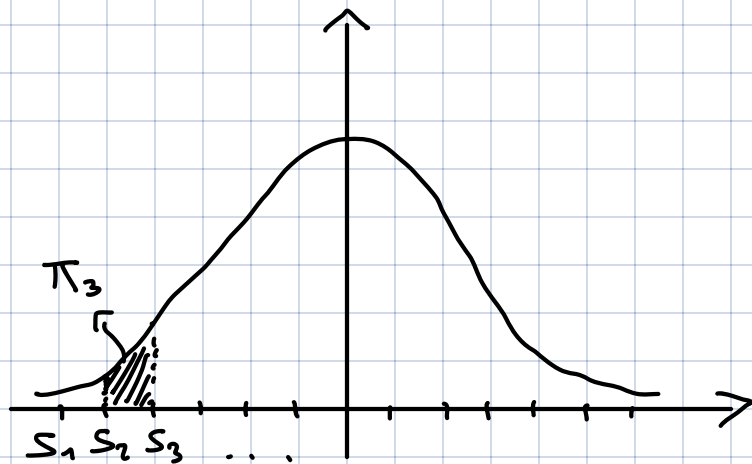
↘  $X^2 \leq \chi^2(d)$  ACCETTO  $H_0$

## TEST DI NORMALITA' ( $\mu, \sigma^2$ NOTI)

$y_i \stackrel{iid}{\sim} N(0, 1)$

$$\pi_i = \frac{1}{\sqrt{2\pi}} \int_{S_{i-1}}^{S_i} e^{-y^2/2} dy$$

$i = 1 \dots D$



**TEST:** 
$$X^2 = \sum_{i=1}^D \frac{(N_i - n\pi_i)^2}{n\pi_i} \sim \chi^2(d)$$

TEO  
BIS

Se le probabilità  $\vec{\pi}$  dipendono da un  
vettore di  $r$  parametri  $\hat{\theta}$

$$\chi^2(\hat{\theta}) = \sum_{i=1}^D \frac{(N_i - n \pi_i(\hat{\theta}))^2}{n \pi_i(\hat{\theta})} \rightarrow \chi^2_{(d-r)}$$

### TEST DI NORMALITA'

$$y_i \stackrel{iid}{\sim} N(\mu, \sigma^2)$$

$$\hat{\theta} = (\bar{y}, s^2)$$

media  
campion.

varianza  
campionaria

←  
come nel caso  
del TEO BIS

$$\pi_i(\hat{\theta}) = \frac{1}{\sqrt{2\pi} s} \int_{S_{i-1}}^{S_i} \exp\left(-\frac{1}{2s^2} (y - \bar{y})^2\right) dy$$

TEST:

$$\chi^2(\hat{\theta}) = \sum_{i=1}^D \frac{(N_i - n \pi_i(\hat{\theta}))^2}{n \pi_i(\hat{\theta})}$$

REFIUTO SE  $\chi^2(\hat{\theta}) > \chi^2_{\alpha}(d-2)$

ACCETTO  $H_0$  SE  $\chi^2(\hat{\theta}) \leq \chi^2_{\alpha}(d-2)$

## DUE VARIABILI CATEGORICHE

A \ C	1	2	...	J	
1	$N_{11}$	$N_{12}$	...	$N_{1J}$	$M_{1+} = \sum_{j=1}^J M_{1j}$
2	$N_{21}$				
$\vdots$	$\vdots$				
I	$N_{I1}$				
	$N_{+1}$	$N_{+2}$			$n$

TEST DI OMogeneITA' (mi chiedo se le probabilità del fattore J sono omogenee lungo fattore i)

$$H_0: \pi_{j|1} = \pi_{j|2} = \dots = \hat{\pi}_{+j} \quad \forall j=1 \dots J$$

Sono parametri stimati a partire dai dati !!!

(Sono nel caso del TEC BIS)

$$\hat{\pi}_{+j} = \frac{N_{+j}}{n} \quad \forall j=1 \dots J$$

$$\text{n° param.} = I(J-1)$$

$$\text{n° vincoli} = J-1$$

$$\chi^2 = \sum_{i=1}^I \sum_{j=1}^J \frac{(N_{ij} - m_{i+} \hat{\pi}_{+j})^2}{m_{i+} \hat{\pi}_{+j}}$$

→ RIFIUTATO SE  $\chi^2 > \chi^2_{\alpha}((I-1)(J-1))$



↙ ACCEPTANCE SE  $\chi^2 \leq \chi^2_{\alpha} ((I-1)(J-1))$