

*Semmelweis Egyetem Budapest, Bőr- és Nemikórtani Klinika
(igazgató: Horváth Attila dr. egyetemi tanár) közleménye*

Matematikai-statisztikai módszerek a tudományos és alkalmazott kutatásban és interpretációjuk kérdései

Mathematical-statistical methods in scientific and applied research and their interpretation

KOVÁCS JÁNOS DR.

ÖSSZEFOGLALÁS

*A tudományos kutatásban alapvető igény, hogy a vizsgá-
lás során kapott eredményeket a matematikai és statisztikai
módszertan törvényszerűségei szerint értékeljék,
hogy eredményeik reprodukálhatók legyenek, és a valósá-
got tényszerűen írják le. Az ilyen tartalmú közlemények
olvasói számára is elengedhetetlen, hogy a kapott ered-
mény alapján az alapvető összefüggések a sokszor bonyo-
lult adathalmazból értelmezhetőek legyenek. A cikkben a
szerző az alapfogalmak tisztázásával, ezek kritikájával és
néhány gyakorlati példával segíti a kérdéskörben való
eligazodást.*

Kulcsszavak:
**matematikai-statisztikai módszerek -
tudományos kutatás - értelmezés**

SUMMARY

*It is strongly recommended to evaluate results in the
scientific research upon principles of mathematical-
statistical methodology, because they must be repeatable
and realistic. Readers of scientific articles also should
understand basic relations obtained from the complex
data. With explanation and critic of basic elements to this
issue and with practical examples the author helps the
orientation in the theme.*

Key words:
**mathematical-statistical methods -
scientific research - interpretation**

A statisztika a valóság tömör jellemzésére szolgáló tudományos módszertan, illetve gyakorlati tevékenység, mindig a tények valamilyen összességét kívánja jellemezni. Kiindulási bázisát adatok vagy mutatószámok jelentik, melyek tulajdonképpen az egyedek egy bizonyos körét összességükben jellemző számszerű információk. A vizsgálat tárgyát képező egységek halmazát statisztikai sokaságnak nevezzük. A sokaság tagjainak besorolását ismérvek (területi, időbeli, minőségi, mennyiségi) segítik, melyek alapján a sokaság egymást át nem fedő részekre bontható.

Az adatszerzési módszerek a következők:

- adatgyűjtések, adatfelvételek (teljes körű, részleges)
- tervezett kísérletek

A rögzítés elsődleges eszköze ezek alapján kérdőív, kutatási napló, egyéb dokumentáció lehet.

Az egyedi adatfelvétel keresztmetszeti, míg az ismétlődő adatfelvétel longitudinális adatokat szolgáltat, utóbbi így a követéses, összehasonlító vizsgálatok alapvető eszköze, és induktív elemzés során általánosításra is lehetőséget ad.

- Exploratív (feltáró jellegű) adathasznosítás esetén a cél az adatokban megnyilvánuló szabályszerűségek feltárása, míg a
- konfirmatív (igazoló) jellegű elemzés a sokaságra előre megfogalmazott hipotézisek igazának tisztázására irányul (pl. adott gyógyszer kiváltja-e a kívánt hatást vagy sem).

Példaként elég csak a Viagra nevű sikergyógyszerre gondolnunk, ahol a konfirmatív elemzés nem igazolta a gyógyszer szívre való jó hatását, és az exploratív analízis tárta fel a mellékhatásként jelentkező, később főhatásként hasznosított áldásos effektust.

A sokaságok egy ismérv szerinti vizsgálatának alapját a gyakorisági eloszlások ismerete képezi, mivel ennek hiányában következtetések levonására alkalmas bizonyos módszertani eljárások nem alkalmazhatók.

A gyakorisági eloszlást helyzete, szóródása és alakja jellemzi. Ezen jellemzők ismerete esetén az alapadatok hiányában is képet kaphatunk az eloszlásról, illetve az általa jellemzett sokaságról.

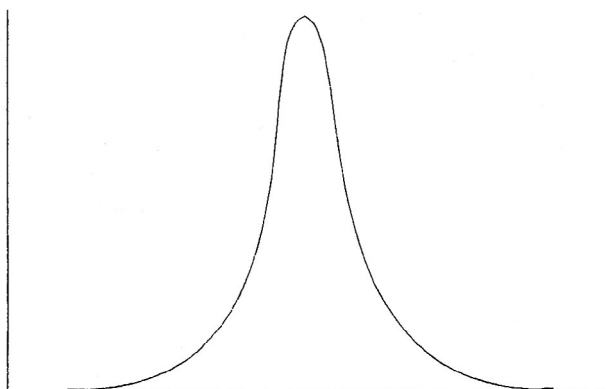
- Az eloszlás **helyzete** a tipikusnak mondható ismérv-értékek előfordulási helyét jelenti az x-tengelyen.

➤ Az eloszlás **szóródása** az ismértékek egymás közötti különbözőségéről, változékonyságáról nyújt információt.

➤ A gyakorisági eloszlás **alakja** poligonjának egy vele azonos elhelyezkedésű és szóródású normális eloszlás gyakorisági görbéjéhez képesti jellegzetességeit jelenti.

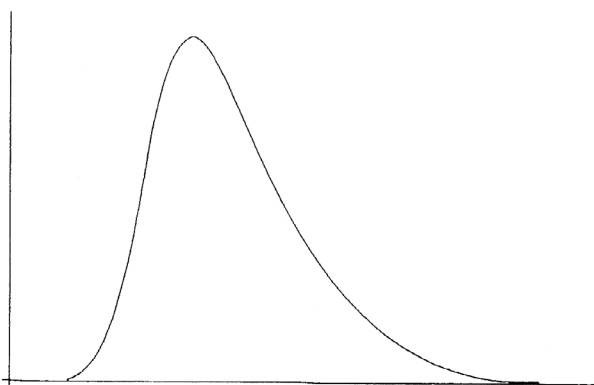
Példaként három, típusosnak mondható eloszlást mutatnak a következő ábrák.

A leggyakrabban előforduló, és legjobban értelmezett a normális eloszlás (ilyen például a munkateljesítmények eloszlása) (1. ábra).



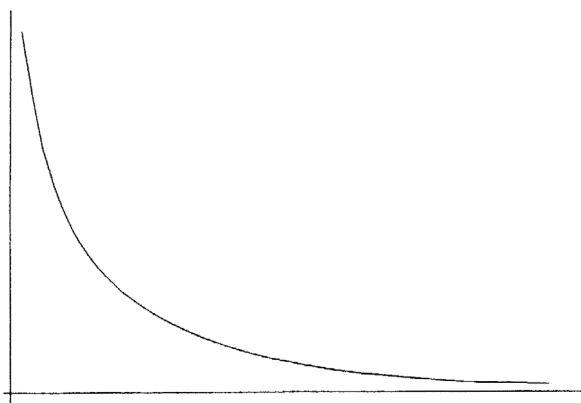
1. ábra
Normális eloszlás (1.)

A következő, tipikus gyakorisági eloszlás a lognormális eloszlás, ilyen például a legtöbb országban a háztartások jövedelmének eloszlása (2. ábra).



2. ábra
Lognormális eloszlás (1.)

Az exponenciális eloszlást kezdetben nagyobb, majd kisebb csökkenés jellemzi, ilyen például a gépalkatrészek élettartamának eloszlása (3. ábra).



3. ábra
Exponenciális eloszlás (1.)

A gyakorisági eloszlás helyzet mutatóit középtértéknek is nevezzük, hiszen az egész gyakorisági sort egyetlen, az ismértékekkel azonos mértékegységű, közepes helyzetű, tipikus, lehetőleg egyértelműen meghatározható, könnyen értelmezhető középtértékkel jellemzik.

Helyzeti mutatók: medián, módusz.

Számított középtértékek: számtani, harmonikus, mértani és a négyzetes átlag.

A **medián** a középső kvantilis, a szó szoros értelmében középtérték, hiszen azt az értéket jelenti, melynél pontosan ugyanannyi kisebb, mint nagyobb érték fordul elő (1. táblázat). Mindig egyértelműen meghatározható (mindig van középső ismérték), ordinális (sorrendi) mérési szint esetén is használható. Közvetlenül nem függ az összes ismértéktől, így a szélsőséges értéktől sem, induktív statisztikai célokra viszont alig alkalmas.

A **módusz** a leggyakrabban előforduló ismértéket jelöli, folytonos ismérték esetén ez a gyakorisági görbe maximumhelye. Ha a gyakorisági sorban egynél több kiugró ismérték van, akkor az eloszlás több módusú, ilyenkor célszerű a sokaságot részekre bontva vizsgálni. A módusz nem mindig határozható meg egyértelműen, nem is mindig létezik, viszont nominális mérési szint (nem számszerűsíthető jellemzők) esetén is értelmezhető. Olyan középtérték, melyet az ismértékek helyébe téve a lehető legtöbbször nem követünk el hibát, de induktív statisztikai célokra is alkalmas.

k	Elnevezés	Jelölés	i lehetséges értéke	Lehetséges kvantilisek
2	medián	Me	1	Me
4	kvartilis	Q_i	1, 2, 3	Q_1, Q_2, Q_3
5	kvintilis	K_i	1, 2, 3, 4	K_1, K_2, K_3, K_4
10	decilis	D_i	1, 2, ..., 9	D_1, D_2, \dots, D_9
100	percentilis	P_i	1, 2, ..., 99	P_1, P_2, \dots, P_{99}

1. táblázat

A legfontosabb kvantilisek elnevezése és jelölése (1.)

A **számtani átlag** a leggyakrabban használt számított átlag, az „átlag” szóról gyakorlatilag mindenkinek ez jut eszébe. Úgy számítjuk, hogy az összes ismérvértéket összeadva az összeget elosztjuk az elemszámmal:

$$\bar{Y} = \frac{Y_1 + Y_2 + \dots + Y_N}{N} = \frac{\sum_{i=1}^N Y_i}{N} = \frac{\sum Y}{N}$$

A sokasághoz tartozó értékösszegeből, illetve a gyakorisági sor adataiból (súlyozással) is számítható. Bármely adathalmazból egyértelműen meghatározható, de meglehetősen érzékeny a szélsőséges (főleg a szélsőségesen nagy) értékekre. Az alapadatok számtani átlaggal való helyettesítése a mért összes hibát minimalizálja, így különösen alkalmas induktív statisztikai célokra.

Példaként tételezzük fel, hogy 11 beteg szérum koleszterin értékei növekvő sorrendben a következők:

Beteg	1	2	3	4	5	6	7	8	9	10	11
Secholest.	2.3	2.5	3.4	4.8	4.8	5.2	5.6	5.7	5.9	6.2	7.4

Az ismérvértékek számtani átlaga: $(2.3+2.5+\dots+7.4)/11 = 4.89$.

A medián értéke: $Me = 5.2$, hiszen ez az érték, melynél pontosan ugyanannyi ismérvérték kisebb, mint nagyobb.

A módusz: $Mo = 4.8$, mert ezt az értéket veszi fel a változó a legtöbbször, példánkban 2-szer. Látható ugyanakkor, hogy nagyjából egyenletes eloszlást feltételezve a három középérték nem tér el lényegesen egymástól.

Elnevezés	Jelölés	Számításmód a	
		súlyozatlan	súlyozott
		esetben	
Harmonikus átlag	\bar{Y}_h	$\frac{N}{\sum_{i=1}^N \frac{1}{Y_i}}$	$\frac{N}{\sum_{i=1}^k \frac{f_i}{Y_i}}$
Mértani (geometriai) átlag	\bar{Y}_g	$\sqrt[N]{\prod_{i=1}^N Y_i}$	$\sqrt[k]{\prod_{i=1}^k Y_i^{f_i}}$
Négyzetes (kvadratis) átlag	\bar{Y}_q	$\sqrt{\frac{\sum_{i=1}^N Y_i^2}{N}}$	$\sqrt{\frac{\sum_{i=1}^k f_i Y_i^2}{N}}$

2. táblázat
A különböző átlagfajták áttekintése

Tegyük fel, hogy a 11. érték nem 7.4, hanem 15.2. Ekor a számtani átlag – mivel érzékeny a szélsőségesen nagy értékekre – nem 4.89, hanem 5.6 lesz, ami 14.5%-kal nagyobb értéket jelent, a medián és a módusz értéke ugyanakkor változatlan!

A harmonikus, mértani (geometriai) és négyzetes (kvadratis) átlag ritkán használt számított átlagfajták, a 2. táblázatban részletezett módon értelmezhetők.

A **harmonikus** és **mértani** átlag olyan esetekben használható, ha nem az ismérvértékek összegének, hanem

azok reciprokából képzett összegnek vagy azok szorzatának van kézzelfogható jelentése. Példa: láncviszony-számok alkalmazásakor mértani átlagot érdemes használni, mert azok szorzata bázisviszonyszámot ad eredményül.

A **négyzetes átlag** alkalmazása akkor célszerű, ha el akarunk vonatkoztatni az átlagolni kívánt értékek előjelétől, hiszen az előjelből adódó különbséget a formula négyzetre emeléssel tünteti el, majd gyökvonással semlegesíti, így a szórás számítás alapját is képezi, hiszen itt az átlagtól való eltérés nagysága érdekes, nem az előjele. Példaként elég csak olyan biológiai jellemzőkre, laboratóriumi értékekre gondolni, ahol adott populáció (beteganyag) vizsgálata esetén a normálértéknél kisebb értékek is kórosnak tekinthetők (ellentétben pl. a májfunkciós értékekkel), így a számtani átlag alkalmazása tévedésre adna alkalmat, ha a mintában mind a normálnál magasabb, mind alacsonyabb értékű elemek vannak, a négyzetes átlag viszont kimutatja a különbséget.

Végezetül álljon itt egy összefüggés, mely minden sokaságra vonatkozathatóan az egyes átlagfajták közötti nagyságbeli összefüggést mutatja (egyenlőség akkor áll fenn, ha a minta minden eleme egyenlő):

$$Y_{\min} < \text{harmonikus} < \text{mértani} < \text{számtani} < \text{négyzetes} < Y_{\max}$$

A gyakorisági eloszlások jellemzésére a szóródási mutatókat használjuk leggyakrabban. A szóródás terjedelme annak az intervallumnak a teljes hossza, amelyen belül az ismérvértékek elhelyezkednek:

$$R = Y_{\max} - Y_{\min}$$

Előbbi példánknál maradván ez azt jelenti, hogy a szóródás terjedelme $7.4 - 2.3 = 5.1$ mmol/l.

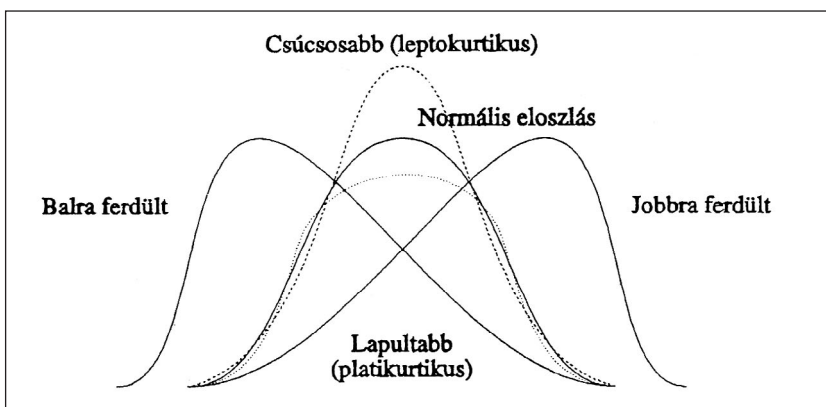
A szóródás legfontosabb mérőszáma azonban a szórás, mely az ismérvértékek számtani átlagtól vett eltéréseinek négyzetes átlaga. Azt mutatja meg, hogy az Y_i ismérvértékek átlagosan mennyivel térnek el saját átlaguktól. Az átlagtól való eltérés irányából adódó előjel eltérést (a négyzetes átlaghoz hasonlóan) négyzetre emeléssel iktatja ki, majd gyökvonással egyenlíti ki az előbbi műveletet:

$$\sigma = \sqrt{\frac{1}{N} \sum_{i=1}^k f_i (Y_i - \bar{Y})^2}$$

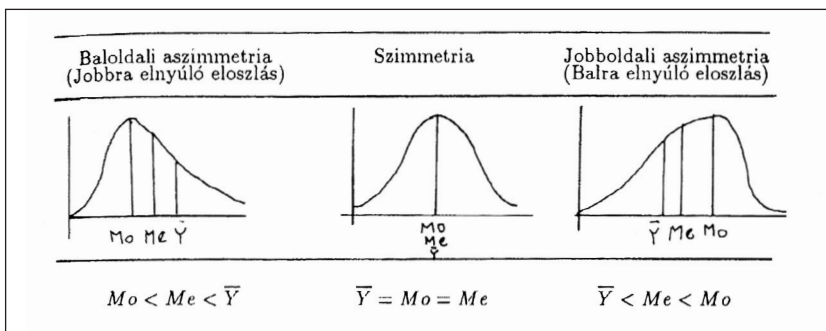
A szórás négyzetét varianciának nevezzük.

A korábbi példa nyomán a képletbe behelyettesítve a szérum koleszterin értékek szórására 1.505 mmol/l adódik, ez azt jelenti, hogy az ismérvértékek átlagosan ennyivel térnek el a saját átlaguktól. A második esetben (a szélsőséges érték fennállta esetén) a szórás 3.37 mmol/l lesz, ami több mint kétszeres értéket jelent!

A gyakorisági eloszlások alak-mutatószámai annak tömör jellemzésére szolgálnak, hogy azok milyen tekintetben és mértékben térnek el a normális eloszlás gyakorisági görbájától. Számításuknak csak egymódusú eloszlás esetén van értelme! Így az adott eloszlás a 4. ábrán bemutatottak szerint a normális eloszlásnál lehet csúcsosabb, lapultabb, ill. ferdült:



4. ábra
A gyakorisági eloszlások alakja (1)



5. ábra
Szimmetrikus és aszimmetrikus eloszlások jellegzetességei (1)

Asszimmetrikus eloszlás esetén a számtani átlag, módusz és medián egymáshoz való viszonyát mutatja az 5. ábra.

Az eloszlás **csúcsosságának** mérésére a

$$K = (Q_3 - Q_1) / 2(D_9 - D_1)$$

formula szolgál, melyben Q és D értékei a kvantilisok közül a kvartilis és decilis megfelelő rendű értékeit jelölik. Értéke normális eloszlás esetén: $K = 0.263$.

Az **aszimmetria** (Pearson-féle) mutatószáma pedig a számtani átlagból, a mediánból és a szórásból indul ki, értéke -3 és 3 között lehet, a következőképpen számítható:

$$P = \frac{3(\bar{Y} - Me)}{\sigma}$$

Példánk adataival $P = 3(4.89 - 5) / 1.505 = -0.618$, ez enyhe jobb oldali aszimmetriát (balra elnyúló eloszlást) jelez, a mintában enyhe túlsúlyban vannak a számtani átlagnál nagyobb értékek.

Az eddigiekben a sokaságot egy ismérv szerint vizsgáltuk, de arra is szükség lehet, hogy a statisztikai sokaság elemeit több ismérv szerint vegyük górcső alá.

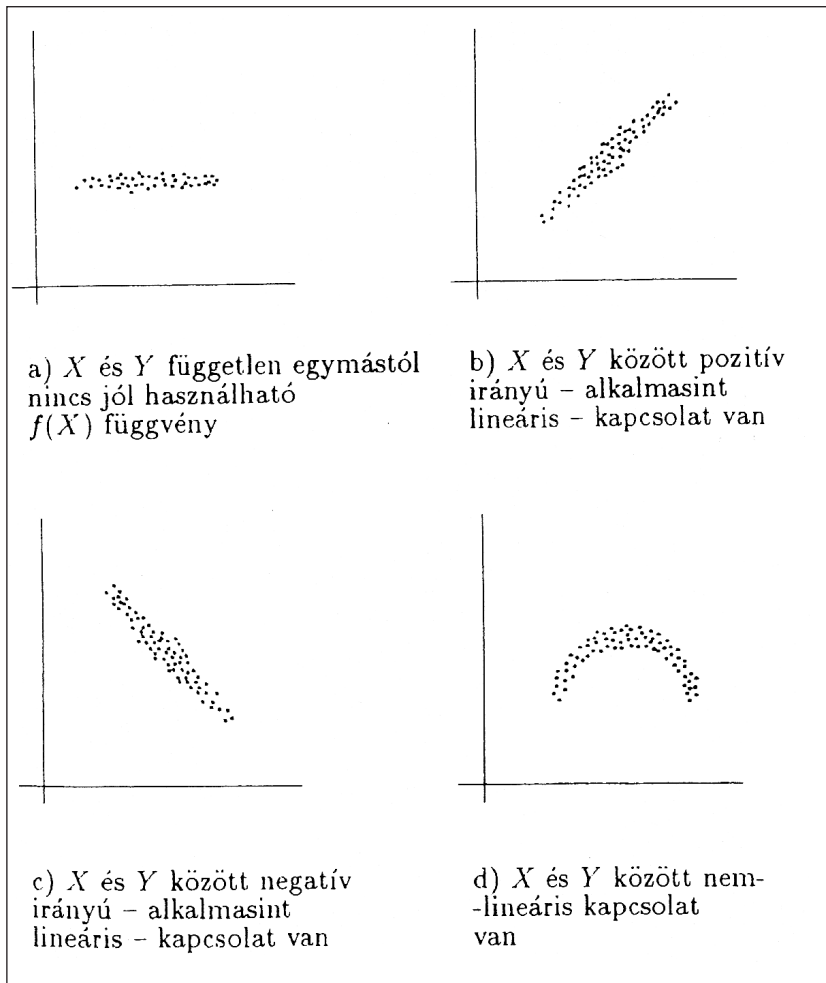
A két ismérv szerinti kapcsolat jellege alapján ennek három alapesete különíthető el:

- > Ha a két csoportképző ismérv egymástól független, akkor az egyik ismérv szerinti hovatartozás ismerete nem ad információt a másik ismérv szerinti hovatartozásról. Például ha a sokaság egy egyedének magas a szérum koleszterin szintje, ennek ismeretében nem tudjuk megmondani, hogy volt-e neki betegsége az utóbbi egy évben.
- > Ha a két ismérv között függvény-szerű kapcsolat van, akkor az egyik ismérv szerinti hovatartozás ismeretében determinisztikusan meghatározott a másik ismérv szerinti besorolhatóság is. Például ha valakit fél éve bizonyított syphilis miatt kezeltek, biztosan megállapítható, hogy a TPHA vizsgálata pozitív lesz.

> A két előbbi köztes állapotot jelenti a sztochasztikus kapcsolat, amikor

az egyik ismérv hatással van ugyan a másikra, de nem határozza meg azt egyértelműen, az egyik ismérv szerinti hovatartozás ismerete csak csökkenti a másik ismérv szerinti hovatartozást illető bizonytalanságot. Ilyen eset például az, ha egy beteg magas szérum koleszterin szintje esetén nagyobb eséllyel várhatjuk, hogy a triglicerid szintje is magasabb lesz, de annak konkrét értékét nem tudjuk megmondani csak a koleszterin szint ismeretében.

A kapcsolat jellegét a 6. ábra szemlélteti:



6. ábra

Ismérvek közötti kapcsolat jellege (1.)

Mennyiségi ismérvek esetén a sztochasztikus kapcsolat erősségét a korrelációs együtthatóval jellemezzük, számításakor az egyedi ismérvértékek saját átlaguktól számított eltéréseiből indulunk ki (a szóráshoz hasonlóan):

$$r = \frac{\sum dx dy}{\sqrt{\sum d_x^2 \sum d_y^2}}$$

Korreláció esetén annak iránya is értelmezhető, hiszen pozitív irányú kapcsolat esetén x értékének növelésével y párhuzamosan nő, negatív irányú kapcsolat esetén csökken, a lineáris korrelációs együttható értéke ennek meg-

tendenciájú kapcsolat van. A lineáris regresszió függvénye a

$$\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_i$$

formulával adható meg, paramétereinek számításához szintén a vizsgált változók értékeinek átlagtól való eltérést veszik alapul, melyek levezetésére és értelmezésére terjedelmi okokból nem érdemes kitérni.

Példánknál maradva a szérum triglicerid szintje (Y) a következő függvénnyel becsülhető a koleszterin szint (X) ismeretében:

$$Y = -0.47 - 0.68X$$

Beteg	1	2	3	4	5	6	7	8	9	10	11
Secholest.	2.3	2.5	3.4	4.8	4.8	5.2	5.6	5.7	5.9	6.2	7.4
Se triglic.	1.4	1.7	1.6	1.8	1.9	3.2	3.0	3.4	3.6	5.6	4.2

3. táblázat

felelően -1 és 1 között lehet. A 0 értéket akkor veszi fel, ha a két ismerv között nincsen korrelációs kapcsolat, két szélső értéke esetén pedig determinisztikus (függvényszerű) kapcsolat áll fenn.

Tegyük fel, hogy az előbbieken vizsgált 11 betegől nem csak a szérum koleszterin, hanem a triglicerid érték is rendelkezésünkre áll, és az a következőképpen alakul (3. táblázat).

Feltételezhetjük (a táblázatból is látszik), hogy a két ismerv egy irányban „mozog”, közöttük kapcsolat van.

Az ismérvértékek átlagtól való eltérése segítségével levezetett is kiszámolt korrelációs együttható értéke a számszaki levezetés nélkül: $r = 0.814$. Ez azt jelenti, hogy a két ismerv egymástól nem független, a szérum triglicerid és koleszterin szintje között meglehetősen erős, pozitív irányú korrelációs kapcsolat van. A szérum koleszterin szintjének ismerete jelentősen csökkenti a triglicerid szintre vonatkozó bizonytalanságot.

A regressziószámítás az összefüggésekben levő sztochasztikus tendenciát vizsgálja, és a kapcsolat természetét valamilyen függvénnyel írja le. A leggyakrabban alkalmazott modellek egyike a **lineáris regressziós modell**, mely akkor ad valószerű képet a két jellemző kapcsolatáról, ha azok között lineáris

Pl. a 6. beteg esetére alkalmazva: $Y = -0.47 + 0.68 \cdot 5.2 = 3.066$ adódik, ami jól közelíti a valóságban talált 3.2-es szérum triglicerid értéket.

Az exponenciális regressziófüggvény alkalmazására akkor kerül sor, ha valamilyen jelenség növekedése függ a je-

lenség már elért színvonalától (infláció, kamatos kamat, a biológiában az önerősítő folyamatok, pl. véralvadási kaskád, komplementrendszer). Általános képlete a következő:

$$\hat{Y} = \beta_0 \beta_1^X$$

A hatványkitevős regressziófüggvényt akkor érdemes választani, ha az X és Y változók logaritmusai között van lineáris összefüggés. Példa erre a szerológiai reakciók területének vizsgálata, számítási formulája alább látható:

$$\hat{Y} = \beta_0 X^{\beta_1}$$

Mi történik akkor, ha – mint ahogy az esetek többségében ez így van – a statisztikai elemzéshez csak részletes, és nem teljes körű adatfelvétel eredményei állnak rendelkezésre? A reprezentatív mintavétel nem csak adatszerzési mód, hanem a statisztikai következtetések alapja is. Azt jelenti, hogy a vizsgált sokaság a vizsgálat szempontjából releváns karakterisztikumát tekintve megoszlásában megegyezik az alapsokasággal.

Csak felsorolásszerűen az alábbi mintavétel módszereket különböztetjük meg:

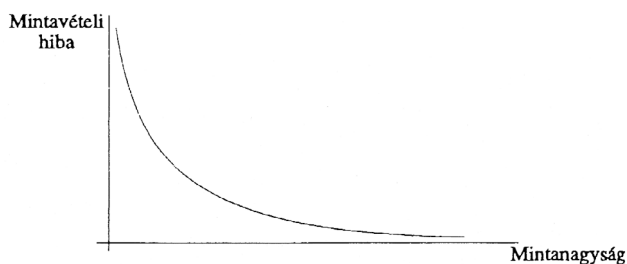
- > Véletlen mintavételi eljárások: független, azonos eloszlású (FAE), egyszerű véletlen, rétegezett, csoportos és többlépcsős mintavétel.
- > Nem véletlen eljárások: szisztematikus, kvótás, koncentrált, hólabda és önkényes kiválasztás.

Nem véletlen eljárások esetén a kapott eredmények torzítottak lehetnek, de kis létszámú, nehezen felderíthető populációk (pl. homoszexuálisok) vizsgálatára, szokásainak feltárására jól alkalmazhatók. A hólabda eljárásnál például a soron következő kikérdezett személyére mindig az előző vizsgálati alany tesz javaslatot.

Ha nem ismerjük a sokaság minden elemét, megadhatjuk azt eloszlásával is, melyet a mintából számíthatunk. A tipikus eloszlások (normális, lognormális, exponenciális) hisztogramját lásd az előzőekben.

Helytálló statisztikai következtetések levonásához korrekt mintavételi technika szükséges, sikerének alapköve a megfelelő mintanagyság, hiszen nagy minták esetén a mintából számított jellemzők jó része normális eloszlásúvá válik, így kezelésük leegyszerűsödik. A mintavételi hiba a mintanagyság növelésével párhuzamosan csökken, ezt szemlélteti a 7. ábra.

De mekkora legyen ez a minta? Minden helyzetre érvényes útmutatást nehéz adni, de általában tanácsos az alábbi szabályt követni: Szimmetrikus, vagy azt közelítő eloszlások esetén viszonylag kis elemszámú minták ($30 < n$) is elegendőek, de a szimmetrikustól erősen eltérő eloszlások esetén több 100-as mintanagyság lehet kívánatos, s ugyanakkor bizonyos származtatott jellemzők becslésének pontossága sem lesz kielégítő. 100 alatti minta-

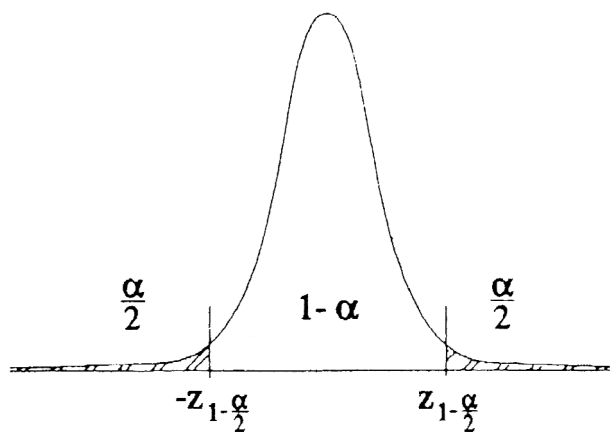


7. ábra

Mintavételi hiba és mintanagyság összefüggése (1.)

nagyság esetén a torzító tényezők miatti százalékszámításnak sincsen értelme, mert véletlen körülmények esetleges torzító hatása miatt hamis eloszlási adatokhoz juthatunk!

A statisztikai módszertan fontos területét képezi a hipotézisek vizsgálata. Lényege az, hogy egy vagy több valódi vagy fiktív sokaságról állítunk valamit, majd a rendelkezésünkre álló véletlen minta vagy minták alapján megvizsgáljuk azt, hogy a szóban forgó állítás mennyire hihető. Menete: a nullhipotézis és a vele szemben álló alternatív hipotézis felállítása után keresünk egy próbafüggvényt, mely alapján az egyik hipotézist elfogadhatjuk, a vele szemben állót pedig elvetjük. Ha a próbafüggvény értéke az elfogadási tartományba esik, úgy nullhipotézisünk elfogadható, a vele szemben álló komplementer (alternatív) hipotézis elvethető. A kritikus tartományba való esés a valószínűségét szignifikancia-szintnek nevezzük, ez a valószínűsége annak, hogy az egyébként helyes kiindulási hipotézist elvetjük. Értékét ennek megfelelően kicsinek (általában 5%) érdemes választani. Szemléltetésére a 8. ábra szolgál.



8. ábra

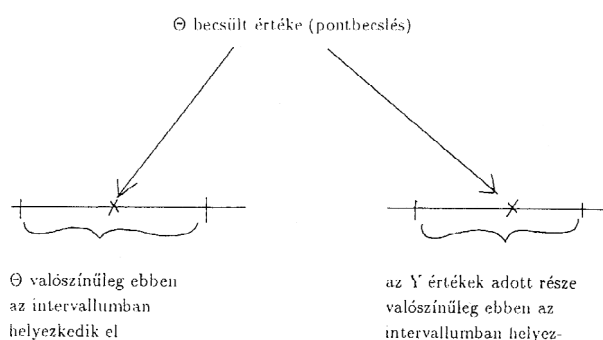
A szignifikancia-szint megválasztása (1.)

Ha a problémát a másik oldaláról közelítjük meg, akkor meghatározhatjuk azt a legkisebb szignifikancia-szintet, melynél a nullhipotézis már éppen elvethető az alternatív hipotézissel szemben. Ez az ún. p-érték.

A statisztikai információgyűjtés célja gyakran a változó várható értékének meghatározása annak érdekében, hogy

abból az egész populációra érvényes következtetéseket vonhassunk le. A várható értékre irányuló próbák az u-próba (másnével z-próba), a t-próba és az aszimptotikus u-próba. A próba alkalmazhatósága a mintanagyságtól, illetve minta szórásának vagy eloszlásának ismeretétől függ. Mindegyik próbának létezik egy, illetve több mintára vonatkozó változata is, a számításokat már erre a célra kifejlesztett szoftverek segítségével végzik.

A statisztikai számításokhoz felhasznált adatok kísérletekből és megfigyelésekből származhatnak. Statisztikai becslésen azt az eljárást értjük, amellyel a mintából számított mutatók segítségével következtethetünk az alapsokaság ismeretlen jellemzőire. Ehhez meg kell határozni egy intervallumot, mely nagy valószínűséggel tartalmazza az eloszlás előre meghatározott részét, míg a gyakrabban használt konfidencia intervallum az adott, ismeretlen jellemzőt (9. ábra).



9. ábra

Konfidencia és tolerancia intervallum (1.)

A különféle becslések nagyon szerteágazó és bonyolult részét képezik a statisztikai tevékenységének, így részletes tárgyalásukra most nem kerül sor.

A továbbiakban néhány típusos hibát érdemes áttekinteni, melyekkel gyakran találkozhatunk statisztikai tárgyú orvosi közleményekben. Gyakorlati példaként álljon itt egy cikk, mely amerikai szerzők tollából származik és több típushibát tartalmaz. Az írás a Journal of the American Academy of Dermatology 2001. júliusi számában jelent meg, a bőrgyógyászati életminőségi indexekkel foglalkozik (5).

A cikkben a bőrgyógyászati betegek életminőségét vizsgáló index használhatóságát elemzik, melynek alapjául egy 1-3 perc alatt kitölthető önkitöltő kérdőív szolgál. Utóbbit 200, egymás után jelentkező bőrbeteggel töltették ki. Az eredményeket egy korábbi vizsgálat adataival vették össze. A 4. táblázatban részletezett eredményekre jutottak:

Melyek azok a bizonytalanságok, melyek felmerülnek az eredmények olvasásakor?

- Annak ellenére, hogy az indexet jól használhatónak találták, lényeges különbségek vannak az egyes betegcsoportokban az index értékét tekintve.
- Konzekutív módon jelentkező betegek vizsgálatába vonásával egyáltalán nem biztos a véletlen kiválasztás.
- Kevés karakterisztikum figyelembe vételével (pár perc alatt kitölthető kérdőív!) kaphatunk-e pontos információkat egy olyan összetett, nehezen kvantálható jellemzőre, mint az életminőség?
- Az összes esethez képes sok volt a nem meghatározott kórképben szenvedők aránya (az első esetben 85/197, az összes eset 43%-a).
- Egyes csoportokban túl kevés beteg szerepel ahhoz, hogy ebből induktív elemzés végezhető legyen.

Diagnózis	Betegszám (IU)	Betegszám (FK)	DLQI-IU középérték	DLQI-IU szóródás	DLQI-FK középérték	DLQI-FK szóródás
Acne	23	18	9.0	0-19	4.3	0-11
Eczema	41	17	8.5	0-29	8.6	2-27
Dermatitis atopica	6	13	5.8	1-9	12.5	6-23
Psoriasis	10	52	4.6	1-14	8.9	0-28
Keratosis seborrhoica	8	5	3.6	0-12	1.8	1-3
Verruca vulgaris	9	12	2.9	0-7	6.7	2-22
Keratosis solaris	2	5	2.9	0-6	3.4	2-6
Naevus	3	7	2.3	1-7	1.0	0-4
Basalioma	6	8	1.8	1-4	2.0	0-6
Pruritus	0	9	-	-	10.5	3-22
Egyéb	85	54	6.3	0-24	6.9	0-28
Ismeretlen	4	0	7.5	0-15	-	-
Összesen	197	200	6.5	0-29	7.3	0-28

4. táblázat

Az életminőség mérése különböző bőrbetegségekből (DLQI: Bőrgyógyászati életminőség index, IU: Indiana Egyetem eredményei, FK: Finlay és Kahn eredményei)

- Nincs információ arról, hogy az egyes betegcsoportokon belül milyen volt a betegek súlyosság szerinti megoszlása, ill. arról sem, hogy a betegség mely fázisában töltötték ki a kérdőívet.

Az előbbi észrevételek fényében nehéz a cikk következtetéseit elfogadni, bár lehet, hogy a vizsgálat a módszertani alapelvek figyelembevételével történt, de ennek részletezése nem szerepel a közleményben.

A statisztikai tevékenység négy fázisa a tervkészítés (előkészítés), adatgyűjtés, feldolgozás-elemzés és az eredmények közzlése. Ezek során a statisztikus és a többi szereplő közötti viszonyt meghatározó magatartási normák összessége a **statisztikai tevékenység etikája**.

Ennek legfontosabb elemei:

- Kompetencia tisztázása, szakmai kérdésekben a statisztikus kapjon szabad kezet!
- Statisztikus és adatszolgáltató közötti jó kapcsolat a válaszmegtagadás és a pontatlan adatszolgáltatás kiszűrésével.
- Részletes jelentés szükséges az adatok megszerzési módjáról, a használt fogalmak pontos definíciójáról, az adatfelvétel körülményeiről, a mintavételi eljárásokról, esetleg kérdőívvel, kísérleti naplóval kiegészítve.
- Az adatok elemzésére használt eszközök helyes alkalmazása, az elemzési módszerek alkalmazási feltételeinek ismertetése és az eredmények interpretálása.
- A rekonstruálhatóság érdekében az adat vagy elemzési eredmény keletkezésének minden lényeges mozzanata jelezve legyen, különösen érvényes ez a közpénzekből fedezett kutatások esetén!

- A statisztikával való visszaélésnek minősül, ha az félreinformálja a felhasználót. A nem elég precíz fogalmak, „laza” definíciók félreértés, téves következtetés forrásai lehetnek.

A statisztikával kapcsolatos visszaélések elleni védekezés alapja a következtetésekkel szembeni egészséges kételkedés, fel kell tenni a kérdést, hogy az adott eredményt milyen adatok alapján, milyen módon eljárva és milyen feltételezések mellett kapták.

Összefoglalva: A tudományos kutatások eredményeinek interpretációjakor megfelelő önmérséklet, ill. a statisztikai módszertan alapelveinek és a különféle mutatószámok használatával kapcsolatos korlátok ismerete szükséges, hogy az adatokat megfelelő kritikával értékelhessük. Ehhez próbált a cikkben a szerző szemléletbeli útmutatást nyújtani.

IRODALOM

1. *Hunyadi László, Vita László, Mudruczó György: Statisztikai I-II.* Aula Kiadó, Budapest 1992.
2. *Hajdú Ottó, Hunyadi László: Statisztikai elemzések.* Oktatási segédlet. Aula Kiadó Budapest 1996.
3. *Zalai Ernő: Matematikai közgazdaságtan.* KJK-Kerszöv Jogi és Üzleti Kiadó Budapest 2000.
4. *Hal R. Darian: Mikroökonómia középfokon.* KJK-Kerszöv Jogi és Üzleti Kiadó Budapest 2001.
5. *Hahn B et al: Use of the Dermatology Life Quality Index (DLQI) in a midwestern US urban clinic.* J. Am Acad. Dermatol 2001; 45, 44-48.

Érkezett: 2003. III. 24.

Közlésre elfogadva: 2003. VII. 23.