
Neparametarski testovi za dva nezavisna uzorka

- ▶ Boris Glišić 208/2010
- ▶ Bojana Ružičić 21/2010



Neparametarski testovi

- ▶ Hipoteze o raspodeli obeležja se nazivaju neparametarske hipoteze, a odgovarajući testovi su neparametarski testovi.
- ▶ Kada se ne može sigurnošću utvrditi da li je raspodela jedne grupe podataka normalna, izračunavanje pojedinih parametara i primena parametarskih metoda daju vrlo nepouzidane zaključke. U tim slučajevima se primenjuju neparametarske metode, koje ne zavise od raspodele posmatranog obeležja.



Neparametrski testovi

- ▶ Pored slučaja kada ne znamo raspodelu obeležja, neparametarske umesto parametarskih testova koristimo i u slučajevima:
 - ▶ 1. Kada je zavisna promenljiva u diskretnom obliku (na primer, skala od 5 nivoa - koliko se mušteriji dopada neki proizvod: od “ne, nimalo” do “da, veoma”)
 - ▶ 2. Kada se nezavisna promenljiva sastoji od dve kategoričke, nezavisne grupe (pol muški - ženski, radni status zaposlen - nezaposlen)
 - ▶ 3. Kada su opservacije međusobno zavisne
-

Testovi

- ▶ Testovi o kojima ćemo govoriti su efikasnije alternative za t - test.
- ▶ Mann - Whitney U test
- ▶ Wald - Wolfowitz Runs test
- ▶ Kolmogorov - Smirnov Z test
- ▶ Moses Extreme Reactions test

Baza

63-baza.SAV [DataSet1] - IBM SPSS Statistics Data Editor

File Edit View Data Transform Analyze Direct Marketing Graphs Utilities Add-ons Window Help

Visible: 11 of 11 Variables

	id	pplata	pol	vrema	splata	obrazovanje	kat_posla	rasa	pol_rasa	starost	r_staz	var	var	var	var	var
1	748	4080	1	70	6300	8	1	0	3	63,83	22,92					
2	832	4080	1	74	6360	8	1	0	3	55,25	3,58					
3	754	3900	1	92	6480	8	1	0	3	55,50	,00					
4	869	4080	1	82	6480	12	1	0	3	60,00	,00					
5	969	4080	1	68	6480	12	1	0	3	62,33	15,00					
6	825	4080	1	66	6540	12	1	1	4	60,50	13,58					
7	1083	4080	1	84	6600	12	1	0	3	62,42	24,00					
8	1107	3900	1	88	6660	8	1	0	3	62,50	34,33					
9	886	4080	1	76	6720	8	1	0	3	59,08	6,33					
10	1127	4080	1	72	6780	8	1	0	3	56,92	26,58					
11	921	3600	1	97	6780	12	1	1	4	60,67	10,33					
12	935	4080	1	72	6780	12	1	1	4	51,50	22,58					
13	831	4080	1	85	6840	12	1	0	3	55,08	6,00					
14	940	4080	1	81	6840	8	1	1	4	51,50	,00					
15	1128	4080	1	82	6900	12	1	0	3	53,92	29,83					
16	749	4080	1	81	6960	8	1	0	3	57,83	32,50					
17	647	4080	1	72	6960	12	1	1	4	46,50	9,67					
18	826	4080	1	72	7080	8	1	0	3	61,50	15,33					
19	995	3900	1	86	7260	12	1	0	3	62,00	6,00					
20	1121	4080	1	85	7380	12	1	1	4	51,00	19,00					
21	1034	4200	1	90	7500	15	1	0	3	58,00	4,50					
22	1096	3600	1	96	7680	15	1	1	4	60,50	1,92					
23	784	5220	1	70	7860	8	1	0	3	55,92	8,50					

Data View Variable View

IBM SPSS Statistics Processor is ready

Pristupanje testovima

The screenshot shows the IBM SPSS Statistics Data Editor interface. The main window displays a data table with 23 rows and 11 columns. The 'Analyze' menu is open, and the '2 Independent Samples...' option is selected. The status bar at the bottom indicates '2 Independent Samples...' and 'IBM SPSS Statistics Processor is ready'.

	id	pplata	plata	obrazovanje	kat_posla	rasa	pol_rasa	starost	r_staz	var	var	var	var	var
1	748	408	6300	8	1	0	3	63,83	22,92					
2	832	408	6360	8	1	0	3	55,25	3,58					
3	754	390	6480	8	1	0	3	55,50	,00					
4	869	408	6480	12	1	0	3	60,00	,00					
5	969	408	6480	12	1	0	3	62,33	15,00					
6	825	408	6540	12	1	1	4	60,50	13,58					
7	1083	408	6600	12	1	0	3	62,42	24,00					
8	1107	390	6660	8	1	0	3	62,50	34,33					
9	886	408	6720	8	1	0	3	59,08	6,33					
10	1127	408	6780	8	1	0	3	56,92	26,58					
11	921	360	6780	12	1	1	4	60,67	10,33					
12	935	408			1	1	4	51,50	22,58					
13	831	408			1	0	3	55,08	6,00					
14	940	408			1	1	4	51,50	,00					
15	1128	408						53,92	29,83					
16	749	408	6960	8				57,83	32,50					
17	647	408	6960	12				46,50	9,67					
18	826	408	7080	8				61,50	15,33					
19	995	390	7260	12				62,00	6,00					
20	1121	408	7380	12				51,00	19,00					
21	1034	420	7500	15				58,00	4,50					
22	1096	360	7680	15				60,50	1,92					
23	784	5220	7860	8				55,92	8,50					

Pristupanje testovima

63-baza.SAV [DataSet1] - IBM SPSS Statistics Data Editor

File Edit View Data Transform Analyze Direct Marketing Graphs Utilities Add-ons Window Help

Visible: 11 of 11 Variables

	id	pplata	pol	vreme	splata	obrazovanje	kat_posla	rasa	pol_rasa	starost	r_staz	var	var	var	var	var
1	748	4080	1	70	6300	8	1	0	3	63,83	22,92					
2	832	4080	1	74	6360	8	1	0	3	55,25	3,58					
3	754	3900	1	92							,00					
4	869	4080	1	82							,00					
5	969	4080	1	68							15,00					
6	825	4080	1	66							13,58					
7	1083	4080	1	84							24,00					
8	1107	3900	1	88							34,33					
9	886	4080	1	76							6,33					
10	1127	4080	1	72							26,58					
11	921	3600	1	97							10,33					
12	935	4080	1	72							22,58					
13	831	4080	1	85							6,00					
14	940	4080	1	81							,00					
15	1128	4080	1	82							29,83					
16	749	4080	1	81							32,50					
17	647	4080	1	72							9,67					
18	826	4080	1	72							15,33					
19	995	3900	1	86	7260	12	1	0	3	62,00	6,00					
20	1121	4080	1	85	7380	12	1	1	4	51,00	19,00					
21	1034	4200	1	90	7500	15	1	0	3	58,00	4,50					
22	1096	3600	1	96	7680	15	1	1	4	60,50	1,92					
23	784	5220	1	70	7860	8	1	0	3	55,92	8,50					

Two-Independent-Samples Tests

Test Variable List: [Empty]

Grouping Variable: [Empty]

Test Type:
 Mann-Whitney U
 Kolmogorov-Smirnov Z
 Moges extreme reactions
 Wald-Wolfowitz runs

Buttons: OK, Paste, Reset, Cancel, Help, Exact..., Options..., Define Groups...

Data View Variable View

IBM SPSS Statistics Processor is ready

Mann - Whitney U test

(Wilcoxon rank-sum (WRS))

- ▶ Ovaj test se primenjuje za testiranje hipoteze o jednakosti neprekidnih raspodela za obeležja X i Y na osnovu dva slučajna uzorka (X_1, X_2, \dots, X_m) i (Y_1, Y_2, \dots, Y_n) pri čemu je $n \geq m$.
- ▶ Efikasniji je od t - testa (0.95) kod raspodela koje su različite od normalne i slične efikasnosti kod normalne raspodele.
- ▶ $H_0(F_X(x) = F_Y(x))$; $H_1(F_X(x) \neq F_Y(x))$
- ▶ Pri testiranju se formira objedinjeni uzorak sortiran u neopadajućem poretku. Definiše se

$$h_{ij} = \begin{cases} 1, & Y_j < X_i \\ 0, & Y_j \geq X_i \end{cases}$$

- ▶ Nulta hipoteza se testira uz pomoć statistike $U = \sum_{i=1}^m \sum_{j=1}^n h_{ij}$.

Mann - Whitney U test

- ▶ Važi da je $E(U) = \frac{1}{2}mn$, $D(U) = \frac{1}{12}mn(m+n+1)$
- ▶ Raspodela statistike U se aproksimira normalnom raspodelom na sledeći način

$$Z = \frac{U - mn/2}{\sqrt{mn \cdot (m+n+1)}} \sqrt{12}$$

- ▶ Kritična oblast se određuje iz uslova $P\{Z \leq c_1\} = P\{Z \geq c_2\} = \frac{\alpha}{2}$
- ▶ Drugi način računanja statistike U je da se saberu svi rangovi elemenata X i svi rangovi elemenata Y. Tada se vrednost test - statistike računa na osnovu jedne od ove dve formule:

$$U = R_X - \frac{m(m+1)}{2} \quad U = -R_Y + mn + \frac{n(n+1)}{2}$$

- ▶ Aproksimacija normalnom raspodelom je dobra već za $m, n \geq 8$. Ako su obimi uzoraka manji od 8, koriste se posebne tablice.
-

Primer

- ▶ Na osnovu baze sa početka ispitujemo:
 - ▶ Da li dužina rada zavisi od pola zaposlenih?

 - ▶ Zavisna promenljiva: Seniornost na poslu
 - ▶ Nezavisna: Pol zaposlenog
-

	Name	Type	Width	Decimals	Label	Values	Missing	Columns	Align	Measure
1	id	Numeric	4	0	Identifikacioni broj	None	None	8	Right	Nominal
2	pplata	Numeric	5	0	Pocetnicka plata	None	0	8	Right	Scale
3	pol	Numeric	1	0	pol zaposlenog	{0, muskarc...	9	8	Right	Nominal
4	vreme	Numeric	2	0	Seniornost na poslu	None	0	8	Right	Scale
5	splata	Numeric	5	0	Sadasnja plata	None	0	8	Right	Scale
6	obrazovanje	Numeric	2	0	Nivo obrazovanja	None	0	8	Right	Scale
7	kat_posla	Numeric	1	0	Kategorija zaposlenog					Ordinal
8	rasa	Numeric	1	0	Rasna klasifikacija					Nominal
9	pol_rasa	Numeric	1	0	Polno-rasna podjela					Nominal
10	starost	Numeric	6	2	Starost radnika					Scale
11	r_staz	Numeric	6	2	radno iskustvo					Scale
12										
13										
14										
15										
16										
17										
18										
19										
20										
21										
22										
23										
24										
25										

Two-Independent-Samples Tests

Test Variable List:
 Seniornost na poslu...

Grouping Variable:
 pol(1 2)

Define Groups...

Test Type

Mann-Whitney U Kolmogorov-Smirnov Z
 Moses extreme reactions Wald-Wolfowitz runs

OK Paste Reset Cancel Help

Two-Independent-Samples Tests

Test Variable List: Seniomost na poslu...

Identifikacioni bro...
Pocetnicka plata [...]
Sadasnja pl...
Nivo obrazov...
Kategorija za...
Rasna klasif...
Polno-rasna...
Starost radn...
radno iskust...

Exact...
Options...

Two Independent Samp...

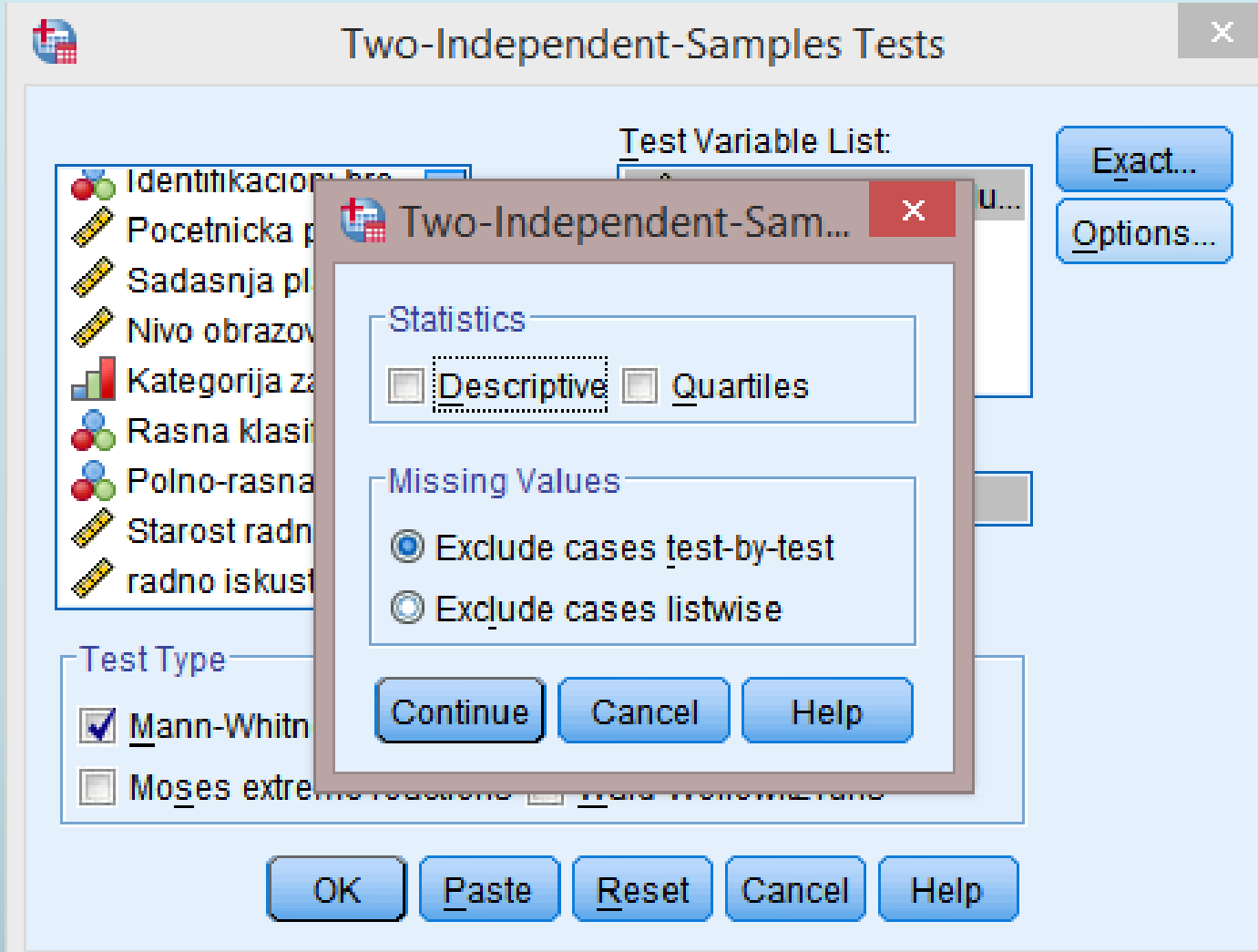
Group 1: 0
Group 2: 1

Continue Cancel Help

Test Type

Mann-Whitney U Kolmogorov-Smirnov Z
 Moses extreme reactions Wald-Wolfowitz runs

OK Paste Reset Cancel Help





- Output
 - Log
 - NPAr Tests
 - Title
 - Notes
 - Active Dataset
 - Mann-Whitney Test
 - Title
 - Ranks
 - Log
 - NPAr Tests
 - Title
 - Notes
 - Active Dataset
 - Mann-Whitney Test
 - Title
 - Ranks
 - Test Statistics

Mann-Whitney Test

Ranks

	pol zaposlenog	N	Mean Rank	Sum of Ranks
Seniornost na poslu	muskarci	258	245,69	63388,50
	zene	216	227,72	49186,50
	Total	474		

Test Statistics^a

	Seniornost na poslu
Mann-Whitney U	25750,500
Wilcoxon W	49186,500
Z	-1,424
Asymp. Sig. (2-tailed)	,155

a. Grouping Variable: pol zaposlenog

Test

- ▶ $H_0(F_X(x) = F_Y(x))$
 - ▶ $H_1(F_X(x) \neq F_Y(x))$

 - ▶ Test - statistika $Z = -1.424$
 - ▶ P - vrednost testa $p = 0.155$

 - ▶ Zaključak: Prihvatamo nultu hipotezu.
-

Wald - Wolfowitz Run test

- ▶ Ili: Test koraka
- ▶ $H_0: F_1(x)=F_2(x)$ Ne postoji značajna razlika između funkcija raspodele. Drugim rečima, populacije iz kojih su uzorci izdvojeni su identičnih raspodela.
- ▶ $H_1: F_1(x)\neq F_2(x)$ Postoji značajna razlika između ovih raspodela.
- ▶ Vrednosti su rangirane u rastućem poretku, svaka vrednost je kodirana sa 1 ili 2 i ukupan broj runs-a (koraka) se sumira i koristi kao test statistika - R. Runs oznacava broj promena u grupi koju posmatramo, tačnije pod jednim korakom podrazumevamo niz elemenata iste kategorije proizvoljne dužine.
- ▶ Male vrednosti govore o tome da se populacije razlikuju, a velike ukazuju na to da se radi o identičnim populacijama u smislu raspodele.



Wald - Wolfowitz Run test

- ▶ Raspodelu test - statistike R i u ovom slučaju aproksimiramo normalnom raspodelom. I to na sledeći način:

- ▶ Ako je $m+n \geq 50$

- ▶ $Z = \frac{R - \mu_R}{\sigma_R}$, gde su $\mu_R = \frac{2mn}{m+n} + 1$ i $\sigma_R = \sqrt{\frac{2mn(2mn - m - n)}{(m+n)^2(m+n+1)}}$

- ▶ Inače $Z_c = \begin{cases} (R - \mu_R + 0.5) / \sigma_R, & |R - \mu_R| \geq 0.5 \\ 0, & |R - \mu_R| < 0.5 \end{cases}$

- ▶ P-vrednost se računa kao $p_1 = P\{Z \leq z\} = \Phi(z)$.

- ▶ Ili, ako je $m+n \leq 30$: $p_1 = P\{r \leq R\} = \sum_{r=2}^R f_R(r)$ gde je $f_R(r) = \frac{2 \binom{m-1}{\frac{r-1}{2}} \binom{n-1}{\frac{r-1}{2}}}{\binom{m+n}{m}}$

- ▶ U oba slučaja važi da ako je $p_1 < \alpha$ odbacujemo nultu hipotezu.

Primer

- ▶ Na osnovu baze sa početka ispitujemo:
 - ▶ Da li plata zavisi od pola zaposlenih?

 - ▶ Zavisna promenljiva: Sadašnja plata
 - ▶ Nezavisna: Pol zaposlenog
-



Two-Independent-Samples Tests



- Identifikacioni bro...
- Pocetnicka plata [...]
- Seniornost na po...
- Nivo obrazovanja ...
- Kategorija zaposl...
- Rasna klasifikaci...
- Polno-rasna pod...
- Starost radnika [s...
- radno iskustvo [r...



Test Variable List:

Sadasnja plata [spl...

Exact...

Options...

Grouping Variable:

pol(0 1)

Define Groups...

Test Type

- Mann-Whitney U
- Kolmogorov-Smirnov Z
- Moses extreme reactions
- Wald-Wolfowitz runs

OK

Paste

Reset

Cancel

Help

*Output2 [Document2] - IBM SPSS Statistics Viewer

File Edit View Data Transform Insert Format Analyze Direct Marketing Graphs Utilities Add-ons Window Help

Output

- Log
- NPar Tests
 - Title
 - Notes
 - Active Dataset
 - Wald-Wolfowitz Test
 - Title
 - Frequencies
 - Test Statistics

Wald-Wolfowitz Test

Frequencies

		pol zaposlenog	N
Sadasnja plata	muskarci		258
	zene		216
	Total		474

Test Statistics^{a,b}

		Number of Runs	Z	Asymp. Sig. (1-tailed)
Sadasnja plata	Minimum Possible	100 ^c	-12,619	,000
	Maximum Possible	196 ^c	-3,720	,000

a. Wald-Wolfowitz Test
 b. Grouping Variable: pol zaposlenog
 c. There are 53 inter-group ties involving 216 cases.

IBM SPSS Statistics Processor is ready

Test

- ▶ $H_0(F_X(x) = F_Y(x))$
 - ▶ $H_1(F_X(x) \neq F_Y(x))$

 - ▶ Test - statistika $-12.619 \leq Z \leq -3.720$
 - ▶ P - vrednost testa $p = 0$

 - ▶ Zaključak: Odbacujemo nultu hipotezu.
-

Kolmogorov - Smirnov test

- ▶ Elementi oba uzorka se sortiraju u rastućem poretku.
- ▶ Računamo empirijske funkcije raspodele oba uzorka.
- ▶ Za svako x_j računamo razliku između uzoraka :

$$m \geq n \quad d_j = \bar{F}_1(x_j) - \bar{F}_2(x_j)$$

$$n \geq m \quad d_j = \bar{F}_2(x_j) - \bar{F}_1(x_j)$$

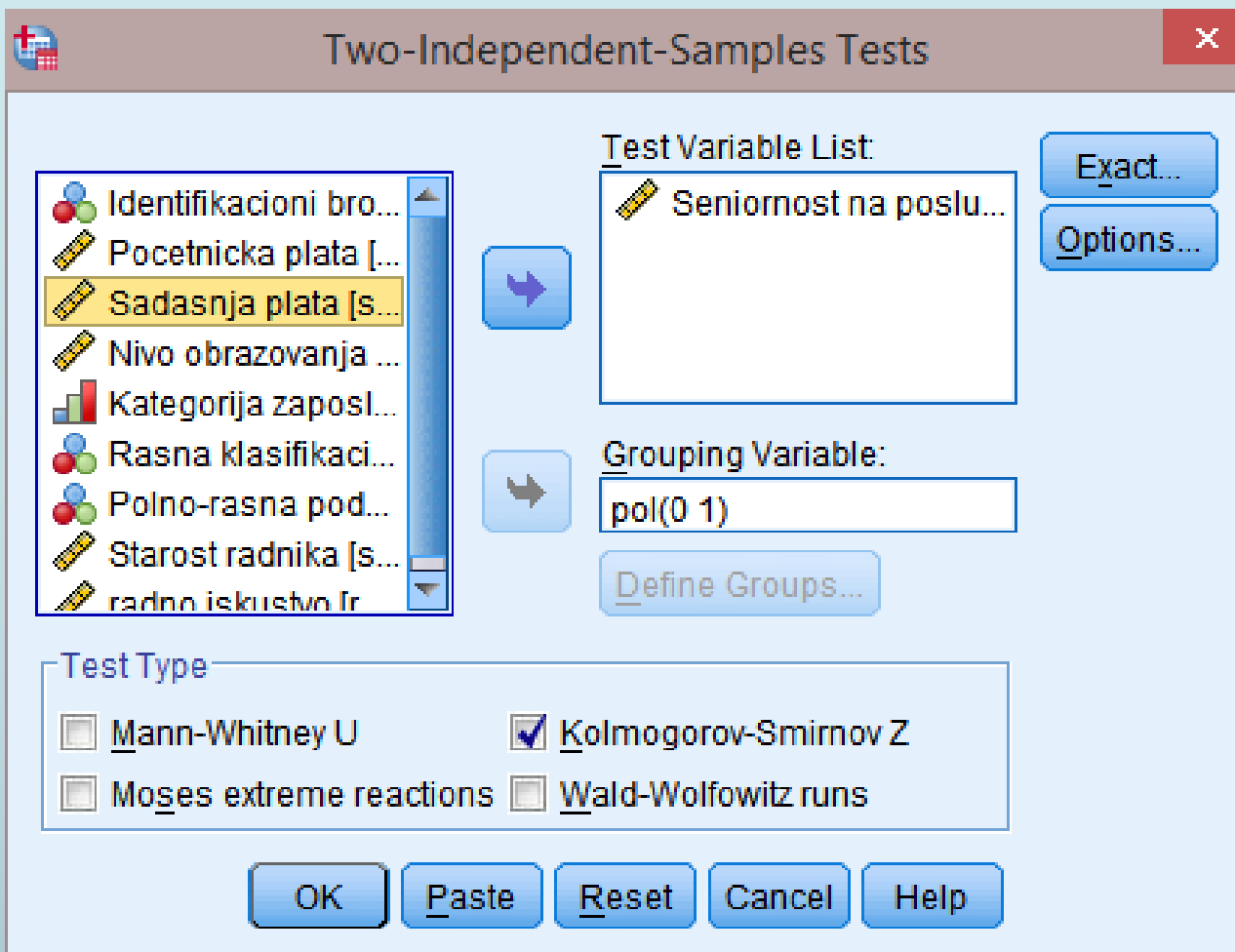
- ▶ Test statistika : $Z = \max_j |d_j|$.
-

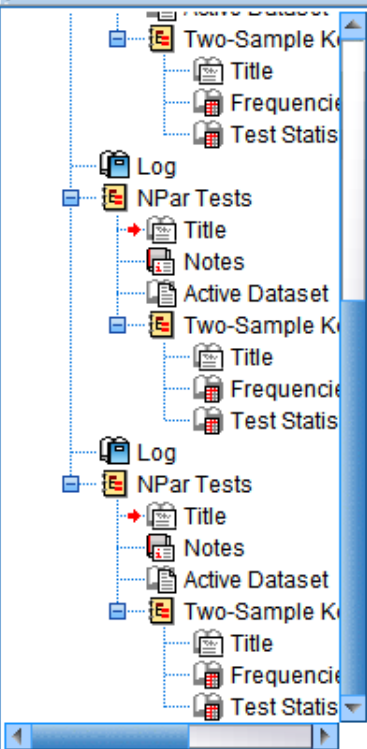
Primer

- ▶ Koristimo isti primer kao kod Mann-Whitney U testa.
 - ▶ Dakle, pitanje je:
 - ▶ Da li dužina rada zavisi od pola zaposlenih?

 - ▶ Zavisna promenljiva: Seniornost na poslu
 - ▶ Nezavisna: Pol zaposlenog
-







Two-Sample Kolmogorov-Smirnov Test

Frequencies

	pol zaposlenog	N
Seniornost na poslu	muskarci	258
	zene	216
	Total	474

Test Statistics^a

		Seniornost na poslu
Most Extreme Differences	Absolute	,094
	Positive	,055
	Negative	-,094
Kolmogorov-Smirnov Z		1,018
Asymp. Sig. (2-tailed)		,251

a. Grouping Variable: pol zaposlenog

Test

- ▶ $H_0(F_X(x) = F_Y(x))$
 - ▶ $H_1(F_X(x) \neq F_Y(x))$

 - ▶ Test - statistika $Z = 1.018$
 - ▶ P - vrednost testa $p = 0.251$

 - ▶ Zaključak: Prihvatamo nultu hipotezu.
-

Moses Extreme Reactions test

- ▶ Moses - ovim testom se testiraju dva nezavisna uzorka sa neprekidnom raspodelom.
- ▶ Nulta hipoteza: Verovatnoće da svaka od populacija sadrži ekstremne vrednosti su jednake.
- ▶ Alternativna hipoteza: Veća je verovatnoća da se ekstremne vrednosti nađu u populaciji iz koje je izdvojen uzorak većeg obima.
- ▶ SPAN (raspon) :
- ▶ Vrednosti iz oba uzorka se spoje, sortiraju i rangiraju. Vrednost najmanjeg ranga određuje kontrolni uzorak.
- ▶ $SPAN = \text{najveći rang} - \text{najmanji rang (kontrolnog uzorka)} + 1.$



Moses Extreme Reactions test

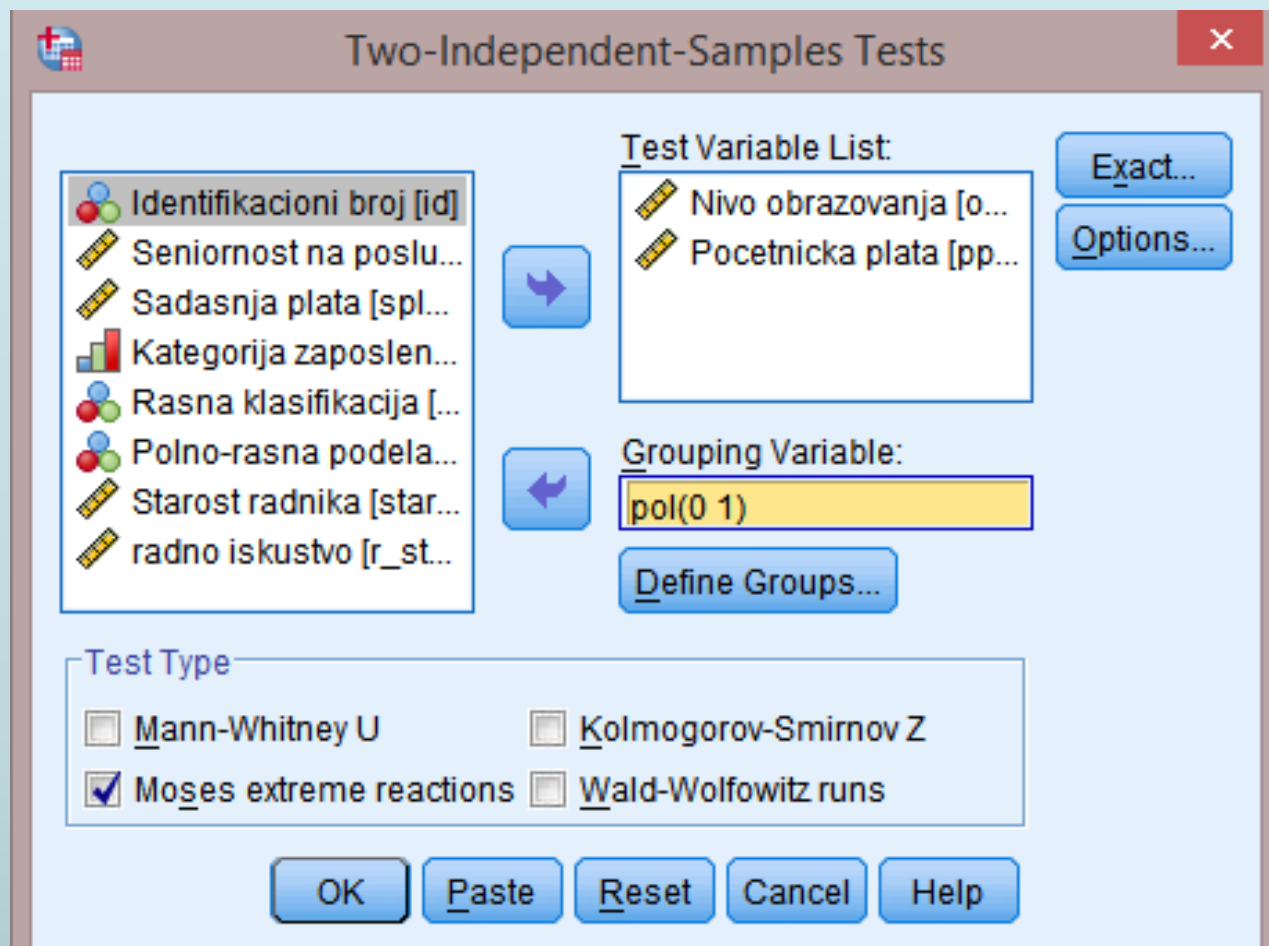
- ▶ Ako SPAN nije ceo broj, zaokružuje se na najbliži ceo broj.
- ▶ Neka m i n predstavljaju ukupan broj članova kontrolnog i eksperimentalnog uzorka, uključujući frekvencije, i $g = SPAN - m + 2h$, h je ceo broj jednak $0.05m$ ili 1 u zavisnosti od toga koji je broj veći.
- ▶ Onda je verovatnoca raspona s jednaka:

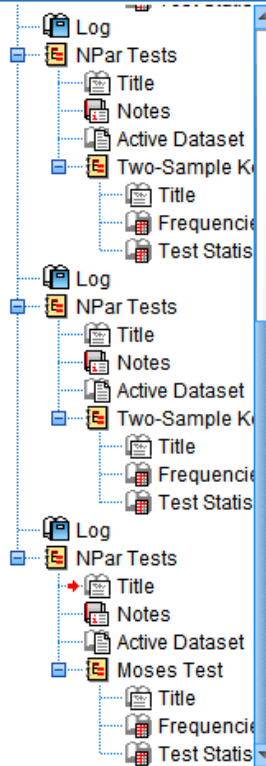
$$p_1 = P\{s \leq SPAN\} = \frac{\sum_{i=0}^g \left[\binom{i+m-2h-2}{i} \binom{n+2h+1-i}{n-i} \right]}{\binom{m+n}{m}}$$

- ▶ Ukoliko je $p_1 \leq \alpha$, odbacujemo nultu hipotezu.
-

Primer

- ▶ U ovom primeru ćemo testirati dve zavisne promenljive: Nivo obrazovanja i Početnička plata.
- ▶ Nezavisna promenljiva je, kao i do sada, Pol zaposlenih.





Moses Test

Frequencies

	pol zaposlenog	N
Nivo obrazovanja	muskarci (Control)	258
	zene (Experimental)	216
	Total	474
Pocetnicka plata	muskarci (Control)	258
	zene (Experimental)	216
	Total	474

Test Statistics^{a,b}

		Nivo obrazovanja	Pocetnicka plata
Observed Control Group		448	473
Span	Sig. (1-tailed)	,000	,704
Trimmed Control Group		432	303
Span	Sig. (1-tailed)	,749	,000
Outliers Trimmed from each End		12	12

- a. Moses Test
- b. Grouping Variable: pol zaposlenog

Test

- ▶ Nivo obrazovanja:
 - ▶ Test - statistika SPAN = 448
 - ▶ P - vrednost testa $p = 0$
 - ▶ Zaključak: Odbacujemo nultu hipotezu.

 - ▶ Početnička plata:
 - ▶ Test - statistika SPAN = 473
 - ▶ P - vrednost testa $p = 0.704$
 - ▶ Zaključak: Prihvatamo nultu hipotezu.
-

