

Non - PARAMETRIC TESTS

DISTRIBUTION

KOLMOGOROV - SMIRNOV

EQUALITY OF CONTINUOUS, 1-D DISTRIBUTIONS, COMPARE SAMPLE W/REFERENCE. H_0 : SAME DISTRIBUTION

EMPIRICAL DISTRIBUTION

$F_n(x) = \frac{1}{N} \sum_{i=1}^n N_i \cdot U(x - a_i)$, N FREQUENZA CLASSE a_i , STATISTIC $D_n = \sup_x |F_n(x) - F(x)|$

COMPARO A DISTRIBUZIONE DA TESTARE
D = DEVIAZIONE MASSIMA, USO QUANTILI
↓
DISTANZA VERTICALE

GOODNESS - OF - FIT TEST

REJECT H_0 IF $\sqrt{N} \cdot D_n > \lambda_\alpha$, $K(\lambda_\alpha) = \sum_{j=1}^{\infty} (j-1)^{-2} \lambda_\alpha^2 = 1 - \alpha$ WHERE λ_α SATISFIES
PER $N \gg 20$

CHI-SQUARED GOF TEST

ALSO NON-CONTINUOUS DISTRIBUTIONS
N SAMPLES, M CLASSES

$T_N = \sum_{i=1}^M \frac{N_i^2}{N \cdot p_{i0}} - M$

$n \gg 50, m \gg 5, H_0: p_1 = p_{10}, p_2 = p_{20}, \dots, M-1 \text{ DOF}$
 $T_N \gg \chi^2_{1-\alpha, M-1 \text{ DOF}}$

CATEGORICAL, BINNED DATA

SCOSTAMENTO FRA FREQ. EMPIRICHE E DENSITA' TEORICHE

$\sum_{i=1}^M \frac{(N_i - N \cdot p_i)^2}{N \cdot p_i} = \sum \frac{(\text{OBSERVED} - \text{EXPECTED})^2}{\text{EXPECTED}}$

CAN USE FOR CONTINUOUS DISTRIBUTIONS
 $p_{10} = P(a_{i-1} \leq X \leq a_i) = \int_{a_{i-1}}^{a_i} f_0(x) dx$

WHEN P DEPEND ON θ

MLE ESTIMATE $\hat{\theta}_1$

$\chi^2 = \sum_{i=1}^M \frac{[N - n \pi_i(\hat{\theta})]^2}{n \pi_i(\hat{\theta})}$

HOMOGENEITY

VERIFICARE SE X, Y ESTRATTI DA STESSA POPOLAZIONE / POP. IDENTICHE

$H_0: F_X(t) = F_Y(t)$

~~SIGN TEST~~ ~~WILCOXON SIGNED-RANK TEST~~

$X_i \neq$ SAME NUMEROSITY, m^\pm COUNT WHERE PAIR DIFF IS $\pm \rightarrow$ SAME DISTRIBUTION
 $N_0 \neq$ COUPLE UGUALI
 $S_N = m^+ - m^-$ ON PAIR

$|S_N| < \sqrt{\frac{n - m_0}{2}} \cdot z_{1-\alpha}$

SE UGUALE: DIFFERENZA PICCOLA

$m_0 = |X_i = y_i|$ FREQUENCY; N° OF PAIRS

~~WILCOXON SIGNED-RANK TEST~~ ^{RANK-SUM} ~~WILCOXON MANN-WHITNEY~~

X, Y DIFFERENT NUMEROSITY. SORT INCREASINGLY. RANK = ORDER NUMBER (AUG IF TIE THEN 1). S_X, S_Y = SUMS OF RANKS | SYMMETRIC DISTRIBUTION

$U_X = M \cdot m + \frac{m(m+1)}{2} - S_X, U_Y = \frac{m \cdot m(m+1)}{2} - S_Y, U = \min(U_X, U_Y)$. IF H TRUE $\rightarrow U$ IS NORMAL WITH $\mu_U = \frac{m \cdot m}{2}, \sigma_U^2 = \frac{m \cdot m(m+1)}{12}$

OK WHEN $m, m \gg 4, m+m \gg 20$

$Z_{m,m} = \frac{U - m \cdot m / 2}{\sqrt{m \cdot m(m+m+1)/12}}$; ACCEPT FOR $|z_{m,m}| < z_{1-\alpha}$

SMIRNOV TEST

COUPLE OF DISTRIBUTIONS $F_X(t), F_Y(t)$, EMPIRICAL DISTRIBUTIONS. IF H_0 TRUE $D \rightarrow 0$ FOR INCREASING SAMPLE

INDEPENDENTLY OF DISTRIBUTIONS

$d_{n,m} < \sqrt{\frac{m+m}{m \cdot m}}, \lambda_\alpha$, SAME OF KS TEST

TWO POPULATIONS

SAME PRINCIPLE AS

KOLMOGOROV-SMIRNOV TEST

CHI-SQUARED HOMOGENEITY TEST

SUPER-GENERAL
MULTIPLE REALIZATIONS, CLASS

$$j = 1 \dots N \quad \sum_{i=1}^M M_{ij} = M_j = N \text{ OF } X_j \quad \sum_{j=1}^M \sum_{i=1}^M M_{ij} = \sum_{j=1}^M M_j = N, \text{ TOTAL OBSERVATIONS}$$

STILL BANNED
DATA NEEDED

RELATIVE FREQUENCIES OF CLASSES ARE SAME FOR ALL SAMPLES, $p_k^0, p_1^0, \dots, p_m^0 \approx 1$

$$T_N(\hat{p}^0) = N \left(\sum_{j=1}^M \sum_{i=1}^m \frac{n_{ij}^2}{n_{j \cdot} \cdot n_{\cdot i}} - 1 \right)$$

$$\hat{p}_N \in \chi^2_{1-\alpha, (m-1)(M-1) \text{ DOF}}$$

ALT FORMULATION

IF $m = 2$, IN/OUT, TRUE/FALSE

$$T_N(M_A) = \frac{N}{N - M_A} \left(\frac{N}{M_A} \sum_{j=1}^M \frac{n_{Aj}^2}{n_{j \cdot}} - M_A \right) \cdot \hat{e}_{ij} = n_{i \cdot} \cdot \frac{n_{\cdot j}}{n}, \quad \chi^2 = \sum_{j=1}^M \sum_{i=1}^2 \frac{(n_{ij} - \hat{e}_{ij})^2}{\hat{e}_{ij}}$$

$$\chi^2 = \sum \frac{(\text{OBS} - \text{EXP})^2}{\text{EXP}}$$

1 POPULATIONS
2 TREATMENTS/BINS/CATEGORIES

INDEPENDENCE TESTS

CHI-SQUARED INDEPENDENCE TESTS

DISCRETE RV, OR DIVIDED IN CLASSES

$$(X_1, \dots, X_{M \times X}) (Y_1, \dots, Y_{M \times Y})$$

$n_{hk} = \#$ COPIES

$$\sum_{h=1}^{M \times X} n_{hk} = n_{\cdot k}, \quad \sum_{k=1}^{M \times Y} n_{hk} = n_{h \cdot}, \quad \sum_{h=1}^{M \times X} \sum_{k=1}^{M \times Y} n_{hk} = n$$

$n_{h \cdot}, n_{\cdot k}$ MARGINAL
ABSOLUTE FREQUENCIES
 $(M-1)(N-1)$ DOF

MLE $\forall h, k \quad \hat{p}_h^0 = \frac{n_{h \cdot}}{n}, \quad \hat{p}_k^0 = \frac{n_{\cdot k}}{n} \rightarrow \text{MINIMIZING}$

$$\frac{n_{hk}}{n} - \hat{p}_h^0 \cdot \hat{p}_k^0 = \left(-\frac{n_{hk}}{n} + n_{hk} \right) \frac{1}{n}$$

OK IF $\forall n_{hk} \geq 5$

ALT FORMULATION

$$T_m(\hat{p}^0) = m \left(\sum_{h=1}^{M \times X} \sum_{k=1}^{M \times Y} \frac{n_{hk}^2}{n_{h \cdot} \cdot n_{\cdot k}} - 1 \right)$$

REJECT H_0 IF $t_m \geq \chi^2_{1-\alpha, (M \times X - 1)(M \times Y - 1)}$

$$\hat{e}_{ij} = n_{i \cdot} \cdot p_{\cdot j} = \frac{n_{i \cdot} \cdot n_{\cdot j}}{n}, \quad \chi^2 = \sum_{j=1}^M \sum_{i=1}^m \frac{(n_{ij} - \hat{e}_{ij})^2}{\hat{e}_{ij}}, \quad \chi^2 = \frac{(O-E)^2}{E}$$

SPEARMAN TEST (RANK CORR COEFF)

X_k, Y_k SORT, COMPUTE RANKS SEPARATELY

$$d_k = R_n(X_k) - R_n(Y_k)$$

$$R_s = 1 - \frac{6}{n(n^2 - 1)} \cdot \sum_{k=1}^n d_k$$

SPEARMAN
CORR
COEFF

NECESSARY, NOT SUFFICIENT CONDITION FOR INDEPENDENCE

$$-1 \leq R_s \leq 1$$

$H_0: X_k \overset{\text{NONCORRELATION}}{\text{INDEPENDENCE}} Y_k \rightarrow \text{AVG}(R_s) = 0$

$$n \geq 10 \quad T_s = R_s \sqrt{\frac{n-2}{1-R_s^2}} \quad \text{IS T-STUDENT, } n-2 \text{ DOF}$$

REJECT H_0 IF $|t_s| > t_{1-\alpha/2}$

RANDOMNESS TESTS

RANDOMNESS \rightarrow INDEPENDENCE

$$F_X(x_1, \dots, x_m) = F_{X_1}(x_1) \cdot F_{X_2}(x_2) \cdot \dots \cdot F_{X_m}(x_m)$$

H_0 : SAMPLES ARE RANDOM

SERIAL CORRELATION TEST

IF RANDOM \rightarrow ANY SUBSET (X_n, X_{n+1}) IS INDEPENDENT

$$R_S = \sum_{n=1}^N \frac{(X_n - \bar{X})(X_{n+1} - \bar{X})}{N S_X^2}$$

CIRCULAR CORRELATION COEFFICIENT

PERMUTATIONS ARE EQUIPROBABLE \rightarrow

$$R_X = \sum_{n=1}^N (X_n - \bar{X})(X_{n+1} - \bar{X})$$

IF X RANDOM, $R_X \sim N$ $n \rightarrow \infty$

IS NORMALE

$$E[R_X] = -\frac{S_2}{n-1}$$

$$\sigma^2[R_X] = \frac{S_2^2 - S_4}{N-1}$$

$$S_q = (X_1 - \bar{X})^q + (X_2 - \bar{X})^q + \dots + (X_n - \bar{X})^q \quad q=2,4,$$

STANDARDIZE \rightarrow

$$Z_X = \frac{R_X - E[R_X]}{\sigma[R_X]}$$

REJECT H_0 $|Z_X| > Z_{1-\alpha/2}$

RUN TEST

RUN: N° OF IDENTICAL SUBSTRINGS

'AABBAAA BA ...'

2 CLASSES ONLY

U = TOTAL N° OF RUNS IN SEQUENCES

IF RANDOM U WILL BE FAR FROM EXTREME VALUES DUE TO RANDOM = MAX INFORMATION

GOOD FOR $N_A, N_B > 10$. $U \rightarrow$ NORMAL

$$E[U] = 1 + \frac{2N_A N_B}{N_A + N_B}$$

$$\sigma_U^2 = \frac{2N_A N_B (2N_A N_B - N_A - N_B)}{(N_A + N_B)^2 (N_A + N_B - 1)}$$

STANDARDIZE \rightarrow

$$Z_U = \frac{U - E[U]}{\sigma_U}$$

REJECT H_0 $|Z_U| > Z_{1-\alpha/2}$

CAN USE ON REALS

IF CONSIDER $\begin{cases} A: x < \hat{x} \text{ (MEDIAN)} \\ B: x > \hat{x} \end{cases}$