

## 2. NONPARAMETRIC INFERENCE

The nonparametric approach is based on the fact that we don't want to make "strong" assumptions on the distribution of the data. However, this may imply heavy computational methods.

"strong" assumptions = (e.g.) gaussian data  
(which are not so frequent in reality)

Let's see for instance what happens if we try to use the t-test on non gaussian data (Example 1.). The dataset is small and nongaussian. ( $\Rightarrow$  we cannot use t-test and we either can rely on central limit theorem / asymptotic results).

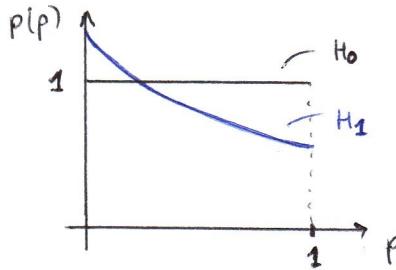
We sample, with sample size  $n=5$ , from different distributions (normal, uniform and t-student). We generate 100.000 datasets and we apply 100.000 t-tests (and we collect the p-values).

We consider the t-test for one population; assumptions: iid normal data

t-test: exact test, i.e., under  $H_0$  the distribution of the p-value (which is a random variable depending on the random sample) is the uniform distr.

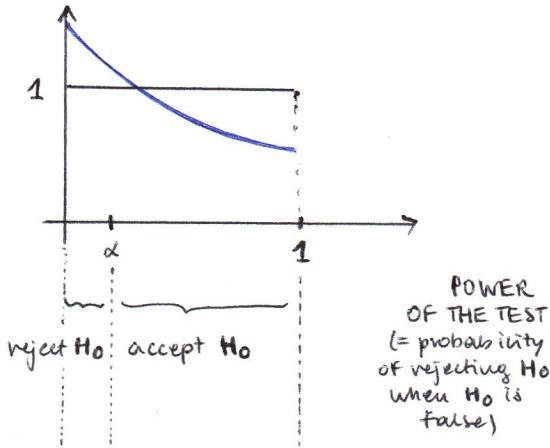
$U([0,1]) \Rightarrow$  under  $H_0$  all the values of the p-values have the same probability to happen. Under  $H_1$  the distribution of the p-value change:

(if we're not under  $H_0$  we're more likely to see smaller values of p-values)

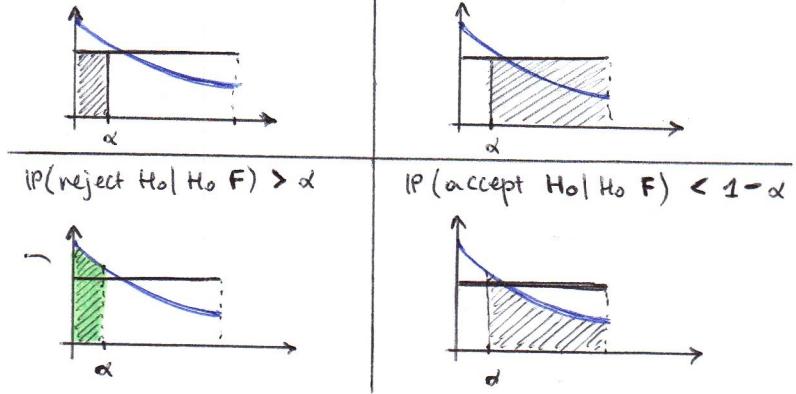


- depends on the violation of  $H_0$  (in this case it depends on the mean): the larger the difference between the true mean and the tested value of the mean and the more concentrated this distribution will be around 0.

so if we fix  $\alpha = \text{IP}(\text{type I error}) = \text{IP}(\text{rejecting } H_0 | H_0 \text{ is true})$



$$\text{IP}(\text{reject } H_0 | H_0 \text{ true}) = \alpha \quad \text{IP}(\text{accepting } H_0 | H_0 \text{ T}) = 1 - \alpha$$



idea of the t-test:

we compare the test statistic, we measure the statistical distance between the hypothesized mean and our estimate of the mean:

$$= \frac{\text{sample mean} - \text{mean under } H_0}{\text{standard deviation of the estimator}} \sim t(n-1)$$

The p-value is just a measure of the extremity of our sample w.r.t. the normal distribution of the sample under the null hypothesis

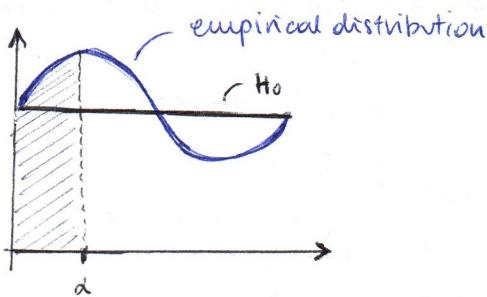
small p-value  $\rightarrow$  reject  $H_0$   
 $\rightarrow$  the data are kind of "extreme" w.r.t. the distribution under  $H_0$

The p-value can be used to run the test: we select a level of significance  $\alpha$  and we reject if the p-value is smaller than  $\alpha$ .

So what are we gonna do with the experiment (Example 1) is:  
 we generate some data under  $H_0$  (= case in which the true mean is identical to the hypothesized value) and we estimate the actual probability of rejecting  $H_0$ .  
 So, for each value of  $\alpha$  ( $:=$  "nominal level") we estimate the empirical level of the test. We'll see if the test is actually rejecting  $H_0$  with probability  $\alpha$  or not.

To make this experiment we have to estimate the distribution of the p-value. (We'll do that in Monte Carlo setting). We'll have an estimate of this true distribution of the p-values under  $H_0$  and to compute the empirical level of the test, for a given value  $\alpha$ , we'll just have to compute how many p-values are smaller than  $\alpha$ .

for instance:



we know that under  $H_0$   $\boxed{\alpha}$  should be  $\alpha$ , if it's much higher/smaller then maybe  $H_0$  is not true!

We can check the empirical  $\alpha$  and the  $H_0 \alpha$  for all the values of  $\alpha \in [0, 1]$ .

( $H_0$  holds  $\Rightarrow \boxed{\alpha} \approx \alpha$ )

In Example 1 we fix  $\alpha = 0.05$ . Since p-values are uniformly distributed under  $H_0$ , having  $\alpha = 0.05$  means that, on average 1 time over 20 that we do the test under  $H_0$  true we wrongly reject  $H_0$ .

(for instance  $\alpha = 0.01$  means 1 time over 100 times that we do the test we wrongly reject  $H_0$ ).

What can we do in practice if we don't have data that satisfies the assumptions of parametric tests?

- we can try to remove outliers
- we can try to transform data (box-cox transformations)
- non parametric way ( $\downarrow$ )

## NON-PARAMETRIC STATISTICS

### <sup>19</sup> SIMULATED EXAMPLE ABOUT THE POSSIBLE FAILURE OF PARAM TESTS. (T-TEST FAILURE.R)

It simply counts the number of deviations

#### SIGN TESTS

(not very used but important since all the other tests are variations of this idea)  
(ARBUTHNOT 1710)

we can do this test with an explicit formula for the calculation of the p-value

#### ONE-SAMPLE (TWO-SIDED) SIGN TEST FOR THE MEDIAN

ASS.  $X_1, X_2, \dots, X_m \sim \text{iid } X$  generic common distribution  
(FOR SIMPLICITY CONTINUOUS)

$$H_0: \text{MED}(X) = c_0 \quad H_1: \text{MED}(X) \neq c_0$$

$$H_0: P(X > c_0) = 0.5 \quad H_1: P(X > c_0) \neq 0.5$$

Test statistic:

$$W = \sum_{i=1}^m I\{X_i > c_0\} \quad \text{NO NUMBER OF DATA GREATER THAN } c_0$$

$$W \in \{0, 1, \dots, m\}$$

$$W \sim B(m, P(X > c_0))$$

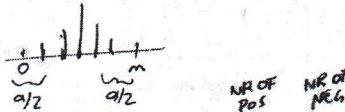
$$\stackrel{H_0}{W \sim B(m, 0.5)} \quad \stackrel{H_1}{W \sim B(m, p)} \quad \text{WITH } p \neq 0.5$$

the test is not based on the mean, it's based on the MEDIAN value of  $X$   
(= value which divides the real axes in two sections of equal probability  
(50%, 50%))

If we want to run:  
 $H_0: P(X > c_0) > 1/2$   
then we rely on the right tail, if the test is left sided we use only the left tail of the distribution

THE CRITICAL REGION IS DEFINED BY THE TWO  $\alpha/2$ -TAIL OF THE  $B(m, 0.5)$

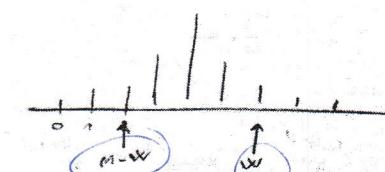
THE P-VALUE IS TRIVIALLY DEFINED ACCORDINGLY



$\Rightarrow$  IF  $W^* = \max(W, m-W)$   
THE P-VALUE IS THE PROBABILITY OF HAVING MORE (OR EQUAL) THAN  $W^*$  HEADS OR MORE (OR EQUAL) THAN  $W^*$  TAILS IN  $m$  INDEPENDENT TOSSES OF A FAIR COIN.

$$P = 2 \cdot \sum_{k=W^*}^m \binom{m}{k} 0.5^m$$

(taken from binomial distribution)



$$\begin{aligned} n-W &\stackrel{H_0}{\sim} Bi\left(n, \frac{1}{2}\right) \\ W &\stackrel{H_0}{\sim} Bi\left(n, \frac{1}{2}\right) \end{aligned}$$

} under  $H_0$  both  $W$  and  $n-W$  are identically distributed  
 $\Rightarrow$  the inference made on the two is the same

The good thing of this test is that we don't have any assumption on the distribution of the data. The test is valid (and exact) for  $\mu$  distribution of data. However there are some common problems in the nonparametric case:

### (NIB) THE DISTRIBUTION OF $W$ IS DISCRETE

$\Rightarrow$  THE SAMPLE SIZE  $m$  DETERMINES THE SO-CALLED ACHIEVABLE SIGN LEVELS OF THE TEST

and so we have some issues with the computation of the quantiles, since it's not always possible to split data in proportions that we want (2 data? We can't have  $\frac{1}{3}, \frac{2}{3}$ )

(i.e. LEVELS OF SIGNIFICANCE FOR WHICH AN SIGN TEST CAN BE BUILT WITHOUT A RANDOMIZATION STRATEGY)

If we reject 10%, we would like to have a tail of 5% and an other tail of 5%. However we cannot do that. We can take a tail of 0%, or 15.67%, or 31.27%, or 3.17%, not others.  
 $\Rightarrow$  Only 4 possible values.  
 These values are the ACHIEVABLE LEVELS of the test (/reachable levels)  
 $=$  values for which we can give an exact test  
 $(P(\text{reject } H_0) = \text{desired value})$

Ex  $m=5$

$$\begin{aligned} P(W=0) &= 3.1\% \\ P(W=1) &= 15.6\% \\ P(W=2) &= 31.2\% \\ P(W=3) &= 31.2\% \\ P(W=4) &= 15.6\% \\ P(W=5) &= 3.1\% \end{aligned}$$

ACHIEVABLE LEVELS OF  $\alpha$  ARE

$$\begin{aligned} \alpha &= 0\% \\ \alpha &= 6.2\% = 2 \cdot 3.1\% \rightarrow \text{reject if } W=0 \text{ or } W=5 \\ \alpha &= 31.4\% = 2 \cdot (3.1\% + 15.6\%) \rightarrow \text{reject if } W=0, 5, 1, 4 \\ \alpha &= 100\% \end{aligned}$$

IN ALL OTHER CASES \*

-  $\alpha$  CAN BE FLOORED TO LOWER AND CLOSEST ACHIEVABLE VALUE (OBTAINING A CONSERVATIVE TEST)

- A RANDOMIZATION STRATEGY CAN BE PUT ON TOP OF THE PROCEDURE \*

### (NIB) IF THE DISTRIBUTION OF $X$ IS SYMMETRIC $\Rightarrow \text{MED}(X) = \text{MEAN}(X)$ AND THE MEAN EXISTS

In this case the sign test is also a test for the mean.

THIS IS GENERALLY NOT TRUE

a.  $X \sim \text{Exp}(\lambda)$

$$\text{MEAN}(X) = 1/\lambda$$

$$\text{MEDIAN}(X) = \log 2/\lambda = \frac{0.693}{\lambda}$$

### (NIB) FOR LARGE VALUE OF $m$ : $B(m, 0.5) \approx N(m/2, m/4)$

### (NIB) $W$ CAN BE ALSO COMPUTED AS

$$W = \frac{m + \sum_{i=1}^m \text{Sign}(X_i - c_0)}{2}$$

### (NIB) ONE-SIDED TESTS CAN BE TRIVIALLY BUILT FOCUSING ON THE CORRESPONDING TAIL.

Ex RIGHT-SIDED

$$H_0: \text{MED}(X) = c_0 \quad H_1: \text{MED}(X) > c_0$$

$$[H_0: P(X > c_0) = 0.5 \quad H_1: P(X > c_0) > 0.5]$$

$$W = \sum_{i=1}^m \mathbb{1}_{\{X_i > c_0\}}$$

$$P = \sum_{k=W}^m \binom{m}{k} 0.5^m$$

ACHIEVABLE LEVELS ARE DIFFERENT FROM THE ONE-SIDED CASE

### (NIB) EXTENSIONS TO THE USE OF DISCRETE OR ORDINAL DATA ARE POSSIBLE. (REMOVAL OF $X_i = c_0$ )

\*  $\alpha = 10\% \rightarrow$  left tail and right tail 5%, we take 3.1%  $\rightarrow$  we obtain a value of 6.2%, which is not exact but  $P(\text{reject } H_0 | H_0 \text{ true}) < 10\%$ , even if we would like to have  $\alpha = 10\%$

\* Used when we want to make continuous a discrete variable.

If we want to obtain  $\alpha = 10\%$ , we consider:

- reject if  $W=0$  (3.1%)
- reject if  $W=5$  (3.1%)
- if  $W=1$  or  $W=4$  we toss a coin with a given probability s.t. the general hypothesis lead us to a global probability of rejection of 10%.

## TWO-SAMPLE PAIRED SIGN TESTS

ASS.  $X_1 - Y_1, X_2 - Y_2, \dots, X_m - Y_m$  r.v. iid  $\mathbb{Z} \in X - Y$  (CONTINUOUS)

$$H_0: \text{MED}(X - Y) = 0 \quad H_1: \text{MED}(X - Y) \neq 0$$

$$[H_0: P(X > Y) = 0.5 \quad H_1: P(X > Y) \neq 0.5]$$

Test statistic:

$$W = \sum_{i=1}^m I(X_i > Y_i) = \frac{m + \sum_{i=1}^m \text{Sign}(X_i - Y_i)}{2}$$

(NB) THE TEST CAN BE ALSO NOT CENTERED ON "0"

(NB) WE ARE NOT TESTING THE EQUALITY IN DISTRIBUTION OF  $X$  AND  $Y$

BUT THE HYPOTHESIS THAT IN A RANDOM PAIR  $X$  HAS THE SAME PROBABILITY OF BEING LARGER OR SMALLER THAN  $Y$ .

(OF COURSE WITH OTHER ASS, FOR EXAMPLE SYMMETRY, SO WE COULD END UP TESTING MEAN, OR OTHER...)

### GENERAL COMMENTS

- SIGN TESTS DO NOT REQUIRE NORMALITY → because it's not like the t-test, we're not comparing means!
- " " ARE ROBUST TO:
  - VIOLATION OF NORMALITY ( $\neq$  t-test)
  - PRESENCE OF OUTLIERS
- SIGN TEST ARE LESS POWERFUL THAN THE "OPTIMAL" PARAMETRIC TEST.

### R EXAMPLE SIGN-TEST.R (Example 2)

the presence of outliers influences the mean, but here we're not dealing with the mean: an outlier simply means that  $W \rightarrow W + 1$

Weakness: the fact that  $X = 1 > 0$  or  $X = 1.000.000 > 0$  doesn't change at all the  $W$  statistics, which in all cases take a value of 1 (which however helps with the robustness with outliers).

This test's power is also a weakness: this test is not able to distinguish a major deviation from 0 from a minor one.

This weakness pushed scientists to work out a solution which takes into account the magnitude of deviation, not only the sign.

They looked for something:

- distribution independent
- taking into account the magnitude of deviation

→ this lead to:

1. MANN-WHITNEY TEST
2. WILCOXON TEST

\* The way to take into account the magnitude of the deviation without introducing assumptions on the distributions is through "contests". We create a contest between the X team and the Y team. We have  $n_1$  players for the first team and  $n_2$  players for the second. We build all the possible contests: (how many contests?  $n_1 \cdot n_2$ )

$$\begin{array}{c} X \\ \text{---} \\ \boxed{\begin{array}{|c|c|c|c|c|} \hline & \text{X}_1 & \text{X}_2 & \dots & \text{X}_{n_1} \\ \hline \end{array}} \end{array} \Rightarrow U_1 = \# \text{ contests won by } X \Rightarrow \frac{U_1}{n_1 \cdot n_2} = \text{max likelihood estimator of } P(X > Y)$$

$$\begin{array}{c} Y \\ \text{---} \\ \boxed{\begin{array}{|c|c|c|c|c|} \hline & \text{Y}_1 & \text{Y}_2 & \dots & \text{Y}_{n_2} \\ \hline \end{array}} \end{array} \Rightarrow U_2 = \# \text{ contests won by } Y \Rightarrow \frac{U_2}{n_1 \cdot n_2} = \text{max likelihood est. of } P(Y > X)$$

### RANK TESTS

non-parametric counterpart for the t-test for two populations

#### HANN-WHITNEY U-TEST (TWO-SAMPLE RANK SUM TEST) (1945 WILLCOXON)

(1942 MANN & WHITNEY)

ASS.  $X_1, X_2, \dots, X_{m_1} \sim \text{iid } X$  (CONTINUOUS)  
 $Y_1, Y_2, \dots, Y_{m_2} \sim \text{iid } Y$  (CONTINUOUS)  $\rightarrow$  INDEPENDENT.

$H_0: P(X > Y) = 0.5 \quad H_1: P(X > Y) \neq 0.5$  probability that the first random variable (X) is larger than the second (Y)

Test statistics:  $m_1, m_2$

$$U_1 = \sum_{i=1}^{m_1} \sum_{j=1}^{m_2} \mathbb{1}\{X_i > Y_j\} = m_1 m_2 \hat{P}(X > Y) \quad \left[ \begin{array}{l} \text{NUMBER OF PAIRWISE "CONTESTS"} \\ \text{WON BY THE FIRST SAMPLE} \end{array} \right]$$

$$U_2 = \sum_{i=1}^{m_1} \sum_{j=1}^{m_2} \mathbb{1}\{Y_j > X_i\} = m_1 m_2 \hat{P}(Y > X) \quad \left[ \begin{array}{l} \text{NUMBER OF PAIRWISE "CONTESTS"} \\ \text{WON BY THE SECOND SAMPLE} \end{array} \right]$$

(NB)  $U_1 + U_2 = m_1 m_2 \rightarrow$  THE INFO PROVIDED BY THE TWO TEST

STATISTICS IS THE SAME!  $\rightarrow$  we can choose one or the other to make inference

$U_1$  AND  $U_2$  CAN BE COMPUTED IN A MORE EFFICIENT WAY AS FOLLOWS:

However sometimes it gets computationally heavy, so we proceed in another way to find  $U_1$  and  $U_2$ :

- DEFINE  $\underline{X} = (X_1, X_2, \dots, X_{m_1}, Y_1, Y_2, \dots, Y_{m_2})$ ,  $\dim(\underline{X}) = n_1 + n_2$

$\uparrow$   
POOLED SAMPLE  
 = concatenation of  
 = the samples from X and Y

- COMPUTE THE RANK OF THE  $m_1 + m_2$  DATA WRT THE POOLED SAMPLE

If a datum is ranked #1 then it's the worst in the game (opposite than intuitively)

$$\underline{R} = (\pi(X_1), \dots, \pi(X_{m_1}), \pi(Y_1), \dots, \pi(Y_{m_2}))$$

$$\pi(X_i) = \sum_{x \in \underline{X}} \mathbb{1}\{x \leq X_i\} \quad i=1, \dots, m_1 \quad \left[ \begin{array}{l} \text{NR OF WINS OR TIES OF} \\ X_i \text{ AGAINST } \underline{X} \end{array} \right]$$

$$\pi(Y_i) = \sum_{y \in \underline{X}} \mathbb{1}\{y \leq Y_i\} \quad j=1, \dots, m_2 \quad \left[ \begin{array}{l} \text{NR OF WINS OR TIES OF} \\ Y_i \text{ AGAINST } \underline{X} \end{array} \right]$$

- COMPUTE THE SUM OF THE RANKS OF THE TWO SAMPLES.

$$\left( \begin{array}{l} R_1 = \sum_{i=1}^{m_1} \pi(X_i) = \sum \text{ranks of the 1st sample} \\ R_2 = \sum_{i=1}^{m_2} \pi(Y_i) = \sum \text{ranks of the 2nd sample} \end{array} \right)$$

$$U_1 = R_1 - \frac{m_1(m_1+1)}{2}$$

$$U_2 = R_2 - \frac{m_2(m_2+1)}{2}$$

these formulas are the ones that give the name "rank sum test"

# wins/ties  
of the first sample  
against the pooled  
sample

# wins/ties of the  
first sample against  
itself

$\Rightarrow U_i \quad (i=1,2)$  is simply  
computed by adjusting  $R_i$

$$(\ast\ast) : \frac{(n_1+n_2)(n_1+n_2+1)}{2} = R_1 + R_2 = \sum \text{of all the ranks}$$

Note: Under  $H_1$   $\underline{X}$  is not a "proper" sample since  $n_1$  elements are from 1 distribution (X) and  $n_2$  are from the 2nd (Y). Under  $H_0$  instead it's drawn from the same distribution.

$H_0$  true

- ⇒ the distribution of  $X$  and  $Y$  are the same. In this case data are iid and so each datum has the same probability to be ranked in  $\mathbf{k}$  position (this happens for all the distributions  $X$  and  $Y$ )
- ⇒ under  $H_0$  the distribution of  $U_1$  and  $U_2$  do not depend on the distribution of the data!

WE DO NOT HAVE AN EXPLICIT FORMULATION OF THE DISTRIBUTION OF  $U_1$  (OR  $U_2$ ) UNDER  $H_0$ !

HOWEVER THE DISTRIBUTION IS ALWAYS THE SAME:

- JUST DEPENDING ON  $m_1$  AND  $m_2$
- NOT DEPENDING ON THE (SAME) DISTRIBUTION OF  $X$  AND  $Y$   
(PROOF:  $\pi(x_i) \stackrel{H_0}{\sim} U(1, \dots, m_1 m_2)$ )

⇒ QUANTILE CAN EASILY ESTIMATED VIA A ONE-TIME MC SIMULATION

MOREOVER UNDER  $H_0$ :

- $U_1$  AND  $U_2$  HAS THE SAME (DISCRETE) DISTRIBUTION ( $U_1 - \frac{m_1 m_2}{2} = -(U_2 - \frac{m_1 m_2}{2})$ )
- $U_1$  AND  $U_2$  ARE SYMMETRICALLY DISTRIBUTED AROUND  $\frac{m_1 m_2}{2}$
- $E(U_1) = E(U_2) = \frac{m_1 m_2}{2}$
- $\text{VAR}(U_1) = \text{VAR}(U_2) = \frac{m_1 m_2 (m_1 + m_2 + 1)}{12}$

$U_1$  and  $U_2$  depends only on  $n_1$  and  $n_2$  (under  $H_0$ )

⇒ CRITICAL REGIONS CAN BE COMPUTED ACCORDINGLY.

⇒ P-VALUE

$$\begin{aligned} \text{RELYING ON THE QUANTITY } U^* &= \max(U_1, U_2) = \\ &= \max(U_1, m_1 m_2 - U_1) \\ &= \max(U_2, m_1 m_2 - U_2) \end{aligned}$$

(NB) BEING  $U_1$  (AND  $U_2$ ) THE SUM OF EXCHANGEABLE RANDOM VARIABLES (i.e. THE RANKS) WITH FINITE MEAN AND VARIANCE AND VANISHING COVARIANCES THE CLT HOLDS:

$$\frac{U_1 - \frac{m_1 m_2}{2}}{\sqrt{\frac{m_1 m_2 (m_1 + m_2 + 1)}{12}}} \xrightarrow[m_1, m_2 \rightarrow \infty]{} N(0, 1)$$

NORMAL APPROXIMATIONS  
OF THE QUANTILES.

(NB) ONE-SIDED TESTS CAN BE TRIVIALLY BUILT RELYING ON THE CORRESPONDING TAIL.

(NB) THE U-TEST CAN BE ADAPTED ALSO TO DISCRETE AND ORDINAL DATA (i.e. AVERAGE RANKS)

(NB) THE U-TEST IS INVARIANT UNTIL MONOTONIC TRANSFORMATION OF DATA. (PROOF: RANKS ARE INVARIANT) NO NEED FOR Box-Cox transformation.

EXAMPLE MANN-WHITNEY U TEST. R

proposal that tries to solve the issue of "weighting" different magnitude of deviations from  $C_0$   
(the idea starts from the sign test)

### WILCOXON SIGNED-RANK TEST

(WILCOXON 1945)

ONE-SAMPLE TWO-SIDED TEST (CENTERED ON "0")

IT IS AN EXTENSION OF THE SIGN-TEST IN WHICH ALSO THE MAGNITUDE (i.e. RANK OF THE ABSOLUTE VALUE) IS TAKEN INTO CONSIDERATION

ASS.  $X_1, X_2, \dots, X_m \sim \text{iid } X$  (CONTINUOUS)

$H_0: P(X > 0) = 0.5$      $H_1: P(X > 0) \neq 0.5$  ← same test of the sign test  
what changes is the test statistics:

$$W^+ = \sum_{i=1}^m \mathbb{I}_{\{X_i > 0\}} R(|X_i|)$$

↓  
RANK OF THE ABSOLUTE VALUES

and SUM OF THE "ABSOLUTE" RANKS OF POSITIVE VALUES.

the weights depend on the magnitude of deviation of data

- UNDER THE  $H_0$  (INDEPENDENTLY ON THE DISTRIBUTION OF  $X$ )

-  $\mathbb{I}_{\{X_i > 0\}} \stackrel{\text{iid}}{\sim} B(1, 0.5)$

-  $R(|X_i|) \stackrel{\text{exc}}{\sim} U(\{1, 2, \dots, m\})$  } INDEPENDENT

SO THE DISTRIBUTION OF  $W^+$  UNDER  $H_0$  DEPENDS ONLY ON  $n$

⇒ QUANTILES CAN BE ESTIMATED VIA A ONE-TIME MC SIMULATION  
(EXPLICIT FORMULA IS NOT AVAILABLE FOR GENERIC  $m$ )

\* exchangeable (+iid):  
a given vector of a random variable is exchangeable if after permutation of components (+permutation)  
we obtain a random vector with exactly the same joint distribution

iid ⇔ exchange.

- UNDER  $H_1$ :

-  $\mathbb{I}_{\{X_i > 0\}} \stackrel{\text{no}}{\sim} B(1, p)$  WITH  $p \neq 0.5$

-  $R(|X_i|) \stackrel{\text{exc}}{\sim} U(\{1, 2, \dots, m\})$  } DEPENDENT

Another way to do the same test:

SIMILARLY AND EQUIVALENTLY:

$$W^- = \sum_{i=1}^m \mathbb{I}_{\{X_i < 0\}} R(|X_i|)$$

and SUM OF THE "ABSOLUTE" RANKS OF THE NEGATIVE VALUES.

$W^+ + W^- = \frac{m(m+1)}{2}$  = summation of the ranks of all data

- $W^+$  AND  $W^-$  HAS THE SAME DISTRIBUTION SYMMETRIC AROUND  $\frac{m(m+1)}{4}$
- SO INFERENCE CAN BE EQUIVALENTLY CARRIED OUT RELYING ON:

If we have that the distribution is around zero then any datum can be ranked anywhere with the same probability. If the distribution has positive median we have that values associated to larger ranking (=associated to larger deviation) are more likely to be positive than negative.

•  $W^+, W^- \sim \text{Poisson}$  } CENTERED ON  $\frac{m(m+1)}{4}$  } P-VALUE BASED  
•  $W^-$  } CENTERED ON "0" } P-VALUE BASED  

$$\begin{aligned} W^+ - W^- &= \frac{m(m+1)}{2} - 2W^- = \sum_{i=1}^m \text{sign}(X_i) R(|X_i|) \\ W^- - W^+ &= \frac{m(m+1)}{2} - 2W^+ = \sum_{i=1}^m \text{sign}(X_i) R(|X_i|) \end{aligned}$$

that's why they're called "sign rank test"

- THE NON-ZERO-CENTERED VERSION CAN BE EASILY IMPLEMENTED  
REPLACING  $X_i \mapsto X_i - c_0$
- general test (we can also change 0.5)

$$H_0: P(X > c_0) = 0.5 \quad H_1: P(X > c_0) \neq 0.5$$

$$W^+ = \sum_{i=1}^m \prod_{\{X_i > c_0\}} R(|X_i - c_0|)$$

- THE ONE-SIDED VERSION CAN BE BUILT RELYING ON THE CORRESPONDING TAIL.

(NB) THE SIGNED-RANK TEST IS USUALLY

- MORE POWERFUL THAN THE SIGN-TEST
- . LESS OPTIMAL PARAMETRIC TEST

(NB) IT IS LESS ROBUST THAN THE SIGN-TEST. (w.r.t. outliers)  
BUT MORE ROBUST THAN THE OPTIMAL PARAMETRIC TEST  
(REBALANCING OF THE ABSOLUTE RANKS)

TWO-SAMPLE (PAIRED) SIGNED-RANK TEST (CENTERED ON 0)

ASS.  $X_1 - Y_1, X_2 - Y_2, \dots, X_m - Y_m \sim \text{ind } Z = X - Y$  (CONTINUOUS)

$$H_0: P(X > Y) = 0.5 \quad H_1: P(X > Y) \neq 0.5$$

$$[ \text{MED}(X - Y) = 0 ] \quad [ \text{MED}(X - Y) \neq 0 ]$$

$$W^+ = \sum_{i=1}^m \prod_{\{X_i > Y_i\}} R(|X_i - Y_i|) \rightsquigarrow \begin{matrix} \text{SUM OF THE ABSOLUTE RANKS} \\ \text{OF THE PAIRS IN WHICH} \\ X_i > Y_i \end{matrix}$$

### R EXAMPLE WILCOXON SIGNED-RANK TEST - R

Conditions:

we know the distributions under  $H_0$  and  $H_1$ .  
(moreover an optimal test can be identified)

Ranking with discrete data?

If  $x_i$  is ranked 3 and  
 $x_j$  is ranked 4

$\Rightarrow x_i, x_j$  both ranked 3.5

Sign rank test with discrete data?

If we have a 0 (an observation which is 0) ~~we cannot assign a sign~~  $\Rightarrow$  we simply remove the data