

1. Modellalapú adatfeldolgozás

1.1. Vizsgálati eredmények statisztikai feldolgozásának módjairól

Megfigyelések és mérések eredményeit, az *adatokat* azért dolgozzák fel, hogy azok között *összefüggéseket* ismerjenek fel, feltéve, hogy ilyenek léteznek. Más szavakkal megkísérlik az adatok között esetleg fennálló valamilyen *adatstruktúra* felismerését.

Ezek az összefüggések azon a szinten, amelyen vizsgálódunk, végső soron bizonyára *okságok*, *kauzálisok*, alkalmasint azonban közvetettek, bonyolultak, kísérletileg végtelenen nehezen vizsgálhatók, a vizsgálat idején gyakorlatilag felismerhetetlenek. Lényeges tehát, hogy a jelenségekről, amelyeket az adatok leírnak, ismeretesek-e bizonyos előzetes, *a priori* ismeretek, vagy nem.

Ha nincsenek előzetes ismeretek, azaz bizonyos oksági összefüggések csak sejtethők, akkor a *feltáró*, *felderítő statisztika* módszereit (leírás, csoportosítás, főkomponens analízis, faktoranalízis, korreláció) alkalmazzák, amelyek eredményei *önmagukban is nagyon hasznosak*, mivel azonban a feltáró statisztika eredményei között az említett kauzális összefüggések nem nyilvánvalóak, értelmezésükhöz nagy figyelem és adott esetben további kutatás kell. Sokszor ezekre az eredményekre támaszkodik a további. már nem feltáró jellegű statisztikai feldolgozás.

Ha adottak előzetes ismeretek, akkor kívánatos, hogy azok *matematikai* összefüggések alakjában legyenek jelen. Ezek az összefüggések szerencsés esetben az *oksgai* viszonyokat tükröző *függvények*, gyakran differenciálegyenletek megoldásai. Egy-egy ilyen matematikai összefüggésről azonban nem tételezhető fel, hogy a vizsgált jelenséget hibátlanul, kimerítően leírja, annál is inkább nem, mivel megfogalmazása során gyakorta hasznossági, egyszerűségi szempontokat is figyelembe vettek. Nem meglepő tehát, hogy a jelenséget a megkívánt vagy elérhető pontossággal több versengő összefüggés is leírhatja.

Bizonyos esetekben (pl. "empirikus formulák" megalkotásánál, folyamat irányító algoritmusoknál) eleve lemondanak arról, hogy a jelenségeket, a vizsgált rendszert leíró összefüggés tagjainak fizikai tartalma legyen.

Mindezek miatt a jelenségek leírására felhasznált matematikai képleteket, ebben a szöveggörnyezetben éppúgy *modelleknek* tekintik, mint ahogyan a például a felderítő statisztika egyes osztályozó módszereinél feltételezett eloszlásfüggvényekkel is tették.

(Megjegyzés: matematikai modellen a matematikusok mást értenek).

Az oksági viszonyok kifejezésére törekvő, fizikai, kémiai, biológiai tartalmú változókkal bíró modelleket az angol irodalom *hard model*-eknek nevezi, szembeállítva azokat a fizikai tartalommal nem bíró *soft model*-ekkel. Az utóbbin jobbjára azokat a modelleket értik, amelyeknél a mért változókat számított *látens* változókkal helyettesítették, pl. főkomponensekkel.

Ami az „adatstruktúrát“ illeti: egymással össze nem függő objektumok mért tulajdonság értékei (feltéve, hogy azok relevánsak) egyenletesen, „statisztikusan homogén“ módon töltik ki azt a tulajdonságtérét, amelyre a kísérletezés kiterjed. Ha a mérési pontok a térben csomósodnak, clusterekbe tömörülnek, bizonyos nemleges összefüggés már létezik, nevezetesen az, hogy objektumok különböznek egymástól, pontosabban tulajdonságaik várható értéke más és más. *Hipergömb* alakú pontfelhők arról árulkodnak, hogy a felhőbe tartozó objektumok tulajdonságai egymás között csak véletlenül különböznek, *hiperellipszoid* alakú (lencseszerű, szivaralakú) felhő eltérő szórású tulajdonságokra, ha pedig az ellipszoid a tulajdonságtér koordináta tengelyeihez képest el is fordul, az objektumok jellemzői valamiképpen korrelálnak. Természetesen ezek a „hiperalakok“ a síkban és a térben vetületeik alakjában szemlélhetők. Ha a vizsgált objektumok végül vonalakba, *síkokba*, *felületekbe* rendeződnek, felderíthető függvénykapcsolat sejthető.

1.2. A modellek jellemzői

A matematikai modell *függő* és *független* változók között állapít meg összefüggést. A továbbiakban egy független változót általában x -szel, vektorukat \mathbf{x} -szel, mátrixukat \mathbf{X} -szel jelöljük, a függő változókat y -nal, \mathbf{y} -nal vagy \mathbf{Y} -nal. A modell legáltalánosabban legyen

$$\mathbf{Y} = f(\mathbf{X}) \quad (1.1)$$

A matematikai modell x független változóit adott esetekben szokás *prediktoroknak* (predictors) vagy *magyarázó* (explanatory) változóknak, a függő változókat *válasznak* (*response*) is nevezni.

Megjegyzés: x azonosító nem szükségszerűen csak magukat a független változókat jelenti, jelölheti azok valamely függvényét (logaritmusát, hatványát, reciprokát is). A két lehetőséget azonban az egyszerűség kedvéért most a képletekben nem különböztetjük meg.

A matematikai modell alakja függvénysseregeket határoz meg, ahol a függvények paramétereikben különböznek. Egy-egy függvényt tehát paramétere is jellemeznek:

$$\mathbf{Y} = f(\mathbf{X}, \mathbf{A}) \quad (1.2)$$

ahol \mathbf{a} , \mathbf{a} , ill \mathbf{A} jelöli a függvény skaláris-, vektor- vagy mátrix paramétereit.

A megismerés első lépése az adott, *ismert* vagy beállított x *változók* mellett *mért* vagy megfigyelt y *változók* alapján a jelenséget legjobban leíró partikuláris modell meghatározása. Ez a tevékenység a modelfüggvény *paramétereinek* megállapítását jelenti. Mivel az y változók kísérleti eredmények, a paramétereket csak pontatlanul lehet meghatározni. Ezért nevezik a megismerésnek ezt a lépését *paraméterbecslésnek*. Bár a *regresszió* eljárása történetileg a korreláció számítás kapcsán, tehát lineáris esetek vizsgálatánál alakult ki, a regresszió elnevezést többé kevésbé a paraméterbecslés szinonimájaként használják. Hasonlóan szinonimája a paraméterbecslésnek a *modellillesztés* vagy *görbeillesztés* (model fitting, curve fitting).

A becsült (más szóhasználattal *számított*) paraméterek jele legyen \hat{a} , $\hat{\mathbf{a}}$, vagy $\hat{\mathbf{A}}$, aszerint hogy skaláris- vektor- vagy mátrixváltozóról van szó. Sok esetben csak a becsült paraméterek meghatározása a cél, természetesen *variánciájukkal* és *kovariánciájukkal* együtt. Utóbbi mennyiségeket $\text{var}(\hat{a})$ ill. $\text{cov}(\hat{\mathbf{a}})$ jellel azonosítják. Vektorba rendezett számított paraméterek variánciáit és kovariánciáit $\hat{\mathbf{C}}$ mátrix tartalmazza.

A becsült paramétereket (2) modellbe helyettesítve megkaphatók a *becsült, illesztett* vagy *számított* y értékek, amelyeket \hat{y} -nal jelölünk. Ezek variánciáját is kiszámítják.

Adott független változók i -edik értéke mellett mért és az illesztett paraméterekkel ugyanott számított y értékek

$$d_i = y_i - \hat{y}_i \quad (1.3)$$

különbségét *reziduális különbségnek* vagy röviden *reziduálisnak* szokás nevezni.

A számított \hat{y} értékek és az d reziduálisok kiszámítása már azért is fontos, mert ezekkel ítélhető meg mennyire *alkalmas (adekvát)* egy modell.

A reziduális különbségek ugyanis két részből állhatnak: *rendszeres eltérésből (bias)* és mérésekre visszavezethető *véletlen hibából*. Más szavakkal az i -edik mért érték a modell alapján várt értéken kívül b_i rendszeres eltérést és a valószínűségi változó jellegű rendszeres hiba e_i realizációját is tartalmazhatja:

$$y_i = f(\mathbf{X}_i, \mathbf{A}) + b_i + e_i \quad (1.4)$$

azaz az i -edik reziduális:

$$d_i = y_i - \hat{y}_i = (f(\mathbf{X}_i, \mathbf{A}) + b_i + e_i) - \hat{y}_i \quad (1.5)$$

Ha a reziduálisok becsült *várható értéke* (átlaga) nem tér el szignifikánsan 0-tól, akkor ez azt jelenti, hogy a rendszeres hiba és a véletlen hibák várható értéke 0, tehát

$$f(\mathbf{X}, \mathbf{A}) = \hat{y} \quad (1.6)$$

tehát *torzítatlan, alkalmas, adekvát*, rendszeres hiba (bias) nincs. Ha nem ez a helyzet, vagy kísérleti hibát kell keresni, vagy a vizsgált jelenséget nem a választott modell írja le.

Az illesztéssel kapott \hat{a} paraméterek és az azokkal számolt \hat{y} értékek egy adott kísérlet eredményei. Nem helyes emiatt az eredményekből levont következtetések kiterjesztése a jövő esetekre. Az eredmények jellemzik az *illesztés jóságát*, de nem jellemzik a modell *jóslási* (predikciós) képességét. A paraméterek akkor lesznek predikcióra alkalmasak, ha alkalmas módon kiszámítják azok *jóslási variánciáját*, a jövőben várható y értékek pedig akkor ismerhetők meg, ha az \hat{y} értékek mellé meg van adva azok *jóslási* szórása. Ezeknek a mennyiségeknek számítási módszereit *keresztellenőrzésnek* (cross

validation) is nevezik. A keresztellenőrzés –amelyről később lesz szó– lineáris modelleknél viszonylag egyszerű, ellenkező esetben általában kísérleti feladat.