

S 145 (M 7a)

Enige opmerkingen over de verdeling van de  
toetsingsgrootte van WILCOXON

Stichting Mathematisch Centrum

Statistische Afdeling

1954

Enige opmerkingen over de verdeling van de  
 toetsingsgrootte van WILCOXON.

In memorandum S 47 (M 7) is vermeld dat de verdeling van de toetsingsgrootte  $\underline{U}$  van WILCOXON, onder de hypothese  $H_0$ , dat de beide steekproeven afkomstig zijn uit dezelfde waarschijnlijkheidsverdeling, voor grote  $n$  en  $m$  bij benadering normaal is. Deze benadering geldt over het algemeen als er weinig of geen gelijke waarnemingen in beide steekproeven voorkomen. Als  $n$  en  $m$  klein zijn, kan men, zo er weinig gelijke waarnemingen voorkomen gebruik maken van exacte tabellen van de verdeling van  $\underline{U}$  (litt. [1], [2]). Als dan echter de gegeven waarnemingsreeksen grote groepen van gelijke waarnemingen bevatten, zijn de exacte tabellen en de normale benadering niet altijd voldoende nauwkeurig. We zullen laten zien hoe men voor dergelijke gevallen de verdeling van  $\underline{U}$  exact kan berekenen.

De gegeven waarnemingsreeksen kan dan worden samengevat als aangegeven in tabel I.

Tabel I

Voorkomende steekproef waarden	Aantal malen dat $z_i$ optreedt bij		totaal
	$x$	$y$	
$z_1$	$a_1$	$b_1$	$t_1$
$z_2$	$a_2$	$b_2$	$t_2$
⋮			
$z_k$	$a_k$	$b_k$	$t_k$
totaal	$n$	$m$	$N$

Hierin zijn  $z_1, \dots, z_k$  de voorkomende steekproefwaarden, gerangschikt naar opklimmende grootte;  $n$  en  $m$  de steekproef-groottes;  $t_1, t_2, \dots, t_k$  de grootten der groepen gelijke waarnemingen en  $a_i$  (resp.  $b_i$ ) het aantal malen dat de waarde  $z_i$  in de eerste (resp. tweede) steekproef optreedt.

De kans dat de combinatie, gegeven in tabel I, optreedt onder de hypothese  $H_0$ , en onder de voorwaarden  $t_1 = t_1, t_2 = t_2, \dots, t_k = t_k$  wordt nu:

$$(1) \quad P[a_1 = a_1, a_2 = a_2, \dots, a_k = a_k | t_1, t_2, \dots, t_k; H_0] = \frac{\binom{t_1}{a_1} \cdot \binom{t_2}{a_2} \cdot \dots \cdot \binom{t_k}{a_k}}{\binom{N}{n}} \quad 1)$$

Nu worden voor de gegeven waarde van  $n$ ,  $m$  en  $t_1, t_2, \dots, t_k$  alle andere mogelijke combinaties van de  $a$ 's en de  $b$ 's opgeschreven en voor ieder hiervan de  $\tilde{u}$  berekend volgens

$$(2) \quad \tilde{u} = \frac{1}{2} \sum_{i < j} (a_i b_j - a_j b_i).$$

Voorbeeld: Gegeven 4 waarnemingen van een stochastische grootte  $x$  en 11 waarnemingen van een stochastische grootte  $y$ . De aangenomen steekproefwaarden zijn gerangschikt naar opklimmende grootte  $z_1, z_2, z_3$  en  $z_4$  en in de twee steekproeven tezamen genomen treden deze waarden resp. 1, 1, 2 en 11 maal op. Dus (zie tabel I):  $n = 4, m = 11, t_1 = t_2 = 1, t_3 = 2, t_4 = 11, N = 15$ .

Wij geven hieronder alle mogelijke combinaties met de corresponderende waarde van  $\tilde{u}$ .

1) $\begin{array}{ c } \hline 0 & 1 \\ \hline 0 & 1 \\ \hline 0 & 2 \\ \hline 4 & 7 \\ \hline -8 \\ \hline \end{array}$	2) $\begin{array}{ c } \hline 0 & 1 \\ \hline 0 & 1 \\ \hline 1 & 1 \\ \hline 3 & 8 \\ \hline -1\frac{1}{2} \\ \hline \end{array}$	3) $\begin{array}{ c } \hline 0 & 1 \\ \hline 1 & 0 \\ \hline 0 & 2 \\ \hline 3 & 8 \\ \hline 0 \\ \hline \end{array}$	4) $\begin{array}{ c } \hline 1 & 0 \\ \hline 0 & 1 \\ \hline 0 & 2 \\ \hline 3 & 8 \\ \hline +1 \\ \hline \end{array}$	5) $\begin{array}{ c } \hline 0 & 1 \\ \hline 0 & 1 \\ \hline 2 & 0 \\ \hline 2 & 9 \\ \hline +5 \\ \hline \end{array}$	6) $\begin{array}{ c } \hline 0 & 1 \\ \hline 1 & 0 \\ \hline 1 & 1 \\ \hline 2 & 9 \\ \hline +6\frac{1}{2} \\ \hline \end{array}$
7) $\begin{array}{ c } \hline 1 & 0 \\ \hline 0 & 1 \\ \hline 1 & 1 \\ \hline 2 & 9 \\ \hline +7\frac{1}{2} \\ \hline \end{array}$	8) $\begin{array}{ c } \hline 1 & 0 \\ \hline 1 & 0 \\ \hline 0 & 2 \\ \hline 2 & 9 \\ \hline +9 \\ \hline \end{array}$	9) $\begin{array}{ c } \hline 0 & 1 \\ \hline 1 & 0 \\ \hline 2 & 0 \\ \hline 1 & 10 \\ \hline +13 \\ \hline \end{array}$	10) $\begin{array}{ c } \hline 1 & 0 \\ \hline 0 & 1 \\ \hline 2 & 0 \\ \hline 1 & 10 \\ \hline +14 \\ \hline \end{array}$	11) $\begin{array}{ c } \hline 1 & 0 \\ \hline 1 & 0 \\ \hline 1 & 1 \\ \hline 1 & 10 \\ \hline +15\frac{1}{2} \\ \hline \end{array}$	12) $\begin{array}{ c } \hline 1 & 0 \\ \hline 1 & 0 \\ \hline 2 & 0 \\ \hline 1 & 10 \\ \hline +22 \\ \hline \end{array}$

$$1) \quad \binom{t}{a} = \frac{t!}{a! (t-a)!}$$

Tabellen voor deze binomiaal coëfficiënten vindt men o.a. in [3].

We berekenen vervolgens de waarschijnlijkheden van al deze combinaties volgens formule (1) (zie tabel II, kolom 5). Deze waarschijnlijkheidsverdeling is in figuur 1 getekend.

De linkseenzijdige overschrijdingskans van een gevonden waarde  $\bar{u}$  van  $\underline{u}$  wordt gedefiniëerd als de som van de waarschijnlijkheden van die waarden van  $\underline{u}$  die kleiner of gelijk zijn aan de gevonden waarde. In tabel II zijn in kolom 6 deze exacte linkseenzijdige overschrijdingskansen gegeven en in kolom 4 de linkseenzijdige overschrijdingskansen die men vindt met behulp van de normale benadering (zie ook figuur 2).

De rechtseenzijdige overschrijdingskans wordt (analoog aan de linkseenzijdige) gedefiniëerd als de som van de waarschijnlijkheden van die waarden van  $\bar{u}$  die groter of gelijk zijn aan de gevonden waarde van  $\bar{u}$ .

Tabel II.

Vergelijking van de linkseenzijdige overschrijdingskans berekend volgens de benadering en de exacte methode.

1	2	3	4	5	6
Combinatie	$\bar{u}$	$\frac{\bar{u}-\underline{u}}{\sigma}$	linkseenzijdige overschrijdingskans m.b.v. de normale benadering	exacte waarde van de combinaties	exacte linkseenzijdige overschrijdingskans
1	- 8	-1,34	0,090	0,242	0,242
2	- 1 $\frac{1}{2}$	-0,25	0,401	0,242	0,484
3	0	0	0,500	0,121	0,605
4	+ 1	+0,17	0,568	0,121	0,726
5	+ 5	+0,25	0,800	0,040	0,706
6	+ 6 $\frac{1}{2}$	+1,09	0,862	0,080	0,847
7	+ 7 $\frac{1}{2}$	+1,26	0,896	0,080	0,927
8	+ 9	+1,51	0,935	0,040	0,967
9	+13	+2,18	0,985	0,008	0,975
10	+14	+2,35	0,991	0,008	0,983
11	+15 $\frac{1}{2}$	+2,60	0,995	0,016	0,999
12	+22	+3,69	1,000	0,001	1,000

Zoals duidelijk blijkt is voor enkele van de genoemde combinaties de overschrijdingskans berekend volgens de methode van de normale benadering belangrijk verschillend van die gevonden met de exacte methode.

Een tweezijdige kritieke zone met onbetrouwbaarheidsdrempel  $\alpha$  kan hier op verschillende manieren gedefiniëerd worden:

B.v. door een linkseenzijdige en een rechtseenzijdige kritieke zone te kiezen, ieder met onbetrouwbaarheidsdrempel  $\frac{1}{2} \alpha$ . Een bezwaar van deze methode is dat hij vaak leidt tot een groot verschil tussen  $\alpha$  en de werkelijke onbetrouwbaarheid. Dit is vooral dan in sterke mate het geval, als de kleinste of de grootste waarde, die  $\tilde{u}$  aan kan nemen, een waarschijnlijkheid bezit, die  $> \frac{1}{2} \alpha$  is. Dan wordt de kritieke zone bovendien éézijdig. Als we b.v. in ons voorbeeld nemen  $\alpha = 0,20$  dan bestaat de op deze wijze gedefiniëerde "tweezijdige" kritieke zone uit de waarden  $\tilde{u} = 9, 13, 14, 15\frac{1}{2}$  en 22 en de onbetrouwbaarheid is  $0,040 + 0,008 + 0,008 + 0,016 + 0,001 = 0,073$ .

Een andere methode is, dat men de waarden van  $\tilde{u}$  met de kleinste waarschijnlijkheden bij elkaar zoekt totdat de gekozen onbetrouwbaarheidsdrempel het toevoegen van een nieuwe waarde verhindert. Deze methode heeft het voordeel boven de vorige dat in het algemeen de onbetrouwbaarheid dichterbij  $\alpha$  ligt, maar bij verdelingen met meer dan één top het nadeel dat de kritieke zone geen aaneengesloten geheel is. Nemen we in ons voorbeeld weer  $\alpha = 0,20$  dan bestaat deze tweezijdige kritieke zone uit de waarden  $\tilde{u} = 5, 9, 13, 14, 15\frac{1}{2}$  en 22 en de onbetrouwbaarheid is  $0,040 + 0,040 + 0,008 + 0,008 + 0,016 + 0,001 = 0,113$ . De waarde  $\tilde{u} = 5$  behoort dus wel tot de kritieke zone maar  $\tilde{u} = 6\frac{1}{2}$  en  $\tilde{u} = 7\frac{1}{2}$  niet en men kan deze kritieke zone bezwaarlijk tweezijdig noemen.

We definiëren nu de tweezijdige kritieke zone als volgt: We beginnen met van de kleinste en de grootste waarde van  $\tilde{u}$  diegene te nemen met de kleinste waarschijnlijkheid. In ons voorbeeld dus  $\tilde{u} = 22$ . Dan kiezen we één van de waarden  $\tilde{u} = -8$  en  $\tilde{u} = 15\frac{1}{2}$  en wel zodanig dat het verschil tussen de som van de waarschijnlijkheden links en de som van de waarschijnlijkheden rechts in absolute waarde minimaal is. Nemen we  $\tilde{u} = -8$  dan is dit absolute verschil  $0,242 - 0,001 = 0,241$  en nemen we  $\tilde{u} = 15\frac{1}{2}$  dan is dit absolute verschil  $0,016 + 0,001 = 0,017$ . We nemen dus  $\tilde{u} = 15\frac{1}{2}$ .

Vervolgens kiezen we één van de waarden  $\tilde{u} = 14$  en  $\tilde{u} = -8$ . Nemen we  $\tilde{u} = -8$  dan wordt het genoemde verschil tussen links en rechts in absolute waarde:  $0,242 - (0,016 + 0,001) = 0,225$  en bij  $\tilde{u} = 14$ ;  $0,008 + 0,016 + 0,001 = 0,025$ . We kiezen weer die waarde van  $\tilde{u}$  die het kleinste absolute verschil tussen

links en rechts geeft, dus  $\tilde{u} = 14$ . Zo gaan we door met telkens links of rechts een waarde van  $\tilde{u}$  bij de kritieke zone te nemen totdat de gekozen onbetrouwbaarheidsdrempel het toevoegen van een nieuwe waarde verhindert. Nemen we weer  $\alpha = 0,2$  dan vinden we in ons voorbeeld voor deze kritieke zone de waarden  $\tilde{u} = 7\frac{1}{2}, 9, 13, 14, 15\frac{1}{2}$  en 22 en de onbetrouwbaarheid is  $0,080 + 0,040 + 0,008 + 0,008 + 0,016 + 0,001 = 0,153$ . Op deze wijze bereikt men dus dat de onbetrouwbaarheid in het algemeen dichterbij  $\alpha$  ligt dan volgens de eerste methode; bovendien verkrijgt men een aaneengesloten geheel voor het linker- en rechter deel der kritieke zone, terwijl de onbetrouwbaarheid zo goed als mogelijk is in gelijke mate over de twee delen verdeeld wordt.

De tweezijdige overschrijdingskans wordt nu gedefiniëerd als de onbetrouwbaarheid van de kleinste kritieke zone van deze aard, die de gevonden waarde van  $\tilde{u}$  bevat. B.v. de kleinste kritieke zone die het punt  $\tilde{u} = 6\frac{1}{2}$  bevat bestaat uit de waarden  $\tilde{u} = 22, 15\frac{1}{2}, 14, 13, 9, 7\frac{1}{2}, -8$  en  $6\frac{1}{2}$  en de tweezijdige overschrijdingskans van  $\tilde{u} = 6\frac{1}{2}$  is dus in ons voorbeeld de som van de waarschijnlijkheden van deze waarden van  $\tilde{u}$  en dus 0,476.

In tabel II zijn voor ieder der mogelijke waarden van  $\tilde{u}$  deze exacte en tevens de benaderde tweezijdige overschrijdingskansen gegeven.

Tabel III.

Exacte en benaderde tweezijdige overschrijdingskansen.

$\tilde{u}$	tweezijdige overschrijdingskans	
	exact	benaderd
-8 (7)*	0,395	0,180
$-1\frac{1}{2}$ (11)	0,879	0,802
0 (12)	1,000	1,000
1 (10)	0,637	0,864
5 (9)	0,516	0,400
$6\frac{1}{2}$ (8)	0,476	0,276
$7\frac{1}{2}$ (6)	0,153	0,208
9 (5)	0,073	0,130
13 (4)	0,033	0,030
14 (3)	0,025	0,018
$15\frac{1}{2}$ (2)	0,017	0,010
22 (1)	0,001	<0,001

\* De cijfers tussen haakjes geven de volgorde aan, waarin de waarden van  $\tilde{u}$  bij de kritieke zones getrokken worden.

Om te illustreren welke invloed de waarden van  $n$ ,  $m$ ,  $t_1$ ,  $t_2$ , ...,  $t_k$  op het resultaat hebben geven we hier nog een tweetal andere voorbeelden met de tweezijdige overschrijdingskansen.

A.	54	32	86
	3	5	8
	0	2	2
	57	39	96

$\tilde{u} = + 139$   
 exact:  $k = 0,06$ .  
 benaderd  $k = 0,04$

B.	13	19	32
	1	4	5
	0	2	2
	14	25	39

$\tilde{u} = + 30\frac{1}{2}$   
 exact  $k = 0,23$   
 benaderd  $k = 0,18$

In het algemeen geldt dat de exacte verdeling van  $\tilde{u}$  symmetrisch is als  $n = m$  en ook als  $t_1 = t_k$ ,  $t_2 = t_{k-1}$ ,  $t_3 = t_{k-2}$ , enz. (zie Tabel I). In deze gevallen zal voor niet te kleine  $n$  en  $m$ , de normale benadering meestal vrij goed zijn.

#### Literatuur:

- [1] H.R.VAN DER VAART: Gebruiksaanwijzing voor de toets van Wilcoxon, Rapport S 32 (M 4) van het Mathematisch Centrum, Statistische Afdeling.
- [2] C.WHITE: The use of ranks in a test of significance for comparing two treatments Biometrics, Vol.8, number 1., March 1952.
- [3] T.C.FRY: Probability and its engineering uses D.van Nostrand Company, New York 1928.
- [4] Rapport S 47 (M7) van de Statistische Afdeling van het Mathematisch Centrum.

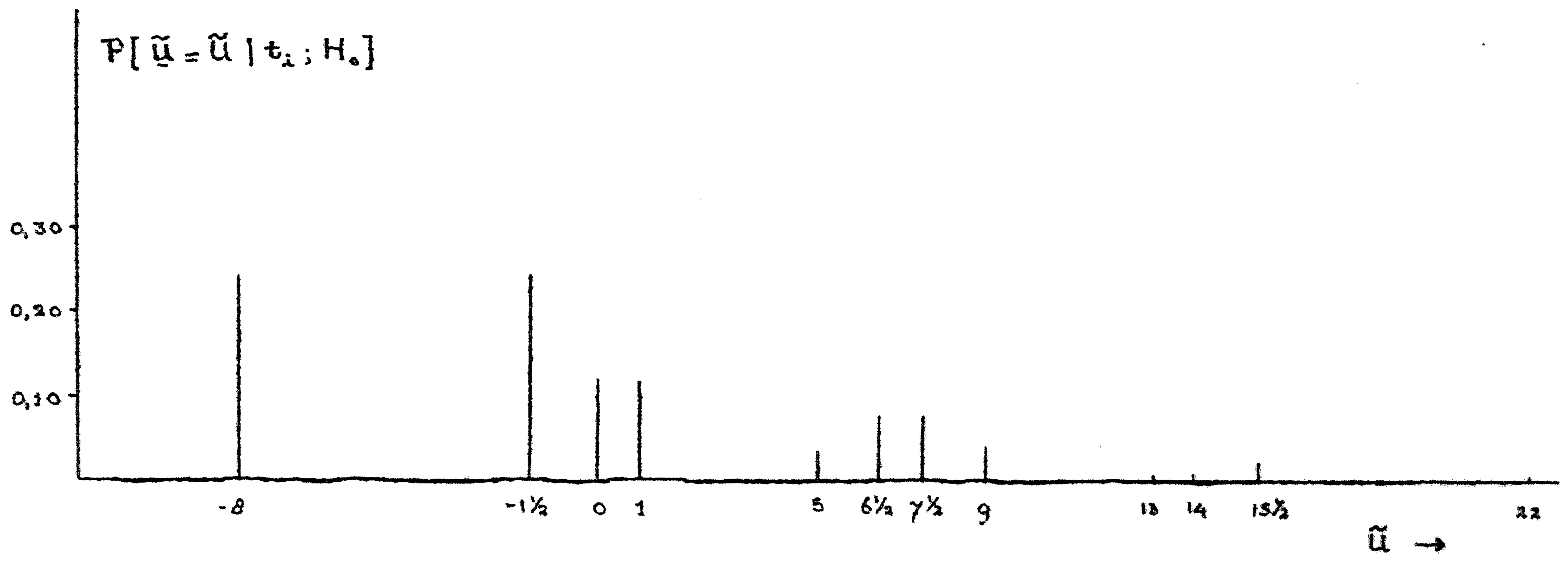


Fig. 1

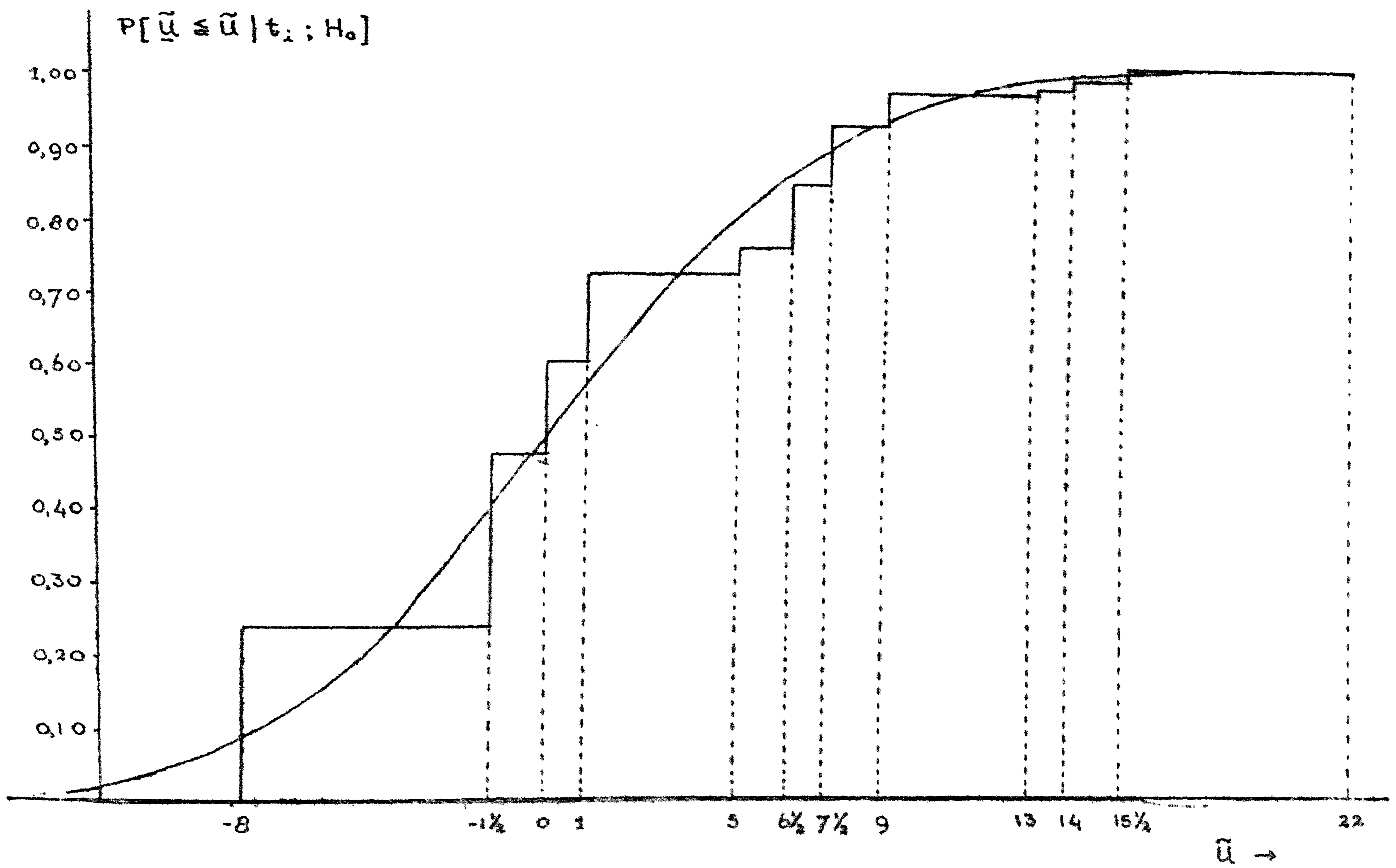


Fig. 2

$P[\tilde{u} = \tilde{u} | t_i; H_0]$  en  $P[\tilde{u} \leq \tilde{u} | t_i; H_0]$  ,  $n=4$  ,  $m=11$  ,  $t_1=t_2=1$  ,  $t_3=2$  ,  $t_4=11$ .