

Országos Korányi Tbc és Pulmonológiai Intézet

Az SPSS statisztikai programcsomag használatának tapasztalatai

Angyal István, Hofhauser Béla, Kiss Péter

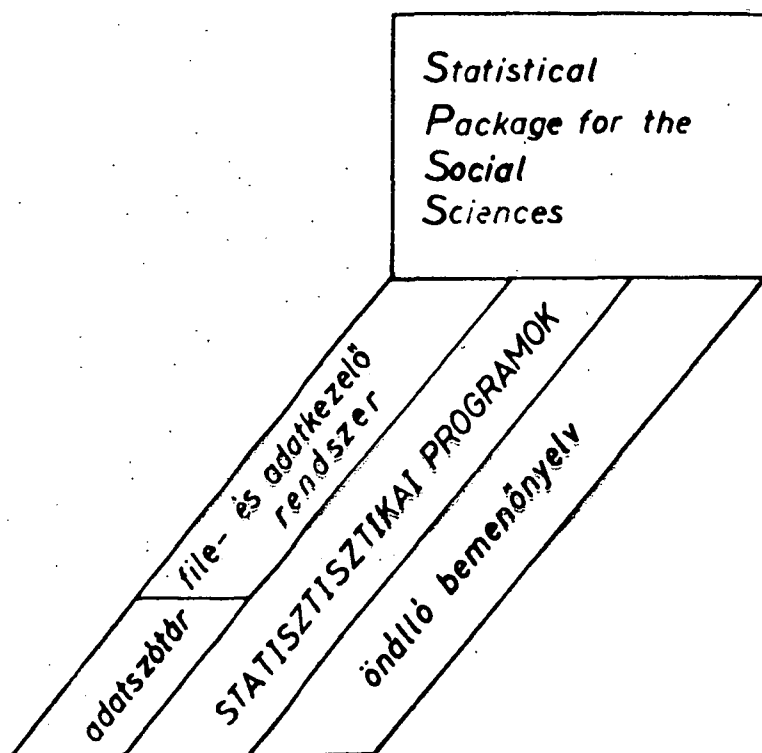
Az Intézetben folyó számítástechnikai tevékenység legfőbb célja az Országos Intézeti Modell kidolgozása. A Modell és a létrehozandó rendszerek feltételezik egy olyan szoftver eszköz mindennapos használatát, mely alkalmas nagytömegű adatok megbízható statisztikai leírására és elemzésére.

Egy ilyen szoftvér eszköz saját erőből való létrehozása a még ma is fennálló létszámhiány miatt kizártnak tekinthető. A hazai programcsomagok között nem találtunk olyan szavatolt és karbantartott terméket, amely statisztikai elemzéseket és egyben tág adat- és file-kezelési lehetőségeket is megenged. A külföldi programcsomagok közül az SSP /Scientific Subroutine Package-et/ és a BMD-t adatkezelési illetve esetszám korlátozási nehézségeik miatt vetettük el.

A külföldről beszerzett ajánlatok közül az SPSS Inc. ajánlata tűnt a legkedvezőbbnek. 30.000 Ft-nak megfelelő devizáért hozzájutottunk az SPSS programcsomag örökös bér-

letéhez. E programcsomagot a világon kb. 1500 helyen installálták már és igen széles felhasználói körben alkalmazzák. Az alapváltozatának kidolgozása kb. 100 emberévi programozói munkát jelentett. Az újabb változatok mind több statisztikai eljárást tartalmaznak, melyek minőségét szavatolják és a teljes programcsomag karbantartását vállalják.

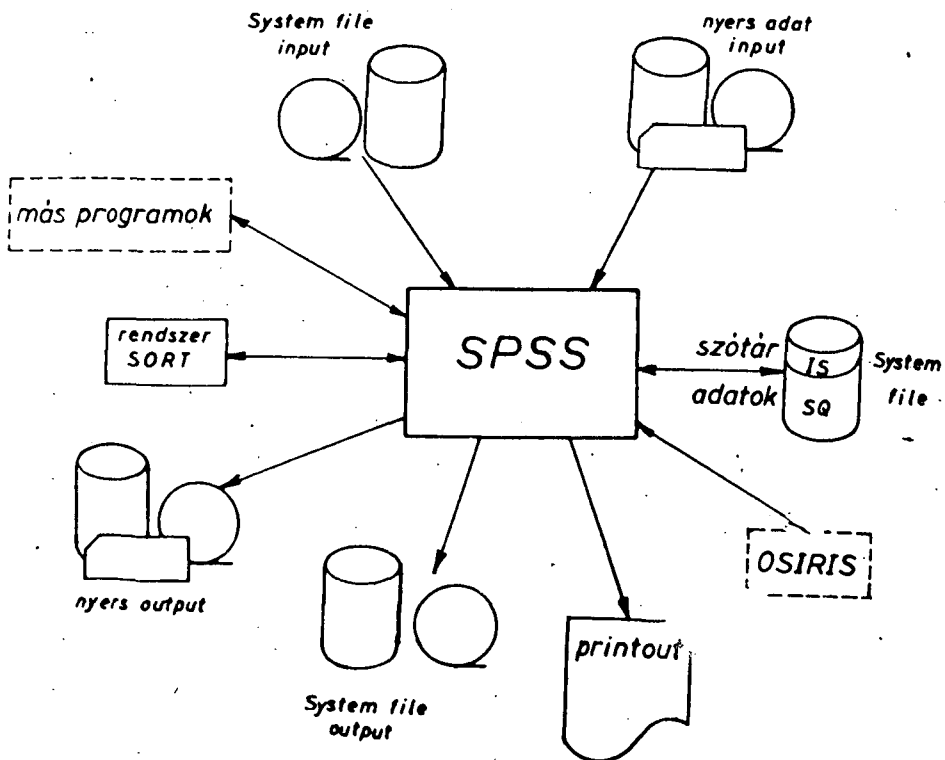
A Statistical Package for the Social Sciences programcsomag /1. ábra/ önálló bemenőnyelven keresztül fix szerkezetű adatállományok statisztikai leírását és elemzését végzi el és egyben adatszótára segítségével bonyolult file és adatkezelési műveleteket is lehetővé tesz.



1. ábra

Az önálló bemenőnyelv minden lehetőséget számítva 80 kulcsszóból áll, szintaxisa és logikája egyszerű, számítástechnikában nem járatos emberek számára is viszonylag könnyen elsajátítható. Részletes hibaüzenetei a felhasználót gyorsan rávezeti a tévesztésekre.

A file és adatkezelés a programrendszer saját adatszótárán és adatállományán /2. ábra/ az u.n. system file-on keresztül történik. Az adatszótár indexszekvenciális szervezésű és a vizsgált adatállomány adatainak fizikai, logikai és egyéb releváns információit tartalmazza. A system file adatrésze az input adatállománynak a rendszer szempontjából standardizált szekvenciális leképezése. A system file magasfokú integráltsága a változók kezelését egyszerűvé, áttekinthetővé teszi és ez többek között a tévesztési és hibalehetőségeket jelentősen csökkenti.



2. ábra

A programrendszer fontos tulajdonsága az adatokkal és más szoftverekkel szembeni nyitottság. Tetszőleges formátumú FORTRAN, PL/1 vagy COBOL programok által létrehozott output adatállományokat tud fogadni, illetve ezek számára input-ot tud generálni. A meglévő statisztikai subrutinokhoz továbbiak is beilleszthetők és az OSIRIS programcsomaghoz kész interface-sel rendelkezik. Az adatállományok rendezése során a rendszer-SORT programot aktivizálja.

Az SPSS szempontjából nyers input adatállományból a VARIABLE LIST és INPUT FORMAT kártyák paraméterei szerint az első statisztikai alprogram hívásakor a rendszer létrehozza az ideiglenes /temporary/ system file-t, melyet a további tetszőleges számban meghívható alprogramok mindegyike végigolvas. A létrehozott temporary file-t a SAVE FILE utasítás hatására a rendszer teljes egészében kimentti és további futtatások során a GET FILE utasítással behívható, nincs szükség további adattleírásra.

A nyers output-ra tetszőleges illetve szelektált adatok vagy a statisztikai programok közbeni eredményei mint pl. korrelációs mátrix, residuumok, decriptiv információk stb. vihetők ki. Ezeket a rendszer standard módon inputállományként is tudja kezelni.

A beolvasott system file-hoz lehetőség van nyers in-

putról esetek, subfile-ok vagy változók hozzáadására. A system file kimentésekor pedig mód van feleslegessé vált változók vagy esetek törlésére. A file-kezelési műveletek tág lehetőséget biztosítanak a feldolgozási szempontoknak megfelelő esetek kialakításához. Erre azért is van szükség, mert az egészségügyi szervezés feldolgozásai gyakran merőben más eset-definiálást igényelnek. Mások a vizsgálandó eset kritériumai, ha egy egészségügyi egység terhelését vizsgáljuk és mások, ha egy bizonyos betegség ápolási, gondozási szükségletét elemezzük.

A tág file-kezelési lehetőségeket sokrétű adatkezelési műveletek egészítik ki. A COMPUTE utasítással meglévő vagy új változóknak értéket adhatunk FORTRAN-szerű értékadó utasítással /3./a. ábra/. A BROCA nevű változó ezen utasítással az u.n. Broca-indexet fogja tartalmazni, melyet a testsúlyból és a testmagasságból lehet kiszámítani. Az IF feltételes értékadó utasítás tetszőleges logikai feltétel teljesülése esetén változtatja meg a hivatkozott változó értékét /3./b. ábra/. A systolés VERNYS és diastolés VERNYD vérnyomás értékeitől függően a vérnyomásmérés minőségét kifejező változó VERNY 1, 2 vagy 3-as értéket vesz fel, mely kifejezi a hypo-, a normo- és a hypertóniát.

A RECODE-utasítás a változók értékészletének átkódolását teszi lehetővé /3./c. ábra/. Az életkort kifeje-

zõ KOR változó értékeit korcsoportnak megfelelően átkódoljuk és ezek után a KOR változó 1-es értéke a tizenéveseket, a 2-es a fiatalokat stb. fogja jelenteni. A COUNT utasítás végigszámolja, hogy egy esetnél hányszor fordulnak elő bizonyos változók megadott értékei /3./d. ábra/. Pl. egy személy betegség-kódjait a BNO1, BNO2, BNO3 változók tartalmazzák. Amennyiben az egyik BNO-kód 250, azaz a diabetes kódja, akkor a DIABET változó értéke 1 lesz.

A programrendszer tetszőleges szelektálást, az esetek súlyozását és mintavételezését is megengedi.

A SELECT IF utasítása alapján csak a megadott logikai feltételnek eleget tevő rekordok vesznek részt a feldolgozásban /3./e. ábra/. Tehát csak a 110-nél kisebb és 90-nél nagyobb Broca-index-el rendelkező u.n. atletikus természetű személyeken végzi el a megadott feldolgozásokat. A WEIGHT utasításban megadott változó értékének megfelelően súlyozza az egyes eseteket. /3./f. ábra/ Példánkban a SÜLYOZ változót úgy határoztuk meg, hogy a diabetes diagnózissal nem rendelkező személyek csak fele súllyal vegyenek részt a feldolgozásokban. Mindezek az adattranszformációk és szelektálási lehetőségek megadhatók permanensen, tehát egy egész futtatásra, vagy pedig időlegesen tehát egy statisztikai program végrehajtásának idejére. Nagy és viszonylag homogén adatállományok elemzését segíti elő a SAMPLE utasítás. Példánk /3./g. ábra/ 10 %-os

mintavételezést ír elő, mely egy beépített és újra indítható véletlenszámgenerátorral történik.

```
a./ COMPUTE      BROCA=TESTS/(TESTM-100)*100

b./ COMPUTE      VERNY=2
   IF            (VERNYS GT 160 OR VERNYD GT 100) VERNY=3
   IF            (VERNYS LT 100 OR VERNYD LT 70) VERNY=1

c./ RECODE       KOR (10 THRU 19=1) (20 THRU 39=2)
                  (40 THRU 59=3) (60 THRU HIGHEST=4)

d./ COMPUTE      DIABET=0
   COUNT         DIABET=BNO1,BNO2,BNO3 (250)

e./ SELECT IF    (BROCA LT 110 AND BROCA GT 90)

f./ COMPUTE      SULYOZ=0.5+0.5*DIABET
   WEIGHT        SULYOZ

g./ SAMPLE       0.1
```

### 3. ábra

Az SPSS 6 leíró és 14 elemző statisztikai eljárást tartalmaz /4. ábra/. Az egyes programok jellemzőinek ismertetése, algoritmusainak, módszereinek leírása, más programcsomagokkal való összehasonlítása túlnő e 15 perces előadás keretein. Az elemző eljárások közül jelenleg még hiányzik a clusteranalysis, a többváltozós varianciaanalysis és mindenek előtt az idősorok elemzésére alkalmas eljárások. Újabb kibocsájtások előrejelzések szerint a többváltozós varianciaanalysisist már biztosan tartalmazni fogják és a további programok beépítését is tervezik.

AGGREGATE	ANOVA
BREAKDOWN	CANCORR
CONDESCRIPTIV	DISCRIMINANT
CROSSTABS	FACTOR
FREQUENCIES	GUTTMAN SCALE
SCATTERGRAM	MULTI RESPONSE
	NONPAR CORR
	NPAR TESTS
	ONEWAY
	PARTIAL CORR
	REGRESSION
	RELIABILITY
	T-TEST

#### 4. ábra

A programcsomag lehetővé teszi a hiányzó értékek egy-  
séges ellenőrizhető kezelését az adat és file kezelés so-  
rán és valamennyi statisztikai programban. Minden változó-  
ra maximálisan három hiányzó érték /missing value/ adható  
meg. Így lehetőség van arra, hogy különbséget tegyünk pl.  
azok között, akik nem akartak, nem tudtak, illetve hely-  
telenül válaszoltak a feltett kérdésekre.



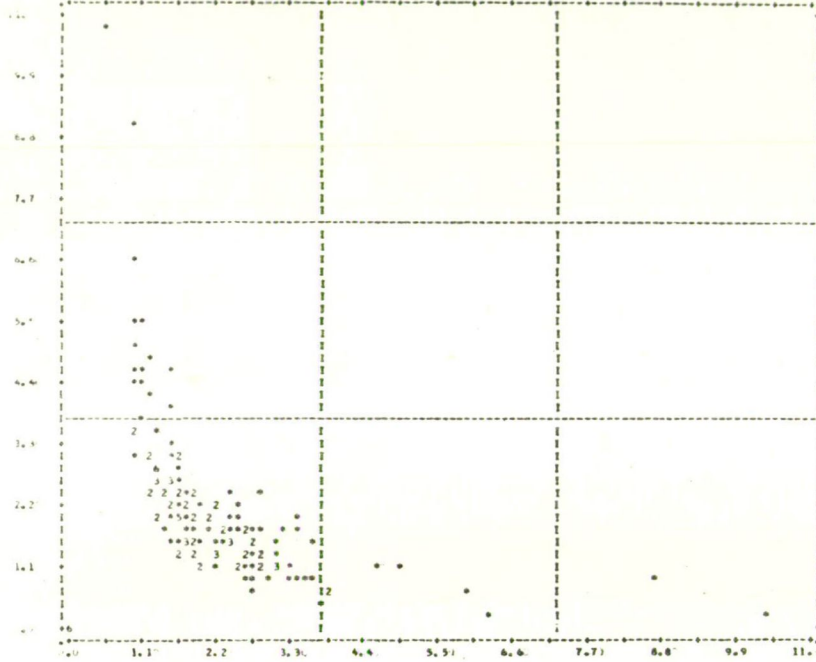
Az 5. ábrán légzésfunkciós értékek SCATTERGRAM-ját láthatjuk. Felette ugyanezen értékek VREZIST illetve SGAV transzformációját és átkódolását az A és B változóba. Az átkódolás érdekessége még, hogy B értékeinek átkódolása során a skálát megfordítottuk és így az A és B kereszttablóján /6. ábra/, az előző összefüggés tükörképe figyelhető meg /CORSSABS/. A kért statisztikák a kontingencia táblában számítható különböző együtthatókat szemléltetik.

Az SPSS installálása óta eltelt másfél év alatt több kisebb, nagyobb feldolgozást végeztünk /7. ábra/. A változók és az esetek száma kétszer logaritmikus koordináta-rendszerében minden egyes pont egy adatállományt illetve feldolgozást jelöl.

- GY : alkalmi kisebb volumenű klinikai-hatástani gyógyszer tesztek, melyek az Intézetünkben folytak és az SPSS-el értékeltünk ki.
- B76 és B78 : az 1976-os Budakeszi lakosságsszűrés és ugyanezen populáció légzési panaszokkal rendelkező részének 1978-ban újraszűrése.
- TOV : az 1976-os Budakeszi lakosságsszűrés légzésfunkciós továbbvizsgáltjairól készült adatállomány. Az 5. és 6. ábrán erről a populációról láthattunk egy kiragadott adatösszefüggési példát.

COMPUTE N=V-1(15)  
 A=SGRY  
 RECORD ALL THRU 2=1117 THRU 3=2113 THRU 4=3116 THRU 5=4115 THRU 6=51  
 16 THRU 7=6117 THRU 8=7118 THRU 9=8119 THRU 11=911 THRU 12=3  
 01 THRU 13=9111 THRU 14=8112 THRU 15=7113 THRU 16=6114 THRU 17=51  
 15 THRU 18=4116 THRU 19=3117 THRU 20=2118 THRU 21=1119 THRU 22=1

HODAKESZI LAKOSSAGSORRES LEGZESFUJNCIOS TOVARVIZSGALTAK 11/23/76 PAGE 6  
 ERZESFUJNCIOS MERSEK SCATTERGRAMJA  
 REGISTRATION DATE = 13/28/78 TOVARVIZSGALTAK  
 FILE TUVZ SPECIFICUS KONDUKTANCIA (ACROSS) VREEST ARHAJAST ELLENALLAS  
 SCATTERGRAM IF 1.55 1.65 2.75 3.85 4.95 6.05 7.15 8.25 9.35 1.45



HODAKESZI LAKOSSAGSORRES LEGZESFUJNCIOS TOVARVIZSGALTAK 11/23/76 PAGE 7  
 ERZESFUJNCIOS MERSEK SCATTERGRAMJA  
 STATISTICS  
 CORRELATION (R) = -1.47596 R SQUARED = 1.16481 SIGNIFICANCE = 0.00001  
 STD ERR OF EST = 1.35691 INTERCEPT (A) = 3.77457 SLOPE (B) = -1.45628  
 PLOTTED VALUES = 155 EXCLUDED VALUES = 0 MISSING VALUES = 0

\*\*\*\*\* IS PRINTED IF A COEFFICIENT CANNOT BE COMPUTED.

5. ábra

RUDAKESZI LARUSSAGSZAMES LEZISEPUNKCIOS TOVAROVIZSGALTAK 11/23/78 PAGE 9  
 ALCOODILT ERTAKPA KERESZETAN.41  
 FILE TUV2 (CREATION DATE = 03/28/78) TOVABBVIZSGALTAK

\*\*\*\*\* CROSSTABULATION OF \*\*\*\*\*  
 A AZ SWAY ERTAKINEK ATKODOLASA BY B A VREZIST ERTAKINEK ATKODOLASA  
 \*\*\*\*\* PAGE 1 OF 2

		B									ROW TOTAL
		10-11	8-9	6-7	5-6	4-5	3-4	2-3	1-2	0-1	
A	COUNT										TOTAL
	ROW PCT										
		C									
		0	1	3	4	5	6	7	8	9	
0-1	0	1	1	1	1	0	7	4	0	6	21
	1	4.8	4.8	4.8	4.8	0.0	33.3	19.2	0.0	28.6	13.5
	1	100.0	133.0	137.0	137.0	0.0	41.2	9.1	0.0	37.5	
	1	0.6	0.6	0.6	0.6	0.0	4.5	2.6	0.0	3.9	
1-2	1	0	0	0	0	2	17	34	22	0	68
	1	0.0	0.0	0.0	0.0	2.9	16.7	5.0	32.4	0.0	43.9
	1	0.0	0.0	0.0	0.0	100.0	58.8	77.3	50.6	0.0	
	1	0.0	0.0	0.0	0.0	14.3	6.5	21.9	14.2	0.0	
2-3	2	0	0	0	0	0	0	6	36	0	42
	1	0.0	0.0	0.0	0.0	0.0	0.0	7.5	14.3	85.7	27.1
	1	0.0	0.0	0.0	0.0	0.0	0.0	13.6	57.1	0.0	
	1	0.0	0.0	0.0	0.0	0.0	0.0	34.9	23.2	0.0	
3-4	3	0	0	0	0	0	0	0	8	3	11
	1	0.0	0.0	0.0	0.0	0.0	0.0	0.0	72.7	27.3	7.1
	1	0.0	0.0	0.0	0.0	0.0	0.0	0.0	11.1	18.8	
	1	0.0	0.0	0.0	0.0	0.0	0.0	0.0	5.2	1.9	
4-5	4	0	0	0	0	0	0	0	5	3	8
	1	0.0	0.0	0.0	0.0	0.0	0.0	0.0	62.5	37.5	5.2
	1	0.0	0.0	0.0	0.0	0.0	0.0	0.0	6.9	18.8	
	1	0.0	0.0	0.0	0.0	0.0	0.0	0.0	3.2	1.9	
5-6	5	0	0	0	0	0	0	0	0	1	2
	1	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	50.0	1.3
	1	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	1.4	6.3
	1	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.6	
6-7	6	0	0	0	0	0	0	0	0	0	1
	1	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	100.0
	1	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	6.3
	1	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.6
9-11	9	0	0	0	0	0	0	0	0	0	2
	1	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
	1	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	100.0
	1	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	12.5
	1	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	1.3
COLUMN TOTAL		1	1	1	1	2	17	44	72	16	155
		0.6	0.6	0.6	0.6	1.3	11.0	28.4	46.5	10.3	100.0

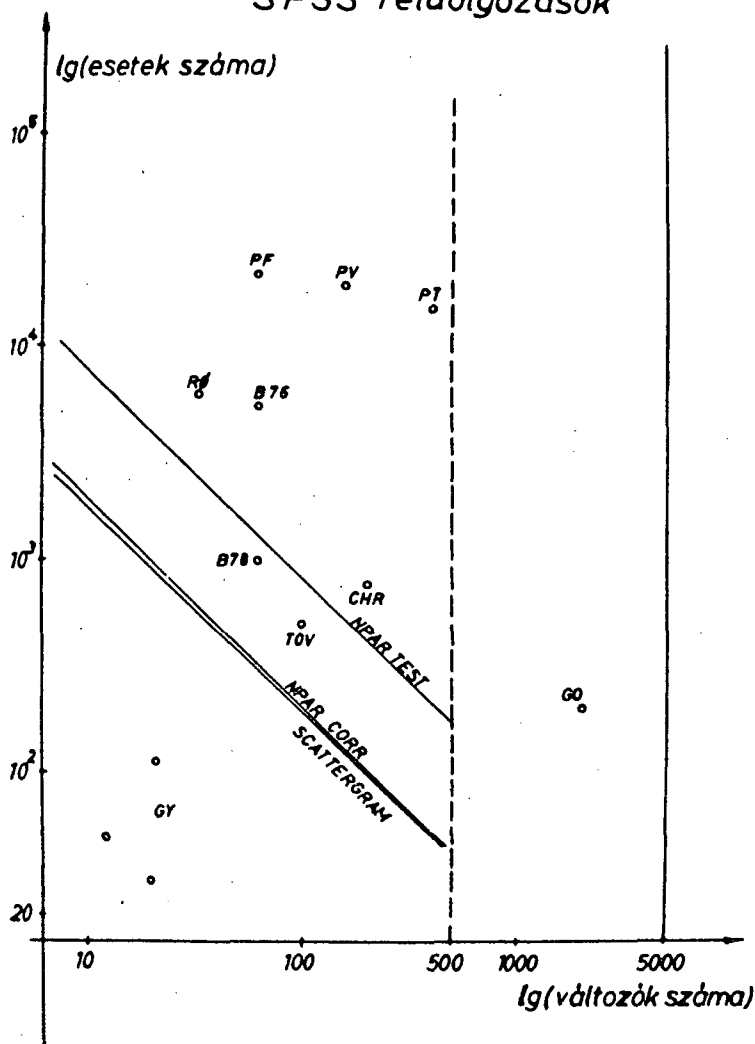
RAW CHI SQUARE = 150.77314 WITH 56 DEGREES OF FREEDOM. SIGNIFICANCE = 0.0000  
 CRAMER'S V = 0.37278  
 CONTINGENCY COEFFICIENT = 0.2722  
 LAMBDA (ASYMMETRIC) = 0.27586 WITH A DEPENDENT. = 0.26506 WITH B DEPENDENT.  
 LAMBDA (SYMMETRIC) = 0.27059  
 UNCERTAINTY COEFFICIENT (ASYMMETRIC) = 0.32622 WITH A DEPENDENT. = 0.34871 WITH B DEPENDENT.  
 UNCERTAINTY COEFFICIENT (SYMMETRIC) = 0.337.9  
 KENDALL'S TAU B = 0.52705. SIGNIFICANCE = 0.0000  
 KENDALL'S TAU C = 0.418.4. SIGNIFICANCE = 0.0000  
 GAMMA = 0.66833  
 SOMERS'S D (ASYMMETRIC) = 0.53743 WITH A DEPENDENT. = 0.51688 WITH B DEPENDENT.  
 SOMERS'S D (SYMMETRIC) = 0.52695  
 ETA = 0.54296 WITH A DEPENDENT. = 0.51998 WITH B DEPENDENT.  
 PEARSON'S R = 0.46650 SIGNIFICANCE = 0.0000

6. ábra

- CHR : a chronicus fertőző tbc-s betegek országos felmérésének adatállománya, mely lényegében még ma Magyarországon található súlyos tbc-s betegeket és egyben a legveszélyesebb fertőzősi potenciált jelentő személyeket tartalmazza.
- RØ : a negatív tüdőrontgen-előzménnyel rendelkező friss tbc-s betegekről készült országos felmérés során a negatív rtg előzmények és a friss tbc-s megbetegedések összefüggését vizsgáltuk.
- GO : a tüdőgondozó intézetek évi statisztikai jelentései /melyek manuális módszerrel készültek/ fontos támpontot jelentenek terveink, megvalósítandó rendszereink kialakításában.
- PF, PV és PT : a pécsi Komplex Szűrőállomás 1977. évi tevékenységének retrospektív vizsgálata, mely több mint 6 millió adatával az eddigi legnagyobb feldolgozásunk.

Az ábrán szemléltetni igyekeztünk az SPSS korlátait. 500 változóig standard kezelést biztosít, speciális kezeléssel 5000 változót lehet a system file-ba bevinni. Egy futtatás során viszont maximálisan 500 változóra lehet hivatkozni. A programok jelentős részénél nincsen esetszám-korlátozás. A hivatkozott változók számától függően csak meghatározott számú eset vehet részt az NPAR TEST, az NPAR CORR és a SCATTERGRAM programokban. Amennyiben több

### SPSS feldolgozások



7. ábra

esetünk van mint a megengedett, akkor az SPSS véletlenszerű mintát vesz a megengedhető maximális esetszámnak megfelelően az összes esetből.

Az SPSS adat és file-kezelésre igen alkalmas, a változók paramétereit tartalmazó adatszótár létrehozása, kezelése, karbantartása egyszerű. Az egyes statisztikai subru-

tinok gyorsaságának, effektivitásának megítélésére még nem teljesen kialakult a véleményünk. Más hasonló szoftvereket használó csoportokkal való megbeszéléseink, esetleg összehasonlító futtatások elemzése során erről realisabb képet kaphatunk. Minden esetre elmondható, hogy kis esetszámoknál /néhány 100 eset/ minden szempontból jónak találjuk a gyorsasági tényezőket, olyannyira, hogy az SPSS akár ilyen adatállományok validálására is alkalmas. Közepes feldolgozásoknál /néhány ezer esetenél/ a file és adatkezelés még kielégítő, de a statisztikai programoknál már körültekintőnek kell lenni. Ebben a tartományban az SPSS-t már nem lehet validálásra ajánlani. Nagy feldolgozásoknál /több tíz- esetleg százezer esetenél/ a végrehajtási idő optimalizálása érdekében körültekintően kell eljárni, pontosan meg kell tervezni a feldolgozásokat és az SPSS kezelésében gyakorlattal kell rendelkezni. Nagy esetszámok statisztikai elemzésénél pedig mindenképpen célszerű a rendszer mintavételezési lehetőségeit kihasználni.

A program közel teljes dokumentálása megtalálható az SPSS Manual-ben /1975 McGraw-Hill, Inc. New York/. E kézikönyv szinte minden szempontból jó, talán egyedüli hiányosságának az adható meg, hogy a matematikai statisztikai eljárások módszereinek, feltételezéseinek, kikötéseinek leírása nem mindig teljes és azok csak a megadott irodalom gondos áttanulmányozása során deríthetők ki. Ez bi-

zonyos mértékig akadályozza az adequát statisztikai elemzések gyors alkalmazását.

Összefoglalva elmondhatjuk, hogy az SPSS széles alkalmazási területe és magasfokú integráltsága alapján a számítástechnika mai fejlettségi szintjén a világszínvonalú, professzionista alkalmazási programcsomagok közé tartozik. A statisztikai programcsomagok közül az SPSS-nek elsőnek sikerült teljesíteni azokat a nagyon nehezen elérhető kritériumokat, melyek alapján 1977-ben megkapta az alkalmazási szoftverek részére adható egyik legnagyobb elismerést a Datapro Software Honor Roll-t.

Közép vagy nagy számítógéppel rendelkező intézményeknek csak javasolni tudjuk e programcsomag beszerzését. Ezen gépek üzemeltetőinek és más érdeklődőknek is szívesen adunk tájékoztatást további részletekről.