

Országos Munkaegészségügyi Intézet, MTA SZTAKI

A cluster analízis alkalmazása respiratorikus  
syndrómák diagnózisában

Csukás Andrásné, Mándi András, Galgóczy Gábor,  
H. Gaudi István

Krónikus légzőszervi aspecifikus betegségek /KALB/ csoportjába a tüdő nem tbc-s, nem daganatos, nem keringési eredetű és nem a mellkas deformitásán alapuló megbetegedései tartoznak. Az utóbbi időben a környezeti ártalmak fokozódása miatt egyre jobban előtérbe kerül ezen betegségek epidemiológiájának vizsgálata. Idősebb korban a lakosság mintegy 30 %-a szenved ilyen típusu betegségben és gyakori halálóki tényező is. Így nagy jelentőségű ezen betegségek időbeni felkutatása és a kiszűrt betegek megfelelő gondozása.

Az epidemiológiai vizsgálatok ismert módszerei közül a KALB felkutatására felhasználhatjuk a standardizált kérdőíves módszert, a beteg orvosi vizsgálatát és a légzőszervek funkcionális állapotának műszeres ellenőrzését. A kérdőívek kitöltése a betegek bemondásán alapszik, így az egyén szubjektivitása nagymértékben érvényesülhet, de a módszer előnye, hogy hosszabb, a felmérés előtti periódusra adhat felvilágosítást. A beteg megtekintésén, valamint a kopogtatáson és hallgatódzásokon alapuló orvosi vizsgálat a beteg pillanatnyi állapotáról ad képet, de az észleletek megítélése a vizsgáló orvos gyakorlatának, szaktudásának és pillanatnyi diszpozíciójának függvénye. A műszeres légzésfunkciós vizsgálatok - amennyiben módunk van a beteg kollaborációs készségétől független paraméterek mérésére - hibalehetőségként csak a

mérés pontatlanságát tartalmazzák, sok esetben ugyan csak a pillanatnyi állapotot tükrözik, de a kóros elváltozásokat bizonyítják. A fentiekből is kitűnik, hogy a beteg minden szempont szerinti kielégítő felkutatása a három módszer együttes alkalmazásával várható. A jelenlegi vizsgálat célja annak eldöntése, hogy a három vizsgálati módszer által szolgáltatott adatok mennyire fedik egymást; ezáltal azt akartuk eldönteni, hogy melyik az optimális vizsgálati módszer és szükséges-e mindhárom vizsgálat együttes elvégzése.

Légzőszervi ambulanciánk 800, 20-70 év közötti férfi betegén párhuzamosan mind a három vizsgálati módszert alkalmaztuk. Ezen összehasonlító értékelésre ez a beteganyag megfelelő volt, mert azon kívül, hogy elegendő számu és egészséges egyén is szerepel, a különböző kórformák is elegendő gyakorisággal fordulnak elő. Ezen utóbbi feltételeket egy tisztán epidemiológiai jellegű felmérésnél csak nagyon nagyszámu egyén megvizsgálása által lehetne biztosítani.

Vizsgálati célkitűzésünk a légzőszervi szempontból egészséges egyének, az idült hörghurutban, /krónikus bronchitis/, hörgőszűkülettel nem járó tüdőtagulatban /nonobstruktív emphysema/ illetőlet hörgőszűkülettel együtt járó tüdőtagulásban /obstruktív emphysema/ szenvedő betegek szétválasztása volt. Ezért kérdőíven rögzítettük a köhögésre, köpetürítésre és nehézlégzésre vonatkozó panaszokat; ezek alapján 5 panasz-kategóriát hoztunk létre. A fizikális vizsgálat során a tüdőtagulásra, a hörghurutra, a szív- és érrendszer betegségeire és a testalkatra vonatkozó adatokat rögzítettük; ezek alapján 9 csoportot hoztunk létre. Minden betegnél részletes testpletizmográfus, spirometriás és vérgáz analitikai vizsgálatot végeztünk a testméretek meghatározásával kiegészítve. Az így nyert 11 funkciós érték közül a most ismerttetendő

értékelési eljárásnál 4 funkciós változót használtunk fel: a léguti áramlási ellenállás [resistance], residualis volumennek a totalis kapacitáshoz viszonyított értéke [RV/TC %], az artériás vér oxigén tensioja [ $\text{PaO}_2$ ] és a relativ testsúly [Broka index]. Ezek voltak azok, melyeket a beteg kollaborációs készsége nem, vagy csak alig befolyásol és csak a relativ testsúlytól függők.

Megfelelő matematikai statisztikai módszerekkel analizáltuk, hogy milyen légzésfunkciós elváltozásokkal járnak együtt az előbbi szempontok szerint kialakított panaszkategóriák, illetve a fizikális vizsgálat során észlelt eltérések. Megállapítottuk, hogy a panaszok és a légzésfunkciós értékek közötti kapcsolat laza és csak a legsúlyosabb panasz-kategória esetén észleltünk kóros légzésfunkciós értékeket. A fizikális leletek és a légzésfunkciós értékek közötti kapcsolat szorosabbnak bizonyult; a súlyosabb fizikális leletek súlyosabb légzőszervi elváltozásokkal járnak együtt, mint a legsúlyosabb panasz-kategória. Ezen eredményeinkről más helyen részletesebben is beszámoltunk.

Mivel a KALB epidemiológiai vizsgálata rendkívül munka- és időigényes, ezért további analiziseket végeztünk annak eldöntésére, hogy a beteg oldaláról szubjektív kérdőíves, illetve az orvos részéről szubjektív fizikális vizsgálat adatai és az objektívnek tartható légzésfunkciós paraméterek mennyire fedik egymást. Ezek alapján kívántunk ugyanis véleményt nyilvánítani a három vizsgálati módszer értékéről. Ilyen problémák megválaszolására alkalmas módszer a cluster analízis.

Green, Frank és Robinson 1967-ben a cluster analízisre a következő definíciót adta: "a cluster analízis egy általános név, amely olyan módszerek összességére vonatkozik, mely-

nek célja, hogy azonosítson hasonló dolgokat azon jellemzőkből [minőségi vagy mennyiségi], amelyek a dolgok tulajdonságai. Eltérően más eljárásoktól - mint például a diszkriminancia analízis - nem ismert előre, hogy mely dolgok tartoznak egy csoportba. Értéke, hogy az adatoknak egy olyan elő-csoportosítása hajtható végre, melyet az adatok természetes csoportosulása sugall".

Egy cluster analízis végrehajtása során két problémát kell megoldani; egyrészt a metrika [távolság/ definiálása, másrészt az algoritmus kiválasztása.

Legyenek  $X_1, X_2, \dots, X_n$  megfigyeléseink, amelyek kölcsönösen független és azonos eloszlású vektor-változók a  $p$ -dimenziós térben, azaz  $X_i' = (x_{i1}, x_{i2}, \dots, x_{ip})$ . Esetünkben  $p=4$ . A cluster analízisben a leggyakrabban azt tételezik fel, hogy a tér Euklidesi a szokásos Euklidesi metrikával, azaz

$$S_E(X_i, X_j) = \sqrt{(X_i - X_j)'(X_i - X_j)}.$$

Itt feltétel, hogy a vektor komponensek függetlenek. /A kovariancia mátrix diagonális./ Nem független vektor komponensek esetén a Mahalanobis metrika a megfelelőbb

$$S(X_i, X_j) = \sqrt{(X_i - X_j)' V^{-1} (X_i - X_j)},$$

ahol  $V$  a kovariancia mátrix. Ha az  $X_i$  vektor komponensei különböző fontosságúak a végső következtetés szempontjából, egy megfelelően választott diagonális mátrix-szal lehet súlyozni

a komponenseket, ekkor a távolság

$$\rho(x_i, x_j) = \sqrt{(x_i - x_j)' U' V^{-1} (x_i - x_j)} .$$

Ha a  $V$  mátrix ismeretlen, akkor helyettesíthető a megfigyeléseink alapján vett becslésével:

$$\tilde{V} = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})' ,$$

ahol

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i .$$

A clusterek képzésére nagyon sok algoritmus ismeretes, ezek mindegyike az  $x_1, x_2, \dots, x_n$   $p$ -dimenziós megfigyeléseinket a mintatér  $S = (S_1, S_2, \dots, S_k)$  cluster-rendszere páronként diszjunkt clustereinek valamelyikébe helyezi. A cluster analízis leggyakrabban használt algoritmusai a hierarchikus algoritmusok. Ezek közül is a legismertebb a Ward-tól származó [1963] eljárás. Az adatok csoportosítása szakaszonként történik. Az első szakaszban meg kell keresni az  $n$  számú  $x_1, x_2, \dots, x_n$  sokaság azon két elemét, melyek a választott metrika szerint legközelebb vannak egymáshoz. Ezt a két pontot középpértékeikkel helyettesítjük. Ezzel  $n-1$  pontunk maradt a  $p$ -dimenziós térben. A második lépésben ezek közül kell kiválasztani a két legközelebbit, és így tovább. Az eljárás egy megfelelő  $k$  értékig folytatódik. Azon pontok alkotnak egy csoportot, amelyek egy cluster középpontot kialakítottak.

A hierarchikus eljárások hátránya, hogy a cluster rendszer kialakításához az összehasonlitások igen nagy száma szükséges. Ez a szám  $\sum_{l=k-2}^{n-1} l!$ . Mivel a mi esetünkben  $n=800$ , így ez az eljárás alkalmazhatatlan. Ezért az ugynevezett k-közepű eljárást használtuk. Ez k csoportból indul, amelyek mindegyike egy  $X_i$  véletlen pontból áll [ $i=1,2,\dots,k$ ] és ezután hozzáadja a többi  $X_{k+1}, \dots, X_n$  pontok mindegyikét ahhoz a csoporthoz, melynek átlaga legközelebb van az adott  $X_i$  [ $i=k+1,\dots,n$ ] ponthoz. Az  $X_{k+r}$ -edik pontot hozzáadva valamelyik csoporthoz, ennek átlagát kiigazítjuk az  $X_{k+r}$ -edik pont figyelembevételével. Így minden egyes szakaszban k csoport van. Egy ilyen eljárás megvalósítása  $k(n-k)$  távolság összehasonlitást igényel. Az eljárás hátránya a rögzített k középszám. Ez feltételezi azt az előzetes információt, hogy az adatok k csoportba tömörülnek.

A k-közepű eljárás egy módosítása kiküszöböli a k értékének előre való megadását, a clusterek végső száma a cluster analízis folyamata közben alakul ki. Ez a módosított algoritmus MacQueen /1967/ cikkében szerepel. A módosítás lényege két előre megadott paraméter: C és R. A C az az érték, melynél közelebbi clustereket összeolvasztunk, míg azokat a pontokat, melyek minden clustertől R értéknél messzebb vannak egyik clusterhez sem vonjuk hozzá, hanem új cluster középként kezeljük. Az algoritmus használata közben a clusterek száma hol csökken, hol növekszik, a clusterek végső számát az adatok belső, sajátos csoportosulása határozza meg.

Ezt a módszert használtuk az adataink analízisére. A program kis módosítással bármilyen számú és dimenziójú adathalmazra használható. Az eljárás a C és az R értékétől függően különböző mértékben függ a kiinduló középpontok megválasztásától. Ha úgy választjuk a C és R konstansokat, hogy eljárásunk független legyen a kiinduló értékektől, akkor megnő a gép-  
idő.

Igy előnyösebb a kettő ésszerű egyeztetése. Nálunk  $C=1$ ,  $R=5$  vagy  $C=1$ ,  $R=4$  volt. Ezért a kiindulási középpontokat mindig teoretikus alapon adtuk meg.

Az első feldolgozás során a légzésfunkciós állapot megítélésére a panasz-kategóriáknál alkalmasabbnak bizonyult fizikális lelet-kategóriáknak megfelelő középpontokból indultunk ki. Az ezek alapján megadott középpontokból kiindulva 22 csoport alakult ki, úgy, hogy az eredetileg megadott középpontok közül néhány összevonódott. A legnagyobb csoportban az esetszám 168 volt, 14 csoportban volt 10, vagy annál több megfigyelés, a többiben 5 vagy ennél kevesebb. A cluster analízis által kialakított csoportok és a fizikális lelet csoportok között az összefüggés igen laza volt. Az egyes clusterek középpontjainak vizsgálatánál úgy tűnt, hogy számos csoport létrejöttében - a főleg fiatalabb egyéneknél fennálló - átlagosnál kedvezőbb légzésfunkciós értékek játszottak szerepet. Ezért úgy gondoltuk, hogy egy transzformációt kell alkalmazni, mely a fiziológiai értékeket a fiziológiai határértékhez vonja össze. Az így transzformált értékeket használtuk a clusterek kiépítéséhez. Az egyes légzésfunkciós paramétereket légzésfiziológiai jelentőségüknek megfelelően a távolság kiszámításában különböző konstansokkal súlyoztuk.

A transzformációk eredményeképpen kialakult clusterek száma nem változott, de feltételezésünknek megfelelően megnőtt a legnagyobb létszámú - légzésfunkciós szempontból egészséges - csoport létszáma, csökkent a 10 vagy annál több esetet tartalmazó csoport, és az előre megadott 9 csoport közül is több olvadt be más, előre megadott középpontba. A fizikális-leletkategoróriákkal való egyezés továbbra is laza maradt.

A súlyozás a clusterok kialakulását lényegesen nem befolyásolta.

Ezen eredményekből azt a következtetést vonhatjuk le, hogy a fizikális leletek alapján választott kiindulási középpontok a funkcionális paraméterek csoportosulásának belső strukturáját nem tükrözik, azaz a fizikális lelet-kategóriák és a légzésfunkciós eltérések nem fedik egymást. Tehát az adatok nem a fizikális leletek szerint válnak szét.

Az egyes cluster középpontok elemzése alapján úgy tűnt, hogy a nagyobb elemszámú csoportok mindegyike különböző légzésfiziológiai szindrómákat tükröz. Ennek alapján 10 új középpontot jelöltünk ki, melyek a tüdőtágulás, a hörgőszűkület, az elhízás, és a vér oxigén tenziójának különböző fokozatait veszik figyelembe. Ezen kiindulási középpontok esetében is alkalmaztuk a változók transzformációját és a súlyozást. Eredményeink a következők:

Transzformáció és súlyozás nélkül 25 cluster alakult ki, a legmagasabb elemszám 144 volt, 12 csoportnál volt az elemszám 10 vagy annál nagyobb. Valamennyi kiindulási középpontnak megfelelő cluster kialakult és valamennyiben az elemszám 10-nél nagyobb volt. Csak két jelentősebb új csoport alakult ki 33, illetve 60-as elemszámmal.

Transzformáció alkalmazásával 26 cluster alakult ki, de ezek közül csupán 10 csoportban nagyobb az elemszám 10-nél. A legnagyobb elemszám - 301 - a fiziológiásnak ítélt csoportban volt. A kiindulási középpontok által meghatározott clusterok ebben az esetben is kialakultak. A különböző súlyozások az itt elmondottakhoz képest lényeges változást nem hoztak.

A légzésfiziológiai szindrómák alapján kijelölt kiindulási középpontok tehát az adatok belső strukturáját jól tükrözik, az analízis során a középpontok nem olvadtak össze, csak néhány előre nem megadott cluster alakult ki. Azonban ezek a csoportok is létező légzésfiziológiai szindrómát adnak meg, melyek egymástól és az előre megadottakból élesen elkülönülnek. Ezen szindrómák előfordulására a csoport közepek megadásánál nem gondoltunk.

Elméletileg a funkciós változóknak négy fokozatát különböztetjük meg: normál vagy fiziológiás, enyhén, közepesen, illetve súlyosan kóros értéket. Teoretikusan tehát  $4^4$  azaz 256 csoport kialakulásával kellene számolni. Ezek közül azonban egyes kombinációk nem fordulhatnak elő, mivel a vizsgált változók részben együtt járhatnak, részben következményei lehetnek egymásnak. Így pl. egy feltűnően kövér, súlyosan obstruktív emphysemás betegnél normális vérgáz érték nem feltételezhető. Egyes csoportok pedig azért nem fordulnak elő, mert az általunk vizsgált ambuláns beteganyagban a súlyos elváltozások együttes előfordulása nem reális feltételezés, ezek a személyek már ágyhoz kötött kórházi betegek.

A légzésfiziológiai szindrómák alapján kialakult nagyon jól elkülönülő és magyarázható clusterok a fizikális lelettel laza kapcsolatot mutatnak. Ennek oka lehet, hogy a fizikális eltérések nem minden esetben járnak együtt funkció károsodással, valamint a fizikális eltérések megítélése - különösen súlyosság tekintetében - igen sok szubjektív komponenst tartalmaz. De ez nem is várható el és nem is feladata a belgyógyászati rutinvizsgálatnak. A fizikális lelet és a funkciós eltérések közötti laza kapcsolat magyarázata egyben annak, hogy a fizikális leletek alapján kialakított kiindulási clus-

ter közepek az analízis során részben összeolvadtak, illetőleg olyan csoportok alakultak ki, melyek fizikális eltérések alapján nem különülnek el.

Vizsgálataink alapján tehát megállapíthatjuk, hogy a cluster analízis alkalmas módszer és jól értékelhető felvilágosítást ad többváltozós, nagytömegű adathalmazok belső kapcsolatairól. Levonhatjuk továbbá azt a következtetést is, hogy a panaszok, a fizikális lelet és a légzésfunkciós állapot felmérésére egyaránt szükség van az epidemiológiai vizsgálatokban. Ezek ugyanis nem pótolják, hanem kiegészítik egymást és csak ezek alapján lehet kialakítani egy olyan szűrési rendszert, mely nemcsak a statisztikai összefüggésekre irányul, hanem lehetőleg az összes beteg embert is kiemeli.