

ALGORITMUS A LEGJOBB PREDIKTORHALMAZ KIVÁLASZTÁSÁRA

Balás Éltés András

Semmelweis Orvostudományi Egyetem Számítóközpont

A klinikai gyakorlatban ma már egyre több diagnosztikus feladat megoldásához létezik szokásos laborvizsgálati rend, un. protokoll. Így például magasvérnyomás, diabetes pontosabb diagnózisához, vagy chronikus vesebajos, művese, kezelt betegek állapotának periodikus ellenőrzéséhez. A jó protokoll korszerű sémát ad a szükséges és elégséges információtartalmu laboratóriumi vizsgálatok megrendeléséhez. A használatos vizsgálati protokollok, szokásos rendek jelentős része "házi" szabály, az adott kórház orvosainak tapasztalataira, véleményére épül, igazodik a helyi lehetőségekhez. Protokolljavaslatokat ugyanakkor nagy számban közöl a hazai és nemzetközi szakirodalom is.

Diagnosztikai protokollokkal szerzett tapasztalatok statisztikai elemzésekor gyakori észrevétel, hogy a vizsgálatok eredményei láthatóan összefüggnek, mégis két vizsgálat között olyan szoros kapcsolat ritkán bizonyítható, amely az egyik elhagyását a protokollból lehetővé teszi. Többváltozós lineáris kapcsolat is gyakran sejthető, de a célszerűen elhagyandó vizsgálatok kiválasztását ez még nem oldja meg.

Orvosi értelemben a feladat egy protokollban szereplő, szokásosan egyszerre megrendelt laboratóriumi vizsgálati halmazból a legköltségesebbek elhagyása, úgy, hogy az elhagyott vizsgálatok eredményét egy alkalmas matematikai formulával, kellő pontossággal kiszámíthassuk a protokollban meghagyott laboratóriumi vizsgálatok eredményeiből.

A feladat matematikai megfogalmazása: $\underline{X} \in \mathbb{R}^n$ valószínűségi vektorváltozó, \underline{c} a hozzárendelt költségvektor és \underline{a} D mátrix kijelöli a prediktorváltozókat:

$$D_{n \times k} = \begin{cases} d_{ij} = 1 & \text{ha az } i\text{-edik prediktor a } j\text{-edik változó} \\ d_{ij} = 0 & \text{egyébként} \end{cases}$$

ahol $k = \sum d_{ij}$. A cél D megválasztása úgy, hogy

$$a/ \quad M(\underline{X}_j) = \theta_j \cdot D \cdot M(\underline{X}) \quad \text{ha } \exists i, \text{ melyre } d_{ij} = 1$$

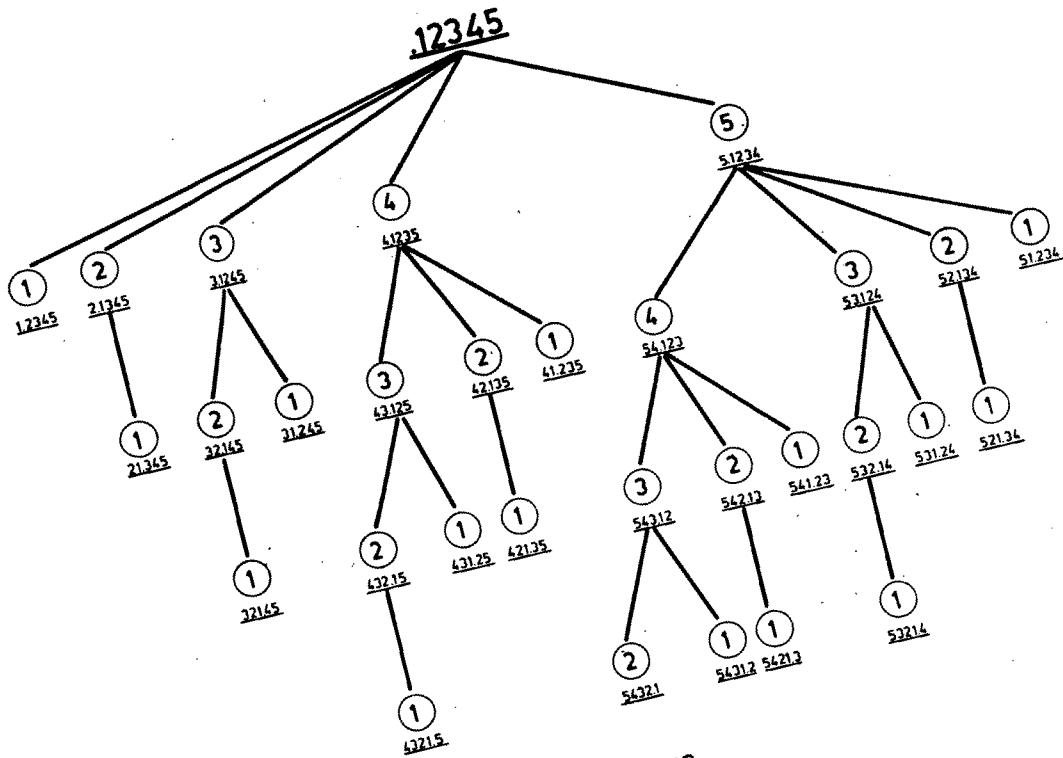
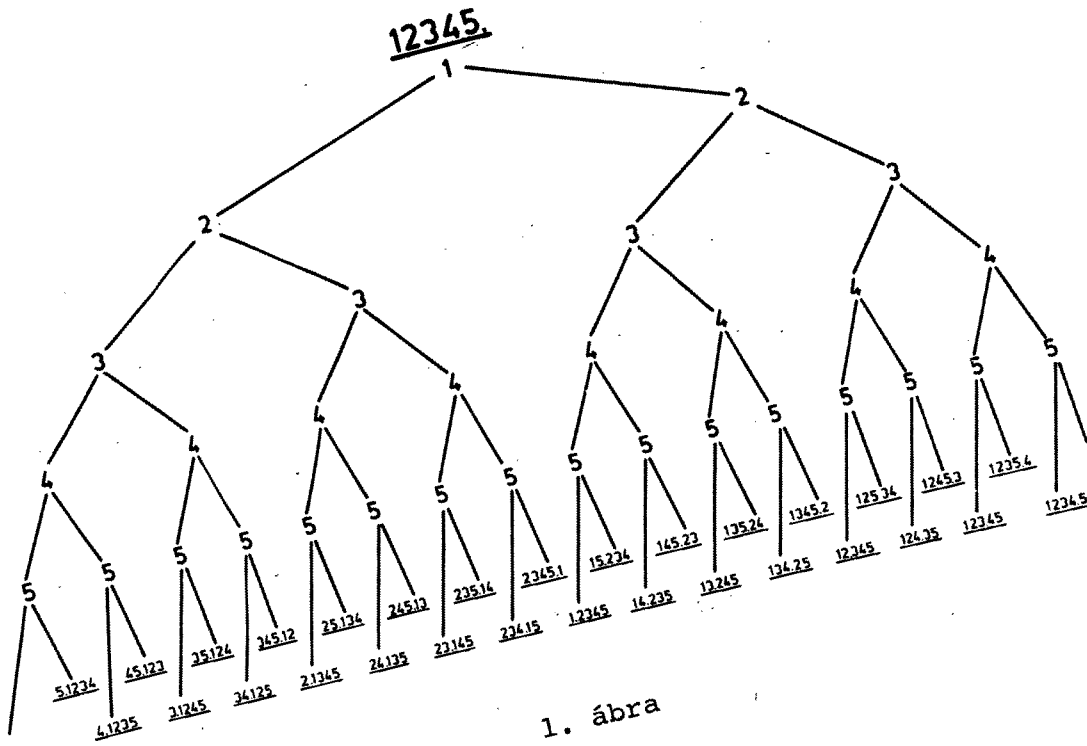
$$b/ \quad \min_i R_{X_j} \cdot D \cdot X \geq R_{\min} \\ d_{ij} = 0$$

$$c/ \quad \underline{l} = \underline{1} \cdot D \cdot \underline{c} \rightarrow \min$$

A θ_j regressziós paramétervektor meghatározása, az R többváltozós korreláció számítása a prediktor változók kovarianciamátrixának invertálását jelenti.

Az egyszerűnek látszó megközelítés, minden lehetséges prediktorhalmaz kipróbálása és így a legjobb kiválasztása rendkívül számításigényes, p változó esetén 2^p invertálást jelent.

Az 1. ábra módszert, fastrukturát ad a prediktorhalmazok megha-



tározásához, a többváltozós analízisben használatos "pont" jelöléssel /a pont mögött a prediktorok, előtte a becsült változók, a prediktandusok/.

A számítást célszerűen lépésenként, azaz változónként sweepeléssel, vagy másnéven pivotálással hajthatjuk végre. A sweep operátor a sweepelt változók kovarianciamátrixának inverzét, a többi változókra pedig a reziduális négyzetösszeget adja. /Ez utóbbit azzal a feltevéssel számítja, hogy a sweepelt változók a prediktorok./ Az operátor invertálható, ezért a továbbiakban be-, illetve kisweepelésről beszélünk, aszerint, hogy az eredeti vagy az inverz leképezésről van szó. A sweepelés azonban különösen nagyobb mátrixok esetén rendkívül számításgényes.

Az ábrán látható fán a szükséges besweepelések száma:

$$n = p \cdot (2^{p-1} - 1)$$

A szükséges sweepelések száma majdnem a felére csökkenthető, ha a fastruktúrát átrendezzük, úgy hogy minden sweepelés értékelhető prediktor részhalmazhoz vezessen. A 2. ábra példája mutatja, hogy a számítás ekkor valamennyi változó besweepelésével indul, majd az ábrán vázolt sorrendben kisweepelünk. A szükséges sweepelések száma:

$$n = (p+1) \cdot 2^{p-2} + p - 1$$

A sweepelések száma tovább csökkenthető, a fa alkalmas átrendezésével. Ennek szabályai:

- a változók költség szerint monoton csökkenő sorrendben vannak /1-es számú a legdrágább/,

- az egyes változók első kisweepelését jelentő főtörzsek alá a tőle balra eső főtörzsek fordított sorrendben kerülnek.

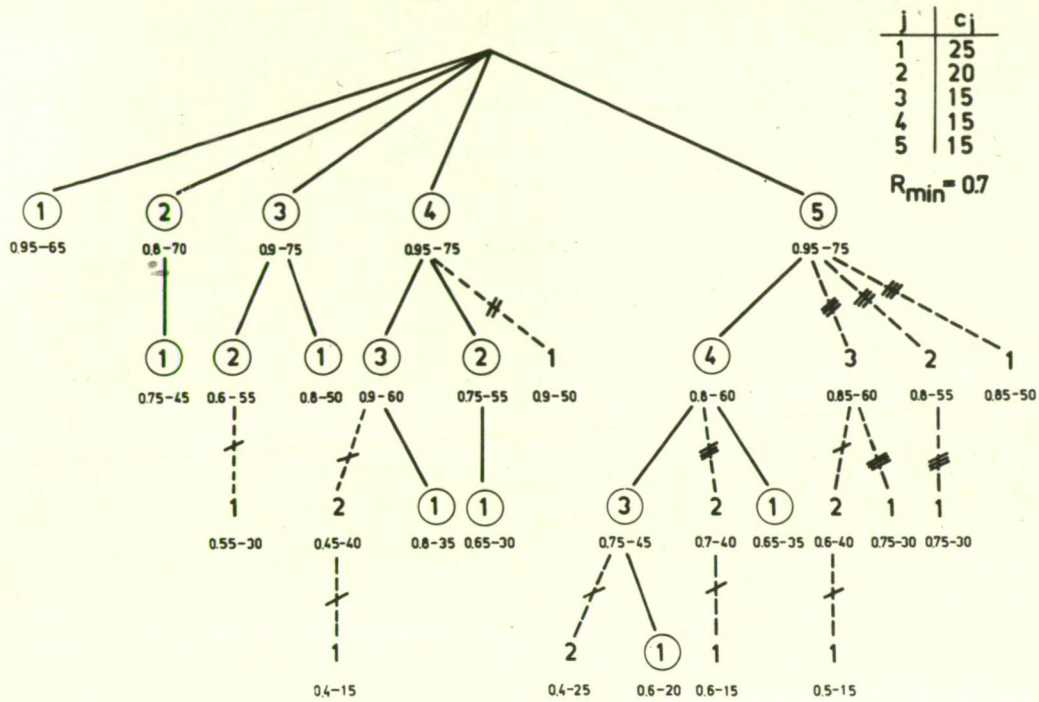
A 3-as ábra 5-változós számpéldával illusztrálja, hogyan hagyhatók el további sweepelések ebben az elrendezésben. /Sweepelni csak a bekarikázott csomópontokban kell./ A csomópontok alatt az adott prediktorhalmazokhoz tartozó minimális többváltozós korreláció és a költség látható. Felülről lefelé haladva az ágakon a vizsgálati költség és a minimális korreláció monoton csökken. Kisweepelés elhagyható:

- ha a minimális többváltozós korreláció az adott limit alá esik; a prediktorváltozók számának csökkentésével ez nem javulhat, így az ág hátralévő része elhagyható a továbbiakban / * jellel/,

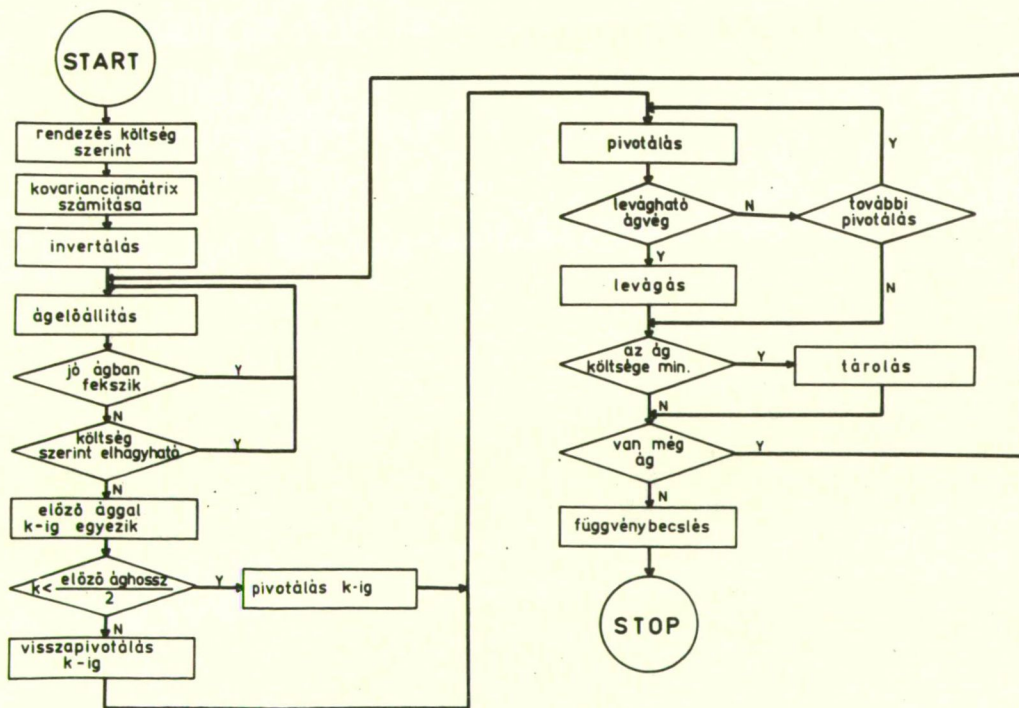
- ha egy ág sweepelési csomópontjai már mind szerepeltek egy korábbi jó ágban, akkor számítása az adott főtörzs alatt elhagyható, mert a pontossági kritériumnak biztosan megfelel, de költsége nagyobb, vagy egyenlő / * jellel/,

- ha egy ágon végighaladva a költség sem az adott főtörzs alatt, sem később nem lesz jobb, mint az addigi minimum, akkor számítása elhagyható / * jellel/.

A 4. ábra az algoritmust foglalja össze. Látható, hogy további megtakarítást tesz lehetővé a sweepelések számában, ha figyelembe vesszük, hogy az új ág a réginek oldalága-e. Ha igen, a csomóponti elágazásig, ahol az új ág indul, az előző ágon besweepelve mehetünk vissza és állíthatjuk elő az oldalág számításának alapjául szolgáló mátrixot.



3. ábra



4. ábra

Az algoritmus alapján FORTRAN nyelvű interaktív programot készítettem HwB 66/20-as gépre, amely outputján megadja az így kiválasztott prediktor subsetét, költségét, a százalékos megtakarítást és az elhagyható változók /betegvizsgálatok/ becslésére szolgáló többváltozós lineáris regressziós egyenleteket.

Megjegyzendő, hogy a feladat megoldására a Furnival és Wilson által 1974-ben leírt "leaps and bounds" algoritmus, illetve a stepwise regressziós módszer /amelyek több statisztikai programcsomagban megtalálhatók, így többek között az BMDP-ben is/ nem alkalmazsak, mert a jóslandó változókat előre rögzíteni kell és a prediktor változók kiválasztásában a költség nem játszik szerepet.