

**NÉHÁNY STATISZTIKAI MÓDSZER
A METEOROLÓGIÁBAN**

Egyetemi doktori értekezés

Faragó Tibor

1977

Tartalomjegyzék

Bevezető	1
1. „Legközelebbi társ” /NN/ becslési eljárás és a Bayes-döntés átlagos veszteségének becslése diszkrét mintatér és folytonos paramétertér esetén	5
2. Adatszugorítás a statisztikai becslésméletben	12
3. Egy „szelektív legközelebbi társ” statisztikai becslési eljárás empirikus vizsgálata	21
4. Maximumértékek eloszlásbecslésének időegység – problémája	25
Zárómegjegyzések	34
Mellékletek	
1. és 2. melléklet: Az 1.3. pontban leírt NN-eljárást realizáló Fortran-program protolja és egy futási eredménye	36
3. melléklet: A 2.3. pontban leírt TSH program protokolja és egy futási eredménye	42
4. melléklet: A 2.3. pontban vázolt HST program protokolja és egy futási eredménye	50
5. melléklet: A 2.3. pontban vázolt KL1 program protokolja és egy konkrét TSH programmal redukált mintasorra való alkalmazása	55
6. melléklet: A 3.2 pontban leírt szelektív NN-eljárást megvalósító program, mely felváltva két gaussi változó realizációit állítja elő	62
Hivatkozások	80

BEVEZETŐ

Alig 50 éve, hogy egy "álmodozó" meteorológus, Richardson felvetette egy "emberi számológépekből" álló "időjárás-gyár" gondolatát, melyben több tízezer ember alkotta volna a mai számítógépek aritmetikáját és operatív memóriáját. Alig 30 év telt el az első számítógép szerkesztése óta, melyet két meteorológussal, Charney-val és Fjörtofttal karöltve a magyar származású von Neumann éppen egy időjárás-előrejelzési modellen próbált ki. Ma pedig már szinte elképzelhetetlennek tűnik, hogy a meteorológiai kutatási és operatív tevékenységek döntő többsége ne éppen számítógépre tervezett programok illetve rendszerek formájában realizálódjon. Különösen igaz ez a megállapítás a statisztikai kutatásokra, melyek például nagyobb skálájú meteorológiai mezők elemzése esetén már meglehetősen memóriaigényesek. Éppen ezért az általánosabb statisztikai vizsgálatok mellett nem egy vizsgálatban fontos szerep jut a többé-kevésbé optimális adattárolásnak, valamint a fizikailag és statisztikailag is alátámasztott adatszugorítási eljárásoknak.

A meteorológia területén a statisztikai módszerek felhasználásának két alapvető területét különíthetjük el. Az első terület a klimatológia, melynek keretében megadott meteorológiai elemek vagy azok függvényei átlagos viselkedésére nézve kívánunk következtetéseket levonni. E téren gyakran találkozhatunk a szakirodalomban meteorológiai állapothatározók rövidebb-hosszabb távú átlagainak statisztikai-fizikai elemzésével, eloszlásbecslésekkel és -illesztésekkel, e becslések torzítottságának vizsgálatával, statisztikai jellegű osztálybasorolásokkal stb. Itt csak egyetlen összefoglaló munkára hivatkozunk ebből a témakörből, amely a különféle meteorológiai mezők statisztikai szerkezetéről ad áttekintést Czelnai, Gandin és Zaharjev szerkesztésében /1976/. A statisztikai prognosztika területe sokszor nem különíthető el élesen a klimatológiai feldolgozások technikájától, a kutatások célja azonban alapvetően más. E téren az egyik legújabb eljárásként a klasszikus dinamikai /hidro- és termodinamikai/ modellekkel

szoros kapcsolatban álló sztochasztikus-dinamikus modelleket kell megemlíteni. Freiburger és Grenander /1965/ ötletadó munkája nyomán több szerző foglalkozott a dinamikus rendszerek ilyen általánosításával, így többek között Gleeson/1966/, Epstein/1969/, Tatarszkij/1969/, Fleming/1971/, Szonyecskin et al./1976/ és Faragó/1977/. Egy másik kutatási ág keretében a vizsgált fizikai folyamatot valamilyen speciális sztochasztikus folyamattal modellezik, melyet statisztikai eljárások segítségével körvonalaznak és utóbb paramétereznek. Többen is kísérleteket végeztek például a csapadékos és nem csapadékos napok sorozatának Markov-folyamattal való szimulálására, illetőleg más vonatkozásban a stacionárius és a Gauss-folyamatok elméletének alkalmazására. Ide sorolható a sztochasztikus automaták elméletének felhasználása is. Ez utóbbival kapcsolatban Kapovicsné-Maller-Titkos /1976/ dolgozatára utalunk. Korábban magam is foglalkoztam ennek az elméletnek prognosztikai alkalmazhatóságával /Faragó, 1973, 1974 /. Igen átfogó, elsősorban szovjet irodalommal rendelkezik a sztochasztikus lineáris elmélet meteorológiai tárgyú felhasználása. Ilyen jellegű eljárásnak tekinthetjük a meteorológus berkekben "természetes empirikus ortogonális sorfejtésnek" nevezett Karhunen-Loève sorfejtés alkalmazását, a kanonikus korreláció becslését vagy a különféle regressziós modelleket /Bagrov, 1966; Thom, 1966; WMO, 1966; Gandin, 1967; Craddock-Colgate, 1974 stb/. A regresszió típusú extrapolációval már korábban foglalkoztunk általános és meteorológiai alkalmazási megközelítésben /Faragó-Gulyás, 1974/. A teljesség igénye nélkül itt végül megemlítjük még a sűrűségfüggvények illetve feltételes sűrűségfüggvények közvetlen vagy közvetett becslését. Ezek vizsgálatának problémája kétirányú: az esetek túlnyomó többségben egyrészt nem áll rendelkezésre a jó hatásfoku becsléshez szükséges kellő mennyiségű adat, másrészt viszont bizonyos feladatoknál szinte egyedül ez a közelítés látszik eredményesnek. E kutatásokhoz sorolható Bagrov diszkriminancia-analízisre támaszkodó, a havi csapadékösszeg előrejelezhetőségét vizsgáló modellje, vagy például a svéd Nyberg /1976/ kontingenciákat alkalmazó modellje.

A becslési és illeszkedési eljárásokkal foglalkozott egy korábbi, konkrét számítástechnikai vonatkozásokat is felölelő tanulmányunk /Faragó-Gulyás,1975/. A "közvetett" becslések csoportjának egyik legegyszerűbb, de a meteorológiai prognosztikában a leggyakrabban hasznosított típusa a "legközelebbi társ" eljárások köre. Ezzel kapcsolatban többek között Craddock/1958/vagy Koppány/1974/ dolgozatára hivatkozhatunk. Az alapelveket és néhány operatív modellt tárgyal Faragó-Kaba-Gulyás/1975/.

A Központi Meteorológiai Intézetben távprognosztikai célokra egy olyan k-NN / k legközelebbi társat felhasználó, azaz k-Nearest Neighbour/ alapú statisztikai modell került kidolgozásra /Faragó et al.,1975; Faragó-Kaba,1976 /, melyben prediktorként egy speciális /diszkrét/ mintatér elemeit alkalmaztuk. Mindenekelőtt a modell "előrejelző képességére" voltunk kíváncsiak, amit az empirikus Bayes-rizikó közvetett becslésével jellemeztünk. A becslések elvégzésére Wagner /1971/ egy tételének általánosítása adott módot. E témával foglalkozik a dolgozat 1. fejezete.

A nagydimenziószámú statisztikai becsléseknél jelentős problémát jelent az előzetes adatszűrités. Ezzel kapcsolatban lényeges kérdés, hogy ez a transzformáció miként változtatja meg a kérdéses statisztikai becslési feladatot: folytonosan viselkedik-e az optimális /bayesi/ becslésen a kitűzött veszteségfüggvény? Korábban speciálisan a sorfejtéses adatredukciókkal foglalkoztunk /Faragó-Gulyás,1973a,1973b/, utóbb pedig a vizsgálat tárgya volt a becslés folytonossága /Faragó,1975a/, valamint általánosabb vetületben az ún. hiba-torzítás függvény viselkedése /Faragó-Györfi,1974,1975/.

A sorfejtéses adatredukció alkalmazása jelenleg egy multi-regressziós extrapolációs feladat kapcsán merült fel. Míg egy korábbi dolgozatban /Faragó,1975b/ egydimenziós mintapéldán a regressziós becslés javításának olyan lehetőségét mutatuk be, amelyik a mintatér particionálására épül, addig ebben a feladatban a nagydimenziószámú minták előzetes adatredukciót tesznek szükségessé. Az optimálisnak számító Karhunen-Loève sorfejtés azonban ismét csak a nagy dimenziószám miatt nem végezhető el közvetlenül és így még egy ezt is megelőző - spektrális - adatredukciót alkalmaztunk.

A 2. fejezetben bemutatjuk a sorfejtéses adatszűrés technikákkal összefüggő elméleti vizsgálatokat, a lineáris multiregresszió alapjait, valamint az ezeket felhasználó meteorológiai modell számítástechnikai tervét. Végül bemutatunk a szuperponált kódolással kapcsolatos néhány konkrét eredményt.

Mint már fentebb utaltunk rá, az "analógiás" / NN / eljárások igen elterjedtek a meteorológiai alkalmazásokban. Osztályozási szabályok statisztikai tanulására, azaz "alakfelismerésre" számos módszer ismeretes a statisztikai szakirodalomból és ezek sorában is az NN-módszerek többféle verziója található meg. Nagyszámú osztály /hipotézis/ és nagyobb dimenziószám mellett általában más algoritmusokkal nem is próbálkozhatunk. Bizonyos, a légköri mozgásrendszerek strukturáját leíró ún. makroszinoptikus kódok felismerésével kapcsolatban éppen ilyen probléma volt adott. Ennek megfelelően memória- és művelettakarékos szelektív NN-eljárás került kidolgozásra. Ennek egy empirikus vizsgálatát mutatja be a 3. fejezet.

A statisztikai kutatások sorában végül maximumértékek becslésével kapcsolatos vizsgálatot mutatunk be /Farágó, 1977/. A vizsgálat célja annak eldöntése, hogy a statisztika alapjául milyen időegységre adjuk meg a maximumértékeket adott mintaszám esetén. A kérdés konzisztencia-vizsgálattal dönthető el. Elemezzük a gyenge függőség és a paraméteres statisztika esetét is. Az ezt tárgyaló 4. fejezet egy konkrét analizissel zárul.

1. "LEGGÖZELEBBI TÁRS" /NN/ BECSLÉSI ELJÁRÁS ÉS A BAYES-DÖNTÉS ÁTLAGOS VESZTESÉGÉNEK BECSLÉSE DISZKRÉT MINTATÉR ÉS FOLYTONOS PARAMÉTERTÉR ESETÉN

1.1. Egy meteorológiai feladat

Hess és Brezowsky /1952/ makroszinoptikus kódok olyan rendszerét dolgozta ki, melyekkel napról-napra többé-kevésbé jól jellemezhető az Európai Szinoptikai Körzet /ESzK/ fölötti levegő troposzférába eső részének áramlási strukturája, a markáns légnyomásképződmények elhelyezkedése és közvetve e régió időjárása. E katalógus 30 elemet tartalmaz. Ha egy adott időszak minden napjához hozzárendeljük a megfelelő kódot /a továbbiakban HB-kódot/, akkor az így nyert sorozat az ESzK fölötti légcirkuláció menetét írja le nagy vonalakban. E kódok hosszú időre visszanyúló archivuma létezik számítógépes adathordozón, amely lehetővé teszi, hogy pl. egy adott hónap x HB-kódsorozatához hasonlókat keressünk a multból prognosztikai céllal. Két tetszőleges HB-kód összehasonlítására adott egy S nyereségmatrix /"scoring"-táblázat/. Ennek segítségével definiálható egy $\Delta(x, \tilde{x})$ hasonlósági mérőszám a kódsorozatok között, amelyik kvázimetrika:

$$\Delta(x, \tilde{x}) \geq 0$$

$$\Delta(x, \tilde{x}) = 0 \iff x = \tilde{x}$$

$$\Delta \text{ - szimmetrikus}$$

Most azt a feladatot tűzzük ki, hogy "analógiás" eljárással megbecsüljük az x_0 aktuális mintaelemet /megfigyelést/ követő hónap \mathcal{V}_0 budapesti középhőmérsékletét. Az analógia-keresés a $T_N = \{(x_k, \mathcal{V}_k)\}_{k=1}^N$ adott megfigyelési sorozatra /archivumra/ épül, ahol \mathcal{V}_k az azt a hónapot követő budapesti havi középhőmérsékleti érték, amelyekre x_k vonatkozik.

1.2. A probléma matematikai megközelítése

Legyen adott az (x_0, \mathcal{V}_0) valószínűségi változópár,

$x_0 \in X$ - metrikus mintatér Δ metrikával, $\mathcal{V}_0 \in \mathbb{H}$ a paraméterter. $T_N = \{ (x_k, \mathcal{V}_k) \}_{k=1}^N$ független, egyforma eloszlású megfigyelések sorozatát jelöli az (x_0, \mathcal{V}_0) párra. Döntésfüggvénynek nevezzük a mintatér mérhető leképezését a paraméterterre:

$$d: X \longrightarrow \mathbb{H} \\ d(x_0) = \hat{\mathcal{V}},$$

ahol $\hat{\mathcal{V}}$ -t a \mathcal{V} -nek a d döntésfüggvénnyel származtatott becslésének tekintjük. E becslés jóságának mérésére bevezetjük a

$$/1/ \quad W(\mathcal{V}_0, \hat{\mathcal{V}}) = (\mathcal{V}_0 - \hat{\mathcal{V}})^2$$

veszteségfüggvényt, itt már kikötve, hogy $\mathbb{H} \subset R_1$.

Az átlagos veszteség:

$$/2/ \quad R(d) = E W(\mathcal{V}_0, d(x_0))$$

A bayesi döntésfüggvény definíciószerűen a legkisebb átlagos veszteségű döntés az összes döntésfüggvények közül:

$$d^* : R^* = R(d^*) \leq R(d) \quad \forall d$$

Ismeretes, /Rényi, 1970/, hogy ebben az esetben a bayesi döntésfüggvény éppen a feltételes várhatóérték /ha az létezik/

$$/3/ \quad d^*(x) = E(\mathcal{V} / x)$$

A feladat minél optimálisabb megoldása, azaz \mathcal{V}_0 minél pontosabb becslése a Bayes-döntés közvetlen közelítésével érhető el. Ez azonban általában nem kivitelezhető, mert bár létezik nem egy, például szekvenciális eljárás többek között a megfelelő feltételes sűrűségfüggvény becslésére, melyek aszimptotikusan kedvező tulajdonságokkal rendelkeznek, de a gyakorlatban a feladatok többsége viszonylagosan kismintás.

A "legközelebbi társ" /a továbbiakban NN/ eljárást a következő döntésfüggvény írja le:

$$d_N: d_N(x) = \mathcal{V}_k,$$

ahol \mathcal{V}_k annak az x_k mintának a paramétere, melyre

$$\begin{aligned} \Delta(x, x_k) &\leq \Delta(x, x_{\mathcal{P}}) & \mathcal{P}=1, 2, \dots, N, \\ \Delta(x, x_k) &< \Delta(x, x_{\mathcal{L}}) & \mathcal{L} < k. \end{aligned}$$

A továbbiakban a

$$d_N(x_0) = \mathcal{V}_{0(N)} \quad ; \quad R_N = R(d_N)$$

jelöléseket fogjuk alkalmazni. Feltéve, hogy $N \rightarrow \infty$ esetén létezik az NN-döntés átlagos veszteségének határértéke, jelölje $R = \lim_{N \rightarrow \infty} R_N$ az NN-döntés aszimptotikus veszteségét.

"Nulla-egy" veszteségfüggvényre

$$/3a/ \quad W(\mathcal{V}_0, \hat{\mathcal{V}}) = \begin{cases} 0, & \text{ha } \mathcal{V}_0 = \hat{\mathcal{V}} \\ 1, & \text{egyébként} \end{cases},$$

amikor is $R(d)$ a $d(x) = \hat{\mathcal{V}}$ döntés hibavalószínűségét jelent, Wagner/1971/ bebizonyította, hogy az

$$/4/ \quad S_N = \frac{1}{N} \sum_{k=1}^N W(\mathcal{V}_k, \mathcal{V}_{k(N)})$$

statisztika aszimptotikusan torzítatlan becslése R -nek, sőt $S_N \rightarrow R$ sztochasztikusan. Itt $\mathcal{V}_{k(N)}$ a \mathcal{V}_k NN-becslése a $T_N = \{(x_k, \mathcal{V}_k)\}$ mintasorozatból. Márpedig ha eképpen R közelíthető egy véges minta segítségével, akkor Cover/1969/ alapján az adott modell előrejelzőképességét kifejező R^* bayesi átlagos veszteség is becsülhető:

$$R^* \in \left[\left(1 - \sqrt{1 - 2R_N}\right) \cdot \frac{1}{2}, R_N \right].$$

Kevésbé kedvező a helyzet az /1/ veszteségfüggvény esetén: $R^* = 2R$.

Ahhoz tehát, hogy az /1/ négyzetes veszteségfüggvény esetén R^* -t megbecsülhessük, szükséges a Wagner-tétel általánosítása. Ehhez mindenekelőtt az NN-módszerek alapvető konvergenciatételére hivatkozunk.

1. Tétel /Cover-Hart, 1967/: Ha az X mintatér szeparábilis és

metrikus a Δ metrikával és $x_1, x_2, \dots, x_N, x_0$ független, egyforma eloszlású valószínűségi változósorozat, akkor $x_{0(N)} \rightarrow x_0$ 1 valószínűséggel.

Ez a tétel igaz marad olyan diszkrét mintatérre is, melyen csak egy Δ kvázimetrika adott. Most már megfogalmazható az általánosított Wagner-tétel. /Faragó-Kaba, 1976/.

2. Tétel: Legyen adott az X mintatér olyan Δ /metrikával vagy/ kvázimetrikával, hogy a független, egyforma eloszlású $x_1, x_2, \dots, x_N, x_0$ mintákra teljesüljön az $x_{0(N)} \rightarrow x_0$ /m.m./ konvergencia. Feltéve, hogy a Θ paraméterterv korlátos, valamint $d^*(x) = E(\mathcal{V}_0 / x)$ és $\sigma(x) = E(\mathcal{V}_0^2 / x)$ feltételes várható értékek folytonosak /m.m./, igaz az $S_N \rightarrow R$ négyzetes középben vett konvergencia.

Bizonyítás.

Írjuk fel a megfelelő négyzetes középben vett eltérést és bontsuk fel a zárójelet:

$$/5/ \quad E(S_N - R)^2 = E(S_N^2) - 2R E(S_N) + R^2 = E(S_N^2) - 2RR_{N-1} + R^2,$$

ahol kihasználtuk, hogy S_N torzítatlan becslés R_{N-1} -re:

$$E(S_N) = \frac{1}{N} \sum_{k=1}^N E(\mathcal{V}_k - \mathcal{V}_{k(N)})^2 = \frac{1}{N} \sum_{k=1}^N R_{N-1} = R_{N-1}.$$

Cover/1968b/ nyomán ismeretes, hogy $R_N \rightarrow R$, következésképpen most azt kell kimutatnunk, hogy $E(S_N^2) \rightarrow R^2$. Itt

$$\begin{aligned} E(S_N^2) &= \frac{1}{N^2} \sum_{k, l} E\{(\mathcal{V}_k - \mathcal{V}_{k(N)})^2 (\mathcal{V}_l - \mathcal{V}_{l(N)})^2\} = \\ &= \frac{1}{N^2} \sum_{k=1}^N E(\mathcal{V}_k - \mathcal{V}_{k(N)})^4 + \frac{1}{N^2} \sum_{k \neq l} E\{(\mathcal{V}_k - \mathcal{V}_{k(N)})^2 (\mathcal{V}_l - \mathcal{V}_{l(N)})^2\} = \\ &= \frac{1}{N} E(\mathcal{V}_k - \mathcal{V}_{k(N)})^4 + \frac{N-1}{N} E\{(\mathcal{V}_k - \mathcal{V}_{k(N)})^2 (\mathcal{V}_l - \mathcal{V}_{l(N)})^2\}, \end{aligned}$$

ahol az utóbbi kifejezés az $(x_1, \mathcal{V}_1), (x_2, \mathcal{V}_2), \dots, (x_N, \mathcal{V}_N)$ párok függetlensége és egyforma eloszlása miatt tetszőlegesen

k és l indexekre értendő. A jobboldal első tagja

⊙ korlátossága miatt tart nullához, tehát ha igazolást nyer, hogy

$$/6/ \quad E (\mathcal{V}_k - \mathcal{V}_{k(N)})^2 (\mathcal{V}_1 - \mathcal{V}_{1(N)})^2 \xrightarrow[N \rightarrow \infty]{} R^2 ,$$

akkor $E \mathcal{S}_N^2 \rightarrow R^2$ is teljesülni fog. /6/ két tényezője nem független, ezért bevezetjük \mathcal{V}_k és \mathcal{V}_1 NN-becslését a $T_N = \{ (x_k, \mathcal{V}_k) , (x_1, \mathcal{V}_1) \}$ mintasorozatból, melyet $\tilde{\mathcal{V}}_k$ ill. $\tilde{\mathcal{V}}_1$ -lel jelölünk. Az így nyert két becslés aszimptótikusan felcserélhető, ugyanis a

$$K = | (\mathcal{V}_k - \mathcal{V}_{k(N)})^2 (\mathcal{V}_1 - \mathcal{V}_{1(N)})^2 - (\mathcal{V}_k - \tilde{\mathcal{V}}_{k(N)})^2 (\mathcal{V}_1 - \tilde{\mathcal{V}}_{1(N)})^2 |$$

jelöléssel K csak az $A = \{ \mathcal{V}_{k(N)} = \mathcal{V}_1 \cup \mathcal{V}_{1(N)} = \mathcal{V}_k \}$ eseményen különbözhet nullától, tehát

$$E(K) = E(K/A) P(A) \leq C \frac{2}{N-1} \xrightarrow[N \rightarrow \infty]{} 0 ,$$

ahol a feltételes várhatóérték ⊙ korlátossága miatt korlátos.

Tekintsük most /6/ baloldalát a módosított NN-becslésekkel - átalakítva :

$$\begin{aligned} /7/ \quad E \left[E \{ (\mathcal{V}_k - \tilde{\mathcal{V}}_{k(N)})^2 (\mathcal{V}_1 - \tilde{\mathcal{V}}_{1(N)})^2 / x_k, \tilde{x}_{k(N)}, x_1, \tilde{x}_{1(N)} \} \right] &= \\ &= E \left[E \{ (\mathcal{V}_k - \tilde{\mathcal{V}}_{k(N)})^2 / x_k, \tilde{x}_{k(N)} \} E \{ (\mathcal{V}_1 - \tilde{\mathcal{V}}_{1(N)})^2 / x_1, \tilde{x}_{1(N)} \} \right] = \\ &= E \left[L(x_k) \cdot L(x_1) \right] , \end{aligned}$$

ahol

$$L(x_k) = \sigma(x_k) - 2d^*(x_k) \cdot d^*(\tilde{x}_{k(N)}) + \sigma(\tilde{x}_{k(N)}) .$$

A Cover-Hart - tétel alapján $\tilde{x}_{k(N)} \rightarrow x_k$, $\tilde{x}_{1(N)} \rightarrow x_1$ /m.m./, tehát d^* és σ folytonossága miatt

$$\begin{aligned} L(x_k) &\rightarrow 2 \sigma(x_k) - 2 [d^*(x_k)]^2 = \\ &= 2 E \{ (\mathcal{V}_k - d^*(x_k))^2 / x_k \} \equiv 2r^*(x_k) \end{aligned}$$

és hasonlóképpen

$$L(x_1) \rightarrow 2 E \left\{ (\mathcal{V}_1 - d^*(x_1))^2 / x_1 \right\} \equiv 2r^*(x_1) \quad .$$

Mivel a paraméterterület korlátos, alkalmazható a Lebesgue-tétel és így a /7/ várhatóérték, azaz /6/ baloldala megfelelőképpen konvergál a

$$4 E \left\{ r^*(x_k) r^*(x_1) \right\} = 4 E r^*(x_k) E r^*(x_1) = 4 (R^*)^2$$

értékhez. /Itt kihasználtuk x_k és x_1 függetlenségét./ Cover/1968b/ tétele alapján $R=2R^*$ tehát valóban fennáll a /6/ konvergencia és ezzel a tételt igazoltuk. \square

E tétel lehetőséget nyújt R^* becslésére, bár a konvergencia sebessége ismeretlen. Csupán végesdimenziós /diszkrét/

⊕ -ra lehettek fel ezzel kapcsolatos eredmények pl.

Cover/1968a/ és Fritz/1975/ munkáiban.

Ha R^* közelítése kellő alapot nyújt a $\hat{\mathcal{V}}$ becslés pontosítására, akkor az NN-eljárásról érdemes áttérni a k-NN - eljárásra. Cover/1968b/ vizsgálatai alapján ugyanis ennek a szimptótikus átlagos vesztesége

$$R^{(k)} = \left(1 + \frac{1}{k}\right) R^* \quad .$$

Ebben az esetben a becslés \mathcal{V}_0 -ra:

$$\hat{\mathcal{V}} = \frac{1}{k} \sum_{j=1}^k \mathcal{V}_{o(N)}^{(j)} \quad ,$$

ahol $(x_{o(N)}^{(j)}, \mathcal{V}_{o(N)}^{(j)}) \in T_N$ és $x_{o(N)}^{(1)}, x_{o(N)}^{(2)}, \dots, x_{o(N)}^{(k)}$ az x_0 minta k legközelebbi társa.

1.3. A meteorológiai feladat megoldása

Olyan programot készítettünk, mely az 1.1 pontban vázolt meteorológiai feladatban egy aktuális x_0 kód-sorozathoz független, egyforma eloszlásúnak tekinthető T_N megfigyelési sorozatból megadott számú analógiát keres ki. N=93 esetén megbecsültük az S_N empirikus átlagos veszteséget akkor, amikor x_0 a január havi HB-kódsorozatot, \mathcal{V}_0

pedig a február havi budapesti átlaghőmérsékletet jellemezte: $S_N=11,7$. Ennek alapján $R^* \approx 5,9$. \mathcal{V}_0 szórásnégyzete

$\mathcal{G}_\mathcal{V} = 9,0$, tehát a Bayes-döntés már jelentékeny információnyereséget nyújthat. A többé-kevésbé optimális $k=9$ esetére $S_N = 8,4$, ami számottevő javulást jelent. /Rögzített N mellett persze k növekedésével az S_N becslés hatékonysága romlik./ Mindamellett $\mathcal{G}_\mathcal{V}$ -val összevetve a kapott S_N még túl nagy, de ezzel a k -NN -technikával csak úgy érhetnénk el jobb eredményt, ha k -t tovább növelhetnénk. Ennek viszont gátat szab a rögzített N érték.

Az 1.,2. mellékletek az "analógia-kereső" program egy olyan újabb variánsát mutatják be, ahol az adott x_0 -hoz már nemcsak a megelőző kerek naptári hónapokból keresünk analógiákat, hanem általánosabban x_k minta lehet a naptári hónapot csak részben fedő hónapnyi hosszúságú időszak kódsorozata is. Ezzel a fogással tovább javult a becslések minősége, a numerikus értékek további konkrét bemutatásától itt azonban eltekintünk.

2. ADATZSUGORITÁS A STATISZTIKAI BECSLÉSELMÉLETBEN

2.1. A lineáris regresszióról a meteorológiában és egy előrejelzési feladat

A dinamikus meteorológia elméletének legnagyobb vívmánya a légkör hidro- és termodinamikáját leíró un. primitív egyenletek rendszerének megalkotása. A nagyteljesítményű számítógépek megjelenése - mint már a bevezetőben említettük - lehetővé tette, hogy e primitív egyenletek segítségével és nemcsak egy elméletileg erősen szűrt egyenletrendszerrel numerikusan is modellezni lehessen a légköri folyamatokat. A numerikus modellezésnél akár differencia-módszereket akár spektrális reprezentációt alkalmaznak, lényegében az egyenletrendszer analitikus megoldásának közvetlen approximációjáról van szó.

Korábban még maga von Neumann /Smagorinsky, 1969 / három csoportba osztályozta a légköri jelenségek előrejelzéseit. Az azóta eltelt időszak tapasztalatai ezt igazolták: alapjaiban eltérő technikával kell közelítenünk a rövid /legfeljebb kb. 1-2 hétre szóló/, a hosszútávú /1-12 hónapos/ és az ultrahosszútávú előrejelzéseket. Míg a rövidtávú előrejelzéseknél uralkodóbb a dinamikai alapú modellezés, addig a hosszútávúaknál szinte kizárólagosan statisztikai eljárásokkal találkozhatunk a szakirodalomban; végül az ultrahosszú távú vizsgálatok ismét csak dinamikai /a termo- hidrodinamikai egyenletrendszer megoldására épülő/ módszerekre támaszkodhatnak. A statisztikai eljárások is akkor lehetnek jobb hatásfokúak ha pl. az egyszerű statisztikai analógiakeresésben kifejezett implicit prediktor-prediktandusz reláción túlmenően a dinamikai eljárásokhoz hasonlóan minél közvetlenebbül "analitikusan" próbálják közelíteni a megoldást - az említett reláció paraméteres alakjának minél jobb megadásával. A legegyszerűbb esetben lineáris multiregressziót választanak e célra és nem ritkán a vele szoros kapcsolatban álló lineáris maximálkorrelációt, az un. kanonikus korrelációt is megbecsülik /Judin, 1967; Bagrov, 1968; Gandin, 1967; Glahn, 1968; Dujceva-Pegy, 1970; Bagrov-Mjakiseva, 1969/. Ennek általános elméletével, számítástechnikai realizációjával és meteorológiai alkalmazhatóságával fog-

lalkoztunk már korábban /Faragó-Gulyás, 1974/. Még hatásosabb lehet a mintatér fizikai vagy statisztikai megfontolásokra épülő particionálása és az egyes particiókon való lineáris regresszió kapcsolat létesítése. Ennek egy alkalmazását Faragó /1975/ mutatja be; az optimális particionálás egy speciális esetével Ter-Mkrtcsan/1970 / foglalkozott.

A hidro- és termodinamikai egyenletrendszerből származtatható legegyszerűbb verzió, amely azonban még nagy vonalakban modellezi a légköri folyamatok fejlődését, az ún. divergencia-mentes barotrop modell. Ez az a modell, melyen von Neumann és társai a bevezetőben említett klasszikus vizsgálataikat elvégezték. Ennek révén vált a meteorológiában kitüntetetté az 500 mb-os /közepes vagy barotrop szint/ topografikus mezeje, mely azóta a statisztikai modellezésben is jelentős szerepet játszik. Ezen a szinten ugyanis már kellőképpen lassúak a változások ahhoz, hogy hosszabb érvényességi időtartamu extrapolációra is gondolhassunk, másfelől viszont még elég jelentős a kapcsolat e szint és a talajfelszínhez közeli időjárás között.

Ezért - irodalmi előtanulmányok után - olyan számítástechnikai rendszer tervét fogalmazzuk meg, mely lehetőséget teremt /pl./ az 500 mb-os topográfia előrejelzésére. Figyelembe véve a légköri mozgások átlagos sebességrendjét a hosszabbtávú előrejelzések már feltétlenül hemiszférikus nagyságu mezők figyelembevételét teszik szükségessé. Magaslégköri megfigyelések már több mint 20 éve folynak rendszeresen és nekünk rendelkezésünkre áll az ezek alapján készült, egy reguláris rácshálózatra interpolált napi megfigyelések archivuma az említett 500 mb-os topografikus mezőkből. E reguláris rács a szélességek mentén 10° -os, a hosszúságok mentén 5° -os osztású /a pólus közelében ritkített és csak a 20° é.sz.-ig terjed/.

Láthatóan a prediktor és a prediktandusz mintatere itt egybeesik és mintegy 500 dimenziós. Ekkora dimenziószámú mintákkal a regressziós extrapoláció feladata gyakorlatilag megoldhatatlan: előzetes adatredukcióra van szükség.

2.2. Sorfejtéses adatszűritések és azok szuperponálása

Az 1.2. ponthoz hasonlóan legyen $(x, \mathcal{V}) \in (X, \mathbb{C})$ a vizsgált valószínűségi változó pár /a zérus indexeket elhagytuk/,

W a veszteségfüggvény, d^* pedig a minimális $R^* = R(d^*) = EW(\mathcal{V}, d^*(x))$ átlagos veszteséget adó /bayesi/ becslés.

A mintatér valamilyen mérhető

$$\begin{aligned} /1/ \quad F: \quad X &\rightarrow Y \\ F(x) &= y \end{aligned}$$

transzformációja nyomán - ahol $\{X, \mathcal{X}\}$, $\{Y, \mathcal{Y}\}$ mérhető terek - az eredeti feladat helyett most \mathcal{V} -nak y -mérhető becslését keressük. /Egészen pontosan, ha $\{\Omega, \mathcal{A}, P\}$ jelöli a valószínűségi mezőt, akkor $\mathcal{A} \supset \mathcal{A}_x \supset \mathcal{A}_y$ és most \mathcal{A}_y -mérhető $\hat{\mathcal{V}}$ -becslést keresünk. Itt \mathcal{A}_x , \mathcal{A}_y az x illetve $y=F(x)$ valószínűségi változó által generált σ -algebra./ Erre feltétlenül igaz az

$$/2/ \quad R^* \equiv EW(\mathcal{V}_0, d^*(x_0)) \leq EW(\mathcal{V}_0, d_F^*(F(x_0))) \equiv R_F^*$$

reláció.

Most bebizonyítjuk, hogy például $\Theta = \{0, 1\}$ esetén, vagyis a legegyszerűbb döntési feladatban bizonyos értelemben R^* folytonosan változik a mintatér transzformációinak függvényében. Ebben az esetben

$$/3/ \quad d^*(x) = \text{sign } E(\mathcal{V}/x) = \text{sign}(P(\mathcal{V}=1/x) - P(\mathcal{V}=-1/x))$$

és természetesen W -nek az 1.1.-ben /3a/ alatt bemutatott 0-1 veszteségfüggvényt választjuk. Ekkor az átlagos veszteség a hibavalószínűséget jelenti. Jelölje $\tilde{\sigma}^*(x)$ a feltételes várható értéket:

$$\tilde{\sigma}^*(x) = E(\mathcal{V}/x),$$

akkor tetszőleges d döntésfüggvényre

$$/4/ \quad R(d) = \frac{1}{2} (1 - E\{d(x)\tilde{\sigma}^*(x)\})$$

Lemma: Legyen F_1, F_2 két mérhető leképezés X -ből Y -ba és

$$D \equiv \{ F_1(x) \neq F_2(x) \} \in \mathcal{A} \quad , \quad P\{D\} < \varepsilon \quad ,$$

akkor

$$/5/ \quad |R(d_{F_1}) - R(d_{F_2})| < 2\varepsilon \quad .$$

Bizonyítás. /A /Faragó, 1975a/-ban adottnál sokkal triviálisabb./

$$\begin{aligned} |R(d_{F_1}) - R(d_{F_2})| &\leq |P(d_{F_1}(x) \neq \mathcal{V}/\bar{D}) - P(d_{F_2}(x) = \mathcal{V}/\bar{D})| P\{\bar{D}\} + \\ &+ |P(d_{F_1}(x) \neq \mathcal{V}/D) - P(d_{F_2}(x) \neq \mathcal{V}/D)| P\{D\} \leq 2\varepsilon \quad . \quad \square \end{aligned}$$

1. Tétel /Faragó, 1975a/: Legyen az X mintatér szeparábilis metrikus mintatér a Δ metrikával és tegyük fel, hogy a /3/ döntésfüggvény folytonos a következő értelemben:

$$\forall \varepsilon > 0 \quad \exists \delta > 0$$

$$/6/ \quad Q\{x: d(x') = \text{const}, \text{ ha } \Delta(x, x') < \delta \text{ /m.m./}\} \geq 1 - \varepsilon \quad ,$$

ahol. Q az x valószínűségi változó által generált valószínűségi mérték X -n, $Q\{B\} = P\{X \in B\}$. Legyen továbbá adott a mintatér mérhető transzformációja önmagába, $F: X \rightarrow X$, melyre

$$/7/ \quad \Delta(x, F(x)) < \delta/4 \quad \text{/m.m./} \quad .$$

Akkor a /2/-ben definiált átlagos veszteségek eltérésére:

$$\lambda(F) \equiv |R_F^* - R^*| < \varepsilon \quad .$$

Bizonyítás. /4/-ből

$$/8/ \quad \lambda(F) \leq \int_{\{d_F^*(x) = d^*(x)\}} |s^*(x)| dP \leq Q\{d_F^*(x) = d^*(x)\} \quad .$$

Bevezetjük a következő jelöléseket:

$$X_0 = \{x; d^*(x') = \text{const}, \Delta(x, x') < \delta \text{ /m.m./}\} \quad ,$$

$$U_r(x) = \{x': \Delta(x, x') < r\} \quad , \quad r > 0 \quad ,$$

$$A = \{d_F^* \neq d^*\} \in \mathcal{A} \quad .$$

Válasszunk egy tetszőleges $x \in X_0$ pontot és ennek $U_{\delta/2}(x)$ környezetét. Akkor /7/ következtében

$$V_x \equiv F^{-1}(F \{ U_{\delta/2}(x) \}) \subset U_{\delta}(x) ,$$

továbbá /6/-nak és X_0 definíciójának figyelembevételével, látható, hogy d^* konstans $U_{\delta}(x)$ -n. Az általánosság csorbítása nélkül kiköthetjük, hogy $d^*(x') = 1$ m.m. $x' \in U_{\delta}(x)$ -re. Mivel V_x eleme az F által generált $\sigma_F \subset \mathfrak{E}$ σ -algebrának, felírható

$$\int_{V_x \cap \{d_F^* = -1\}} E(\delta^*/F) dQ = \int_{V_x \cap \{d_F^* = -1\}} \delta^* dQ .$$

Ha itt az integrálási tartomány Q -mértéke nullától különböző, akkor a baloldal negatív, a jobboldal viszont pozitív előjelű. Következésképpen minden $x \in X_0$ -ra $d^*(x') = d_F^*(x')$ m.m. $x' \in U_{\delta/2}(x)$ elemre, hiszen $U_{\delta/2} \subset V_x$. Tehát X szeparabilitásából azonnal következik, hogy $d^*(x) = d_F^*(x)$ m.m. $x \in X_0$ -ra, azaz $Q\{X_0 A\} = 0$. Így /8/-ból:

$$\lambda(F) \leq Q\{\bar{X}_0 A\} \leq Q\{\bar{X}_0\} < \varepsilon ,$$

amit éppen bizonyítani kellett. \square

Akkor is igaz marad a tétel állítása, ha a transzformáció csak átlagosan kismértékben torzítja az eredeti mintaelemeket.

2. Tétel /Faragó, 1975a/: Ha a fenti tétel feltételei között /7/ helyett azt tesszük fel, hogy valamilyen valós $\delta_1^* > 0$ -ra, melyre $\delta_1^* \leq \varepsilon \delta^*/8$, fennáll az

$$/9/ \quad E \Delta(x, F(x)) < \delta_1^*$$

egyenlőtlenség, akkor $\lambda(F) < 2\varepsilon$.

Bizonyítás. Alkalmazzuk a Csebisev-egyenlőtlenséget:

$$Q\{\Delta(x, F(x)) < \frac{\delta_1^*}{4}\} > 1 - 4 \frac{E \Delta(x, F(x))}{\delta_1^*} > 1 - \frac{4\delta_1^*}{\delta_1^*} > 1 - \frac{\varepsilon}{2} .$$

Legyen $X_1 = \{\Delta(x, F(x)) < \delta^*/4\}$ és vezessük be az

$$F_1(x) = \begin{cases} F(x) & , \quad x \in X_1 \\ x & , \quad x \notin X_1 \end{cases}$$

operátort. Akkor a lemma alapján

$$|R(d_F^*) - R(d_{F_1}^*)| < \varepsilon$$

és F_1 kielégíti az előző tétel feltételeit, következésképpen

$$\lambda(F) \leq \lambda(F_1) + |R(d_F^*) - R(d_{F_1}^*)| < 2\varepsilon$$

és ezzel a tételt igazoltuk. \square

A fenti tételek körülményeinek megfelelő tipikus transzformációk a csonkított sorfejtések. Legyen x_t , $t \in T$ véges szórású folyamat, $x_t \in L_2(\Omega, \mathcal{A}, P)$ és $\varphi_k(t)$ egy teljes ortonormált rendszer az $L_2(T)$ térben. Ha most a folyamat sorbafejthető /m.m./

$$x_t = \sum_{k=1}^{\infty} c_k \varphi_k(t)$$

akkor e sorfejtés kellő számú első együtthatója felfogható a folyamat reprezentációjaként és ha a mintatér megfelelő

$$F: x_t \longrightarrow \sum_{k=1}^K c_k \varphi_k(t)$$

transzformációjára teljesülnek valamelyik fenti tétel feltételei, akkor remélhetjük, hogy az eredeti statisztikai feladat is csak kismértékben torzul.

A bizonyos értelemben optimális sorfejtéses adatredukciót az ún. Karhunen-Loève sorfejtés adná meg /Fu, 1960; Watanabe, 1972; Fritz, 1971/, de elvégzéséhez többek között ugyancsak a mintaanyag statisztikai vizsgálata /nevezetesen a kovarianciafüggvény becslése/ szükséges, ami nagy dimenziószám mellett igen körülményes /memória- és gépidőigényes/.

A hemiszférikus /és más/ meteorológiai mezők dinamikai és statisztikai vizsgálataiban igen elterjedtek a spektrális módszerek. Ezek alapján jól ismert, és a magassági /topografikus/ mezőkre is jellemző, hogy a mozgásformák közül a prognosztikában kiemelkedő szerepet játszanak az egészen nagy hullámhosszak /az ún. Rossby-hullámok/, mert ezek egyben a legnagyobb megmaradási hajlamúak is. A hidrodinamikai egyenletrendszerek egyik szűrési kritériuma is éppen az, hogy az ún. meteorológiailag káros /kisebb hullámhosszúságu, meteorológiailag érdektelen, pl.

gravitációs/ hullámokat elkülönítik. Az ezek elhagyásával megadott approximáció tehát nem feltétlenül eredményez analitikusan kitűnő becslést, de fizikailag jól megalapozott és várhatóan elég jó előrejelezhetőséget biztosít.

Mindezek alapján kézenfekvő, hogy a kérdéses mintaelemeket /lásd 2.1. pont/ előbb Fourier-sorba fejtsük.

Jelöljön $G(\varphi, \theta)$ egy a gömbfelületen gömbi koordinátákban adott véletlen függvényt, $0 \leq \varphi \leq 2\pi$, $0 \leq \theta \leq \pi$. Az un. tesszerális szférikus harmonikusok /TSH/ teljes ortonormált rendszerét a

$$/10/ \quad U_n^m(\varphi, \theta) = \sqrt{\frac{2n+1}{2\pi} \frac{(n-m)!}{(n+m)!}} P_n^m(\cos \theta) \cos m\varphi$$

$$V_n^m(\varphi, \theta) = \sqrt{\frac{2n+1}{2\pi} \frac{(n-m)!}{(n+m)!}} P_n^m(\cos \theta) \sin m\varphi$$

függvények alkotják, ahol P_n^m , $n=0,1,\dots$ $m=0,1,\dots,n$ az n -edfoku, m -rendű Legendre-függvény.

Ujabban egyre intenzívebben hasznosítják a TSH-k szerinti sorfejtéses reprezentációt a differenciámódszerek helyett is a hidrodinamikai differenciálegyenlet-rendszerek megoldásában /Jefimov, 1976; Machenhauer, 1974/ többféle lezárási /sorcsontítási/ eljárást ajánlva. Mi itt a legtermészetesebbet választva a G mező approximációját a következőképpen adjuk meg:

$$/11/ \quad G(\varphi, \theta) = \sum_{n=0}^M \left[\sum_{m=0}^n \alpha_{n,m} U_n^m(\varphi, \theta) + \sum_{m=1}^n \beta_{n,m} V_n^m(\varphi, \theta) \right]$$

$$\alpha_{n,m} = \int_0^{2\pi} \int_0^\pi G(\varphi, \theta) U_n^m(\varphi, \theta) \sin \theta \, d\theta \, d\varphi; \quad \beta_{n,m} = \int_0^{2\pi} \int_0^\pi G(\varphi, \theta) V_n^m(\varphi, \theta) \sin \theta \, d\theta \, d\varphi.$$

Hemiszférikus mezőinket az Egyenlítőre szimmetrikusan egészítjük ki, ezáltal mindazon együttthatók zérusok lesznek, melyekre $n+m$ =páratlan.

A kapott $C = (c_{11}, c_{12}, \dots, c_{2M, 2M})$ valószínűségi vektorváltozó most már /csökkentett dimenziószáma révén gyakorlatilag is/ Karhunen-Loève sorba fejthető:

$$/12/ \quad C = \sum_{k=1}^{2M} \sqrt{\lambda_k} f_k \psi_k,$$

ahol $\psi_k \in R_{2M}$, λ_k a C kovarianciamátrixa (Γ) sajátérték-

egyenletének megoldásai:

$$\Gamma \varphi_k = \lambda_k \varphi_k \quad \text{és} \quad f_k = \frac{1}{\sqrt{\lambda_k}} C \varphi_k .$$

E sorfejtés optimális tulajdonságait itt nem részletezzük.

Összefoglalva, nagydimenziószámú minták esetén két egymástkövető /szuperponált/ sorfejtéses adatredukció alkalmazható, melyek eredményeképpen végülis az x_t folyamat reprezentációját az $f = (f_1, f_2, \dots, f_K)$ valószínűségi vektorváltozó adja meg valamilyen alkalmas K esetén.

Konkrét feladatunkban lineáris multiregressziós extrapoláció kiépítése a cél, pontosabban a

$$W(f, f') = \sum_k (f_k - f'_k)^2$$

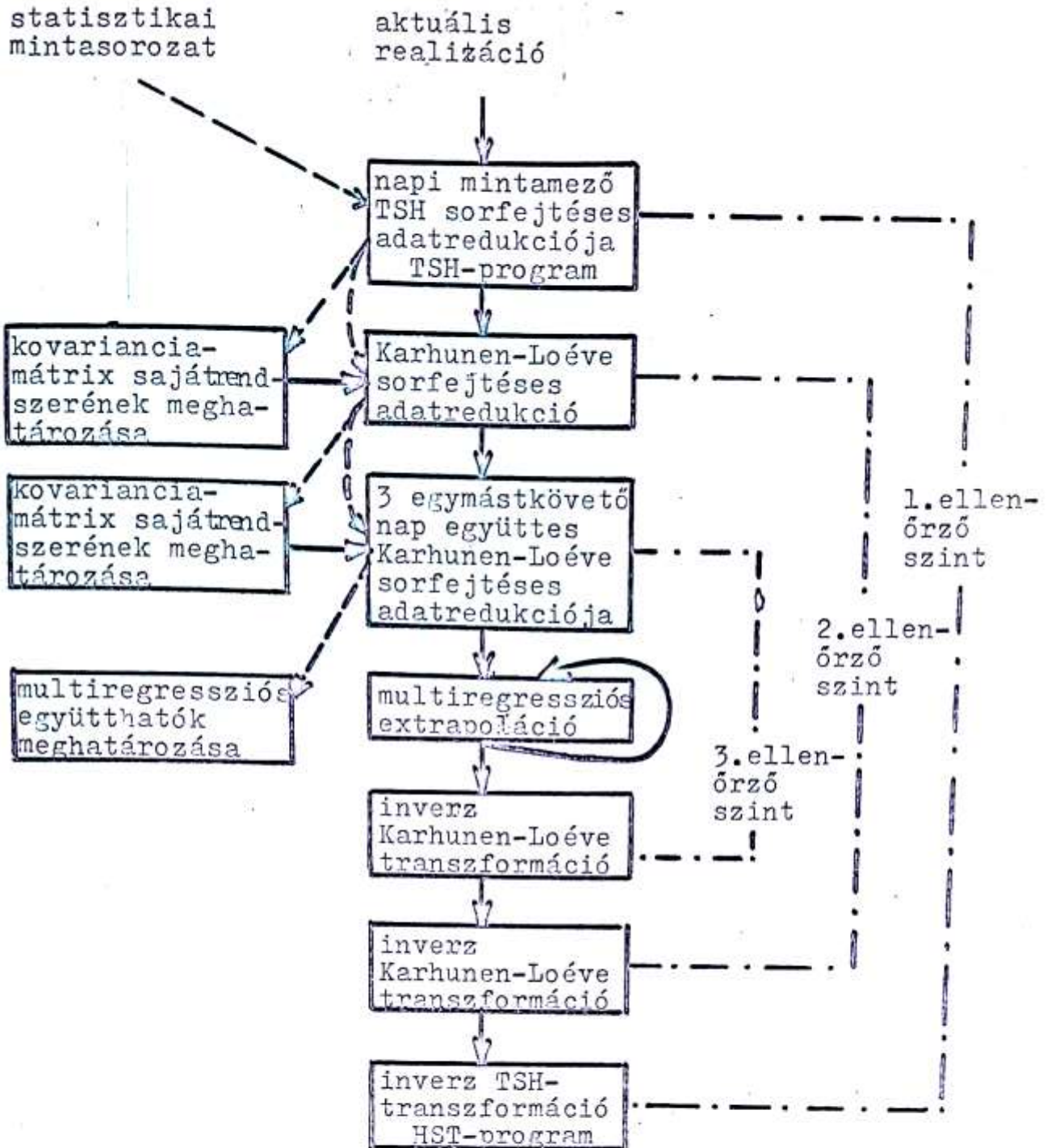
négyzetes veszteségfüggvény mellett - ahol f' az f -t követő időszakhoz tartozó transzformált minta -, az optimális lineáris f' becslés előállítása a feladat.

2.3. Konkrét példa a szuperponált adatszugorításra egy multiregressziós extrapolációs modell esetén

A 2.2.-ben vázolt adatredukciós eljárások egy olyan modellben nyernek alkalmazást, melynek célja az 500 mb-os /abszolút/ topográfia multiregressziós bázisu előrejelzése. A modellt a [1] egy programrendszer realizálja; melyet az 1. ábra mutat be.

E fejezetben az adatszugorításokkal foglalkoztunk és ennek megfelelően a 3., 4., 5. mellékletekben bemutatjuk a spektrális ill. Karhunen-Loève sorfejtéssel összefüggő TSH, HST és KLI programokat. egy-egy futási eredménnyel illusztrálva.

A mintavételezés az 1967-68 -as évekből történt. A TSH-HST programpárral az 1. ellenőrző szinten vizsgálhatjuk meg konkrét realizációkon a spektrális sorfejtéses adatredukció minőségét az egyszerű normabecslésen túlmenően. A példák elemzését a 2. mellékletben adjuk meg.



1. ábra : Szuperponált adatredukciót
megvalósító multiregressziós extrapolációs modellt
realizáló programsorozat

3. EGY "SZELEKTIV LEGKÖZELEBBI TÁRS" STATISZTIKAI BECSLÉSI ELJÁRÁS EMPIRIKUS VIZSGÁLATA

3.1. Makroszinoptikus kódok statisztikai felismerésének problémája

Az 1.1. pontban említett Hess-Brezowsky-féle makroszinoptikus kódok /HB-kódok/ kapcsán adott a következő alakfelismerési feladat.

Az európai-atlanti térség fölött bevezetünk egy 5×7 -es rácshálózatot és ennek rácspontjaiban megadjuk a /megfelelőképpen interpolált/ talajszinti nyomási értéket egy-egy kiválasztott napon. Ilyen 35-dimenziós multbeli mintákból tetszőlegesen hosszú sorozat készíthető úgy, hogy mindegyik mező mellett ismeretes az arra a napra szinoptikusok által megadott HB-kód. A feladat ennek a hozzárendelésnek a statisztikai tanulása.

A meteorológiai szempontból ugyancsak szerény 35 rácspont a matematikai modell kitűzésénél már igen komoly problémát jelent. Ez annál is inkább igaz, mivel az osztályok száma ebben a konkrét feladatban 30. Ezért annak ellenére, hogy számos eljárás ismeretes a feltételes eloszlások közvetett becslésére, szeparáló függvényeket előállító tanuló-felismerő algoritmusokra, mégis csupán "Nearest Neighbour" /NN/ -tipusu eljárásokra támaszkodhatunk. Ebben az esetben a legnagyobb gond a tárolandó mintaelemek nagy száma és ezzel összefüggésben az analógiák várhatóan nagy keresési ideje. Éppen ezért olyan módszerre van szükség, mely eleve csökkenti az archiválásra kerülő minták számát.

3.2. Egy szelektív NN-eljárás /SNN/

Legyen $T_N = \{ (x_k, \mathcal{V}_k) \}_{k=1}^N, (x_k, \mathcal{V}_k), (x_0, \mathcal{V}_0) \in (X, \mathcal{C})$ független, azonos eloszlású mintasorozat; $X \subset R_M$ -metrikus mintér a Δ metrikával, \mathcal{C} pedig véges /diszkrét/ paramétertér /ld. 1.2. pont/.

A közönséges NN-eljárásra támaszkodó alakfelismerési modell-

ben, a "tanulási fázisban" minden beérkező (x_k, ν_k) mintapárt megőrzzük és ennek megfelelően T_N mintasorozat mellett N távolságot kell kiszámítatunk x_0 és x_1, x_2, \dots, x_N között, illetve a legközelebbi társ kikereséséhez $O(N^2)$ összehasonlítást kell tennünk. Itt N rendszerint elég nagy értékeket érhet el és ezért egy-egy osztályozás meglehetősen "költséges".

E gondolat jegyében az NN-módszer különféle javításai láttak napvilágot /Hart,1968; Gates,1972/.

Tipikus az olyan módosítás, mikoris a teljes T_N mintanyagot önmagán tesztelik és ennek nyomán bizonyos elemeit elhagyják. Mi most egy olyan eljárást mutatunk be, melyben az archívum számossága általában véve nem növekszik olyan gyorsan, mint N .

Minden mintapár mellé hozzá fogunk rendelni egy ν_k természetes számot, amelyet osztályozási gyakoriságnak nevezünk. Az algoritmust a következőképpen adjuk meg: az első lépésben az archívum egyetlen elemet tartalmaz

$$S_1 = \{ (x_1, \nu_1) \}, \quad \nu_1 = 1;$$

a $/k+1/-$ edik $(x_{k+1}, \nu_{k+1}) \in T_N$ minta módosítja a k -edik lépésben adott S_k archívumot

$$/1/ \quad S_{k+1} = \begin{cases} S_k \cup \{ (x_{k+1}, \nu_{k+1}) \} & \text{ha } \nu_{k+1} \neq \nu \\ [S_k - \{ (x, \nu) \}] \cup \{ (\tilde{x}, \tilde{\nu}) \} & \text{ha } \nu_{k+1} = \nu \end{cases}$$

és megfelelőképpen az új mintaelem osztályozási gyakorisága $\nu_{k+1} = 1$ lesz, ha $\nu_{k+1} \neq \nu$, illetve a módosított (x, ν) minta osztályozási gyakorisága is változik eggyel $\tilde{\nu} = \nu + 1$, amikor $\nu_{k+1} = \nu$. Itt (x_{k+1}, ν_{k+1}) legközelebbi társa az S_k mintasorozatból (x, ν) , illetve ν jelöli ennek osztályozási gyakoriságát; $\nu_{k+1} = \nu$ esetén a módosított minta:

$$/2/ \quad \tilde{x} = \frac{\nu \cdot x + x_{k+1}}{k+1}, \quad \tilde{\nu} = \nu.$$

Másszóval, ha az eddigi archívum alapján az új minta osztályozása helytelen, akkor azt felvesszük az archívumba $\nu_{k+1} = 1$ -

gyel; ha az osztályozás helyes, akkor csak az ezt osztályozó mintát korrigáljuk /2/ szerint, valamint az osztályozási gyakoriságát növeljük eggyel. Nyilvánvaló, hogy olyan szeparálható osztályok esetén, amikoris tetszőleges $x \in X$ esetén legfeljebb egy olyan $\vartheta \in \Theta$ létezik, hogy az $f(\vartheta/x)$ feltételes sűrűségfüggvény pozitív és az $X_{\vartheta} = \{x: f(\vartheta/x) > 0\}$ jelöléssel $\forall \vartheta$ -ra és $\forall x \in X_{\vartheta}$ -re teljesül a

$$\forall x' \in X_{\vartheta} \quad \forall \vartheta'' \neq \vartheta \quad \forall x'' \in X_{\vartheta''} \quad \Delta(x, x') < \Delta(x, x'')$$

feltétel, akkor S_N összesen legfeljebb $\text{card}(\Theta)$ mintapárból fog állni és ezek első tagjai /2/ miatt rendre konvergálnak a megfelelő $E(x/X_{\vartheta})$ feltételes várható értékhez. /Hiszen ekkor minden NN-osztályozás korrekt lesz /m.m./ és így minden X_{ϑ} -ra külön-külön alkalmazható a nagy számok tétele./ Ettől az egészen speciális esettől eltekintve, ennek az SNN-eljárásnak a konvergenciájára nézve semmilyen elméleti eredmény nem ismeretes. A klasszikus NN-konvergenciák érvénytelenek, hiszen itt az archivum elemei változnak N növekedésével. Ezért egy egyszerű empirikus vizsgálat hivatott az SNN-eljárás egy előzetes elemzésére. Ennek során megvizsgáljuk azt a verziót is, amikor /2/ a következőképpen módosul: $\tilde{x} = x$, vagyis az archivum elemei nem változnak meg. Nevezzük ezt módosított NN-eljárásnak /MNN/. Ennek konvergenciája triviálisan következik a közönséges NN-módszerek konvergenciájából.

Legyen $\Theta = \{0, 1\}$ és legyen $f(x/\vartheta)$ gaussi sűrűségfüggvény $\vartheta = 0$ esetén $m_0 = 0$, $\vartheta = 1$ esetén pedig $m_1 = 2$ várhatóértékkel, valamint mindkét esetben $\sigma_0 = \sigma_1 = 1$ szórással. Legyen továbbá $P\{\vartheta = 0\} = P\{\vartheta = 1\} = 1/2$. Ebben a statisztikai felismerési feladatban az elméleti Bayes-döntést nyilvánvalóan a

$$/3/ \quad \vartheta^* = \begin{cases} 0 & , \text{ ha } x_0 < 1 \\ 1 & , \text{ ha } x_0 \geq 1 \end{cases}$$

valószínűségi változó adja meg. A Bayes-hiba

$$\begin{aligned} /4/ \quad P\{\vartheta_0 \neq \vartheta^*\} &= P\{\vartheta_0 = 0, \vartheta^* = 1\} + P\{\vartheta_0 = 1, \vartheta^* = 0\} = \\ &= P\{x_0 \geq 1 / \vartheta_0 = 0\} P\{\vartheta_0 = 0\} + P\{x_0 < 1 / \vartheta_0 = 1\} P\{\vartheta_0 = 1\} = \end{aligned}$$

$$= \frac{1}{2} \left\{ \left[1 - \int_{-\infty}^1 \varphi(x; 0, 1) dx \right] + \int_{-\infty}^1 \varphi(x; 2, 1) dx \right\} = 0,16.$$

ahol φ a normális eloszlás sűrűségfüggvénye az adott paraméterek-ké. $N=100$ esetén egy konkrét numerikus példában.

	NN	MNN	SNN
az archivum számossága	100	32	28

valamint 100 újabb /véletlen gaussi számgenerátorral előállított/ mintával számítva az empirikus hibavalószínűségek

	NN	MNN	SNN
az empirikus hibavalószínűség	0,20	0,18	0,18 .

A mintaanyagból párhuzamosan megbecsültük a Bayes-hibát is és eredményül 0,13-t kaptunk, ami igen jól közelíti a /4/ alatti elméleti értéket. /A realizáló programot és a három futási eredményt a 6. melléklet mutatja be./

3.3. Az SNN-eljárás alkalmazása

A 3.2.-ben vázolt szelektív NN-módszer első közelítésben sikerrel vizsgázott a 3.1.-ben leírt feladat kapcsán. A 30 osztály közül összesen 26-ból gyűjtöttünk mintákat a becsült a priori valószínűségekkel arányos számban. A szekvenciális tanulást megelőzően mindegyik osztályból kisorsoltunk 1-1 elemet és előre archiváltuk. Ezt követően 110 mintával "tanítottunk": 37 esetben az új minta paramétere egybeesett annak NN-becsülésével a megelőző mintákból. /Mivel a szinoptikai gyakorlatból ismeretes, hogy némely HB-osztályok elég nagy hasonlóságot mutathatnak, tehát feltehetően a Bayes-hiba eleve meglehetősen nagy, ezért figyeltük, hogy a második legközelebbi társ, illetve a szinoptikus szempontból hasonló NN-döntések száma mekkora. Ebben az értelemben még további 28 esetről volt elmondható, hogy elég jó becsülésnek minősült./ Természetesen a módszer konvergenciája és minőségének becslése még további vizsgálatokat igényel.

4. MAXIMUMÉRTÉKEK ELOSZLÁSBECSLÉSÉNEK IDŐEGYSÉG-PROBLÉMÁJA

4.1. Egy "rövid" szélsőérték-maximum megfigyelési sorozat

Az extrémális értékek statisztikai vizsgálata igen fontos szerepet játszik a klimatológiában. A csapadék-maximumokkal foglalkozott többek között Hershfield/1973/, míg mások a szélsőérték maximális értékeinek eloszlását elemezték /Péczely, 1965; Bozó, 1964/. A téma kapcsán általában a legnagyobb gond a kellő számú, statisztikailag homogén minta hiánya.

Miskolcra például 1955 óta rendelkezünk olyan szélsőértéki adatokkal, melyek a kérdéses statisztikai vizsgálat alapjául szolgálhatnak. A tervezésben /pl. épületek, kémények stb. tervezésében/ általában az a kérdés merül fel, hogy mekkora a valószínűsége egy-egy meghatározott küszöbértéket túllépő szélsőértéknek adott időszakban - pl. két év - alatt. Tehát ha az 1955-től 1974-ig terjedő időszakból éves maximumokat képezünk a szokásos módon, akkor összesen csak 20 minta áll rendelkezésünkre a kérdéses statisztika becslésére. Mivel általában az említett küszöbértékek az eloszlás "szélén" helyezkednek el, ez azt jelenti, hogy viszonylag kis valószínűségű eseményekről van szó. Ennek fényében méginkább beátható, hogy a megfigyelések említett száma meglehetősen kicsi.

A kétéves maximum eloszlásnak becslése azonban nem feltétlenül csak éves maximumértékekből végezhető el. Ha statisztikailag homogén marad a mintasorozat, akkor lecsökkenthetjük ezt az elemi időszakot például félévre. A fő kérdés most az, miként változik meg ekkor a becslés minősége.

4.2. A maximumeloszlás az időegység függvényében

Legyen $X_1, X_2, \dots, X_M, \dots$ független, azonos eloszlású valószínűségi változók sorozata. Célunk a

$$/1/ \quad p_T = P \left\{ \max_{t: t-M_1=1, 2, \dots, T} X_t < x \right\}$$

valószínűség becslése valamilyen rögzített x , M_1 és T -re.

/A gyakorlatban $M_1 > M$, de ennek a modell szempontjából nincs különösebb jelentősége./ Tegyük fel, hogy a T "időintervallum" w egyenlő, egymástkövető τ hosszúságú részre osztható, akkor

$$/2/ \quad p_T = \prod_{k=0}^{w-1} P \left\{ \max_{t: t-M_1=k\tau+1, k\tau+2, \dots, k\tau+\tau} X_t < x \right\} = p_\tau^w .$$

Tegyük fel továbbá, hogy M számú megfigyelésünk van /M rögzített/ és az egyszerűség kedvéért M is osztható τ -val, $M=N\tau$, akkor p_τ becslése:

$$/3/ \quad \hat{p}_\tau = \frac{1}{N} \sum_{k=0}^{N-1} I_k$$

ahol

$$/4/ \quad I_k = \begin{cases} 1 & , \text{ ha } \max \{ X_{k\tau+1}, \dots, X_{k\tau+\tau} \} < x \\ 0 & , \text{ egyébként} \end{cases} .$$

Következésképpen p_T becslése

$$/5/ \quad \hat{p}_T = \hat{p}_\tau^w .$$

Mindenekelőtt belátjuk az /5/ becslés aszimptótikus torzítatlanságát /Faragó, 1977/.

1. tétel: Fennáll a következő aszimptótikus reláció:

$$/6/ \quad E \hat{p}_T = p_T + \frac{1}{N} \left\{ -\frac{w(w-1)}{2} p_\tau^w + \frac{(w-1)^2}{2} p_\tau^{w-1} \right\} + O(N^{-2}) .$$

Bizonyítás. /3/-ből

$$/7/ \quad E \hat{p}_T = \frac{1}{N^w} E \left\{ \sum_{\sum k_i = w} I_1^{k_1} I_2^{k_2} \dots I_N^{k_N} \frac{w!}{k_1! k_2! \dots k_N!} \right\} =$$

$$= \frac{1}{N^w} \sum_{k=1}^w J_k p_\tau^k ,$$

hacsak $N \geq w$. Itt

$$J_k = \sum_{\exists k_{i_1}, \dots, k_{i_k} \neq 0} \frac{w!}{k_1! k_2! \dots k_N!} =$$

és $k_i = 0$ egyébként,
 hogy $k_{i_1} + \dots + k_{i_k} = w$

$$= \binom{N}{2} \sum_{\substack{\sum k_{ij} = w \\ k_{ij} \geq 1}} \frac{w!}{k_{i_1}! k_{i_2}! \dots k_{i_w}!} \cdot$$

Speciálisan

$$/8/ \quad J_w = \binom{N}{w} \cdot w! \quad , \quad J_{w-1} = \binom{N}{w-1} (w-1) \frac{w!}{2} \quad .$$

/7/-ben : csak a leglassabban konvergáló első két tagot megtartva, majd /8/-ből a behelyettesítéseket elvégezve

$$\begin{aligned} E\hat{p}_T &= \frac{1}{N^w} \left[J_w p_\tau^w + J_{w-1} p_\tau^{w-1} + o(N^{w-2}) \right] = \\ &= p_\tau^w + \frac{1}{N} \left[- \frac{w(w-1)}{2} p_\tau^w + \frac{(w-1)^2}{2} p_\tau^{w-1} \right] + o(N^{-2}) \quad , \end{aligned}$$

amit /2/ figyelembevételével éppen bizonyítani kellett. \square

A fenti tételből azonnal következik a p_T becslés aszimptotikus torzítatlansága.

A következő tétel e becslés aszimptotikusan érős konzisztenciáját mutatja meg.

2. tétel: Az /5/ becslés szórására teljesül, hogy

$$/9/ \quad D\hat{p}_T = \frac{1}{N} p_\tau^{2w-1} \left(w^2 - w^2 p_\tau - \frac{1}{2} \right) + o(N^{-2}) \quad .$$

Bizonyítás. Fejtsük ki a szórást:

$$D\hat{p}_T = E \left\{ \frac{1}{N} \sum_{k=1}^N I_k \right\}^{2w} - (E\hat{p}_T)^2 \quad .$$

Az első tagra alkalmazható az 1. tétel, ha $N \geq 2w$:

$$E\hat{p}_T^2 = p_T^2 + \frac{1}{N} \alpha(2w; p_\tau) + o(N^{-2})$$

az

$$\alpha(w; p_\tau) = - \frac{w(w-1)}{2} p_\tau^w + \frac{(w-1)^2}{2} \cdot p_\tau^{w-1}$$

jelöléssel. Ugyancsak az 1. tétel alapján

$$(E\hat{p}_T)^2 = p_T^2 + \frac{2}{N} \alpha(w; p_\tau) \cdot p_T + o(N^{-2}).$$

Következésképpen

$$D\hat{p}_T = \frac{1}{N} [\alpha(2w; p_\tau) - 2 \alpha(w; p_\tau) \cdot p_T] + o(N^{-2}),$$

ami α behelyettesítésével éppen /9/-t adja. \square

Vizsgáljuk most meg /9/ függését az időegységtől. Pontosabban, mivel /9/ aszimptotikáját a

$$/10/ \quad \beta(\tau; p_1, T) = p_\tau^{2w-1} \left(w^2 - w^2 p_\tau - \frac{1}{2} \right)$$

együtthető határozza meg, helyettesítsük be ide az

$$N=M/\tau, \quad w=T/\tau, \quad p_\tau = p_1^\tau$$

kifejezéseket. Akkor

$$/11/ \quad \beta(\tau; p_1, T) = \tau \cdot p_1^{2T-\tau} \left[(1-p_1^\tau) \frac{T^2}{\tau^2} - \frac{1}{2} \right]$$

egyenlőség adja meg a vizsgálandó függvénykapcsolatot.

Abban a speciális esetben, amikor a $\tau_1=1$ és $\tau_2=2$ időegységeket kívánjuk összehasonlítani, azt kapjuk, hogy

$$/12/ \quad |\beta(2; p_1, T)| - |\beta(1; p_1, T)| = p_1^{2T-2} \left[2 \left| (1-p_1^2) \frac{T^2}{4} - \frac{1}{2} \right| - p_1 \left| (1-p_1) T^2 - \frac{1}{2} \right| \right].$$

Ennek most csupán előjelét vizsgálva a szögletes zárójelben kifejezés értékeit $T=2,3,4$ esetben az 1. ábrán vázoltuk fel. Ebből azt a következtetést vonhatjuk le, hogy $\tau_1=1$ mellett aszimptotikusan jobb az /5/ becslés, mint $\tau_2=2$ -re egy x_1 küszöbérték alatt, illetve egy $x_2 / x_1 < x_2 < 1$ /küszöb felett.

4.3. Gyengén függő minta és egy paraméteres becslés esete

A /4/-ben bevezetett indikátorfüggvényeket /illetve a megfelelő eseményeket/ erősen keverőnek mondjuk /Rosenblatt, 1956;

Ibragimov-Rozanov, 1970/, ha minden $q > 0$ -ra

$$/13/ \quad \vartheta_n(q) = \sup_{\substack{n_1, \dots, n_q \\ n(n_1, \dots, n_q) = n}} \left| E \left(\prod_{i=1}^q I_{n_i} \right) - \prod_{i=1}^q E I_{n_i} \right| \xrightarrow{n \rightarrow \infty} 0,$$

ahol $n = n(n_1, n_2, \dots, n_q) = \min_{1 \leq i, j \leq q} |n_i - n_j|$.

3. tétel: Feltéve, hogy

$$/14/ \quad \frac{1}{N} \sum_{n=1}^N \vartheta_n(w) \rightarrow 0$$

az /5/ becslés aszimptóti kusan torzítatlan és konzisztens marad.

Bizonyítás. Jelölje \hat{r}_T az /5/ becslést a /13/-mal adott függő mintára, akkor

$$E \hat{r}_T - E \hat{p}_T = \frac{1}{N^w} \left\{ \sum_{n=1}^{\left[\frac{n-1}{w-1} \right]} \sum_{\substack{n_1, \dots, n_w \\ n_i \neq n_j \\ n(n_1, \dots, n_w) = n}} [E \left(\prod_{i=1}^w I_{n_i} \right) - p_T^w] w! + \right. \\ \left. + \sum_{\nu=1}^{w-1} \sum_{\substack{n_1, \dots, n_\nu \\ n_i \neq n_j \\ k_i \neq k_j, \sum k_i = w}} [E \left(\prod_{i=1}^{\nu} I_{n_i}^{k_i} \right) - p_T^\nu] \frac{w!}{k_1! k_2! \dots k_\nu!} \right\} .$$

Az egyszerűség kedvéért tegyük fel, hogy $N-1$ osztható $(w-1)$ -gyel, $L(w-1) = N-1$, valamint A_n jelölje az összes olyan elrendezést az $1 \leq n_1, n_2, \dots, n_w \leq N$ ($n_i \neq n_j$) számokból, melyekre $n(n_1, n_2, \dots, n_w) \geq n$

$$A_n = \binom{w + (N-1 - (w-1)n)}{w}$$

és B_n azokat, melyekre $n(n_1, n_2, \dots, n_w) = n$

$$B_n = A_n - A_{n+1} = \frac{1}{w!} (w-1)^w N^{w-1} + o(N^{w-2}) \quad n \leq L-1$$

$$B_L = A_L = 1$$

Ekkor

$$|E \hat{r}_T - E \hat{p}_T| \leq \frac{1}{N^w} \cdot w! \sum_{n=0}^L B_n \vartheta_n(w) \leq$$

$$1. \leq \frac{1}{N} (w-1) w \sum_{n=0}^L \xi_n \leq \frac{1}{N} (w-1) w \sum_{n=0}^N \xi_n \rightarrow 0 .$$

Az 1.tétel és /14/ alapján ebből már közvetlenül belátható \hat{r}_T aszimptótikus torzítatlansága. Részben a 2.tétel bizonyítására támaszkodva a konzisztencia hasonlóképpen igazolható. \square

Következésképpen az /5/ becslés alkalmazható a /13/-szerint gyengén függő mintákra is, de az erősen keverő minták definíciója túl általános ahhoz, hogy itt konkrétan megvizsgálható legyen az időegység-probléma.

Paraméteres becslés esetén már többet mondhatunk. Legyen az $X_1, X_2, \dots, X_M, \dots$ mintaelemek eloszlása a priori ismertén gaussi m_1 várható értékkel és σ_1 szórással. A paraméteres becslés p_T -re:

$$/15/ \quad \hat{p}_T = \hat{p}_T^w = \left[\int_{-\infty}^x \varphi(z; \hat{m}_\tau, \sigma_\tau) dz \right]^w$$

feltéve itt, hogy σ_τ előre ismert, \hat{m}_τ pedig a megfelelő empirikus várhatóérték:

$$\hat{m}_\tau = \frac{1}{N} \sum_{k=0}^{N-1} \max \{ X_{k\tau+1}, \dots, X_{k\tau+\tau} \} .$$

/ φ a normális sűrűségfüggvény $\hat{m}_\tau, \sigma_\tau$ paraméterekkel./

4.tétel: A /15/ becslés szórásának aszimptótikus alakja:

$$/16/ \quad D\hat{p}_T = \frac{1}{M} p_1^{(T-\tau)} T^2 \frac{1}{2\tau\sigma_\tau} \exp \left\{ -\frac{(x-m_\tau)^2}{2\sigma_\tau^2} \right\} + o(N^{-2}) .$$

Bizonyítás. Fejtsük Taylor-sorba a sűrűségfüggvényt a várhatóérték paramétere szerint:

$$\begin{aligned} \varphi(z; \hat{m}_\tau, \sigma_\tau) &= \varphi(z; m_\tau, \sigma_\tau) - \varphi(z; m_\tau, \sigma_\tau) \frac{z-m_\tau}{\sigma_\tau} (\hat{m}_\tau - m_\tau) + \\ &+ \frac{1}{2} \varphi(z; m_\tau, \sigma_\tau) \frac{1}{\sigma_\tau^2} \left[\frac{(z-m_\tau)^2}{\sigma_\tau^2} + 1 \right] (\hat{m}_\tau - m_\tau)^2 + o(\hat{m}_\tau - m_\tau)^3 . \end{aligned}$$

Ezt /15/-be helyettesítve:

$$E\hat{p}_T = p_T + \left\{ w p^{w-1} \frac{1}{2} \int_{-\infty}^x \varphi(z; m_\tau, \sigma_\tau) \frac{1}{\sigma_\tau^2} \left[\frac{(z-m_\tau)^2}{\sigma_\tau^2} + 1 \right] dz \right\} +$$

$$+ \binom{w}{2} p_{\tau}^{w-2} \left[\int_{-\infty}^x \varphi(z; m_{\tau}, \sigma_{\tau}) \frac{(z-m_{\tau})}{\sigma_{\tau}^2} dz \right]^2 \left\} \frac{\sigma_{\tau}^2}{N} + o(N^{-2}) .$$

Jelölje $\alpha(w; p_{\tau})$ ismét az N^{-1} tag együtthatóját. Vegyük észre, hogy

$$\int_{-\infty}^x z \varphi(z; m_{\tau}, \sigma_{\tau}) dz = - \sigma_{\tau}^2 \varphi(x; m_{\tau}, \sigma_{\tau}) + m_{\tau} p_{\tau}$$

és így

$$\begin{aligned} D\hat{p}_{\tau} &= \frac{1}{N} \left[\alpha(2w; p_{\tau}) - 2 \alpha(w; p_{\tau}) p_{\tau} + o(N^{-2}) \right] = \\ &= p_{\tau}^{2w-2} w^2 \varphi^2(x; m_{\tau}, \sigma_{\tau}) \frac{\sigma_{\tau}^2}{N} + o(N^{-2}) = \\ &= \frac{1}{M} p_1^2 (T-\tau) T^2 \frac{1}{2\pi\tau} \exp\left\{-\frac{(x-m_{\tau})^2}{2\sigma_{\tau}^2}\right\} + o(M^{-2}) . \end{aligned}$$

Ezzel a tételt bizonyítottuk.

A 3. tételből azonnal következik a /15/ becslés aszimptótikus torzítatlansága és könnyen levezethető ennek aszimptótikus konzisztenciája is.

Jelölje /16/-ban M^{-1} együtthatóját $\beta(\tau; p_1, T)$. Megvizsgálva az $\tau_1=1$, $\tau_2=2$ speciális esetet, azt kapjuk, hogy

$$|\beta(2; p_1, T)| - |\beta(1; p_1, T)| \geq 0$$

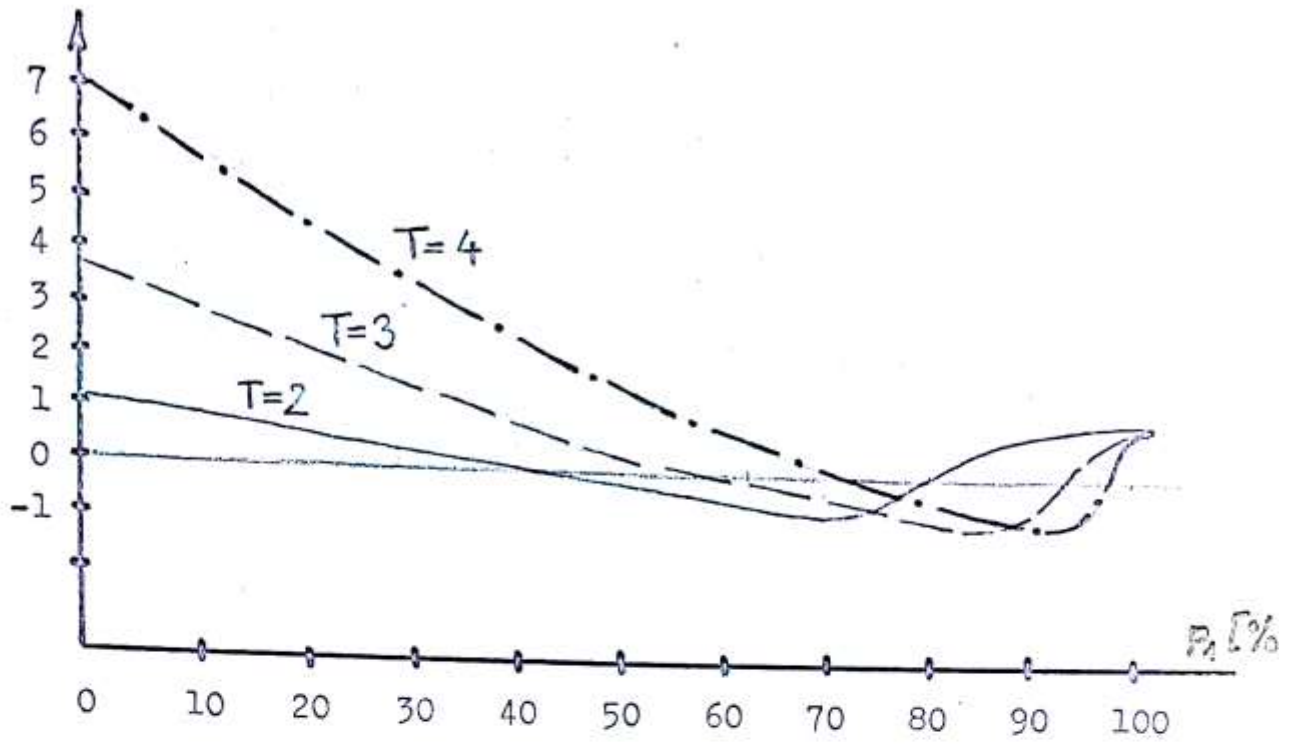
/felhasználva az $m_2=m_1 + \sigma_1/\sqrt{\pi}$ és $\sigma_2^2=(1-1/\pi) \sigma_1^2$ relációkat/, ami azt jelenti, hogy általában a kisebbik időegységű paraméteres becslés jobb aszimptótikusan.

4.4. Egy szélsőbességi mérési sorozat analízise féléves elemi maximumértékek alkalmazásával

A 4.1.-ben említett mintasorból féléves maximumértékeket számítottunk a szokásos éves értékek helyett. A 4.2. és 4.3. pontban bemutatott eredmények alapján helyesebb ennek az időegységnek az alkalmazása, ha az $N_1=20$ éves érték helyett kapott $N=40$ féléves minta továbbra is statisztikailag homogén marad.

Ennek érdekében /Faragó,1977/ -ben először az első illetve második félévi értékekből álló kétdimenziós mintasorozat normalitását igazoltuk statisztikailag /Hald,1967/, majd a korrelációs együtthatóra t-statisztikát számítottunk / $\hat{\rho} = 0,22$; $t = 0,96$; 18 szabadsági fok/, ami alapján 95%-os szignifikanciaszinten elfogadható volt a korrelálatlanság, azaz a függetlenség hipotézise. /Az itt alkalmazott statisztikai eljárásokkal kapcsolatosan Lehmann/1966/, Vincze/1973/ és Prékopa/1972/ munkáira hivatkozunk./

A homogenitás ellenőrzése érdekében előbb F-teszt következett / $F = 1,93$ /, majd az empirikus várható értékek egyenlőségének igazolására ismét Student-próba / $t = 1,07$; a szabadsági fokok száma 38/, amelyek következtében 90 illetve 95%-os szinten elfogadtuk a két szóásnégyzet illetve a két várhatóérték azonosságának hipotézisét. Mindezek alapján a féléves időegységekre kiszámítottuk az empirikus várhatóértéket és szórást és a továbbiakban az e paraméterekkel adott normális eloszlásnak feltételezett féléves időegységű szélmaximum képezte a tetszőleges számú év alatt a tetszőleges küszöbérték túllépési valószínűségek becslésének alapját.



1. ábra

ZÁRÓMEGJEGYZÉSEK

A meteorológiai megfigyelések rendszeres archiválása, a meteorológiai szolgáltatások összessége, a kutatások sorában pedig elsősorban a statisztikai vizsgálatok széles köre számítógépes meteorológiai adatbank és az ezzel kapcsolatban álló adatbankszervező és feldolgozó programok rendszerének létrehozását igényli. Ebbe a témakörbe vágó kérdésekkel foglalkozik többek között a korábbi helyzetről átfogó képet adó WMO/1964/ kiadvány, vagy például Boriszenkov-Romanov/1969/ és Craddock/1970/ munkája. Az említett rendszer megteremtésével kapcsolatos alapelveket, illetve az eddig megvalósult lépéseket egy tanulmányban foglaltuk össze /Faragó et al., 1976/.

Az adott típusú feldolgozás szerint orientált adatbankok általános ismérvein - gyors adathozzáférhetőségen, az adatok különféle rendezettségű elérésének egyszerű szerkesztésén, állománykiegészítési tevékenység kiépítésén stb. - túlmenően a statisztikai természetű általános /például alkalmazott klimatológiai/ és egyedi /általában klimatológián vagy prognosztikán belüli kutatási/ célok egyéb szempontokat is előírnak. A legfontosabb ezek sorában, hogy különféle - többé-kevésbé szekvenciális - mintavételezéseket könnyen, gyorsan és mi több valamilyen speciális "lekérdező" nyelv segítségével egyszerűen megadhatóan végezhesünk el. Említésreméltó ezen túlmenően a "származtatott" állományok problémája is. Nagykiterjedésű adatmezők, vektorok esetén, mint például a 2. fejezetben láttuk, a statisztikai vizsgálatokat az eredeti minták előzetes transformációja /adatredukció/ előzheti meg. Ez pedig gyakorlatilag azt jelenti, hogy a megadott "eredeti" adatállományok mellett az adatbankba felvesszük ezek bizonyos változatait is.

A statisztikai feldolgozás-orientált /pl. távprognosztikai, klimatológiai/ adatbank mellett várhatóan más jellegű meteorológiai adatbankok is kidolgozásra kerülnek majd az Országos Meteorológiai Szolgálat néhány éven belül esedékes nagykapacitású R-kategóriájú számítógépén. A jelenlegi előkészítő szakaszban,

korábban az SzKI IBM 370/125 típusu gépén dolgoztunk DOS operációs rendszerben, ujabban pedig a SzÜV 370/145 -ös gépén OS rendszerben. A statisztikai feldolgozások szempontjából elsősorban a rendelkezésünkre álló SSP /Scientific Subroutine Package/ tudományos programkönyvtár érdemel külön említést, mely egy sor statisztikai becslési és hipotézisvizsgálati eljárást tartalmaz számítógépes /Fortran/ szubrutin formájában. Például a 2.3. pontban vázolt modell messzemenően hasznosítja az SSP-ben adott különféle eljárásokat.

A jelenleg folyó statisztikai vizsgálatok szemszögéből is már alapvető jelentőségű a számítógépes feldolgozási lehetőség, a későbbiek során várható számítógép installálás után pedig a jelenleg sokszor csak elméletileg kidolgozott, vagy más kül- és belföldi kutatók által analizált eljárások tég köre lesz nagyobb adatseregre alkalmazható, illetve adaptálható.

A távprognosztikával kapcsolatos kutatások az Országos Meteorológiai Szolgálatnál jelenleg a Távprognosztikai Csoport keretében folynak, melyek szerves részét képezik az 1.1., 2.1. és 3.1. pontban leírt feladatok. E vizsgálatokat Czövek Istvánné tudományos munkatárssal, meteorológussal közösen végezzük; a globális regressziós modell kidolgozása és számítástechnikai realizálása során Szlachányi Kornélné tudományos segédmunkatárs is bekapcsolódott a kutatásokba. A programozási /főként PL1-nyelvű programszerkesztési / munkákban nagy szerepe van Salamon Lászlóné rendszerszervezőnek. Az adatszugorításokkal kapcsolatos elmélet és a regressziószámítási alkalmazások terén korábban Gulyás Ottó kandidátussal, a Módszertani Csoport vezetőjével végeztünk közösen kutatásokat. A 4.1. pontban leírt feladatot Szalma Jánosné tudományos munkatárs, az Éghajlati Tájékoztató Osztály tagja vetette fel.

MELLÉKLETEK

Hivatkozások

- Bagrov, N.A., 1959: Analiticeszkoje predstavljenije poszledovatelnyosztyi meteorologiceszskih polej poszredsztvom jeszjesztvennih ortogonalnix szosztovljajuscsih. Trudi CIP, 74
- Bagrov, N.A., 1966: Predszkazanyije meszjacsnovo koliceszstva oszadkov. Meteorologija i Gidrologija, 7.
- Bagrov, N.A., 1968: O nyekotorih oszobennosztyah korrelacionovno analiza i ih primenyenyija k prognozam pogodi, Meteorologija i Gidrologija, 1.
- Bagrov, N.A., Mjakiseva, N.N., 1969: Prognoz tyemperaturi na meszjac. Trudi GMC, 44.
- Boriszenkov, E.P., Romanov, M.A., 1969: Algoritmi i programma sztatisticeszkoj obrabotki informacii na EVM. Gidrometeorizdat, Leningrad.
- Bozó P., 1974: A szélmaximum eloszlásának becsléséről. Kézirat. Orsz. Meteorológiai Szolgálat, Budapest.
- Cover, T.M., Hart, P.E., 1967: Nearest neighbour pattern classification. IEEE Trans. on Inf. Theory, 1.
- Cover, T.M., 1968a: Rates of convergence of nearest neighbour decision procedures. First Annual Hawaii Intern. Conf. on Systems Theory.
- Cover, T.M., 1968b: Estimation by the NN rule. IEEE Trans. on Inf. Theory, 1.
- Cover, T.M., 1969: Learning in pattern recognition. Metodologies of pattern recognition, S. Watanabe ed. Academic Press, New York.

- Cradock, J.M., 1958: Research in the Meteorological Office concerned with long-range forecasting. *Meteorological Magazine*, 3.
- Cradock, J.M., 1970: Work in synoptic climatology with a digitized data bank. *Meteorological Magazine*, 2.
- Cradock, J.M., Colgate, M.G., 1974: The use of eigenvectors for smoothing and prediction. *Bull. of the Inst. of Math. and its Appl.*, 5-6.
- Czelnai R., Gandin, L.S., Zaharev, V.I., 1976: Sztatiszticeszkaja sztruktura meteorologiceszkikh polej. Orsz. Meteorológiai Szolgálat, Budapest.
- Dujceva, M.A., Pegy, D.A., 1970: O mnogoparametriczeszkoj szheme prognoza anomalii szrednej meszjacsoj temperaturi vozduha. *Trudi GNC*, 64.
- Epstein, B.S., 1969: Stochastic dynamic prediction. *Helv.*, 6., 739-759.
- Faragó, T., 1973: K posztanovke zadaci o predszkazivanyii sz pomoscsju konyecsnovo verojatnosztnovo avtomata. *Vies. tehnika i voproszi kibernetiki*, 9., 49-61.
- Faragó, T., 1974: Kresenyiju zadaci o predszkazivanyii sz pomoscsju konyecsnovo verojatnosztnovo avtomata. *Vies. tehnika i voproszi kibernetiki*, 11., 109-150.
- Faragó T., Gulyás O., 1973a: Some method of expansion type feature extraction in pattern recognition. *Proc. of Prague Symp. on Asymptotic Statistics*, 69-87.
- Faragó T., Gulyás O., 1973b: The feature extraction in pattern recognition. *Proc. of Conf. on Information Theory, Tallin.*

- Faragó T., Gulyás O., 1974: A regressziós típusú extrapoláció. Meteorológiai Tanulmányok, 2., Országos Meteorológiai Szolgálat, Budapest.
- Faragó T., Gulyás O., 1975: A valószínűségi sűrűségfüggvény becslés és illesztése. Meteorológiai Tanulmányok, 3., Országos Meteorológiai Szolgálat, Budapest.
- Faragó T., 1975: A mintatér bizonyos transzformációról. Problemi peredacsi informacii, 1., 102-107.
- Faragó T. et al., 1975: Az analógia elvén alapuló prognosztikai módszerek matematikai modellje. Időjárás, 3., 166-176.
- Faragó T., 1975: Egy kombinált regressziós becslésről. Időjárás, 6., 366-371.
- Faragó T. et al., 1975: A hosszabb érvényességi idejű prognosztikákkal kapcsolatos kutatások és azok számítástechnikai vonatkozásai, "Az OMSz adatfeldolgozó és adattároló rendszer alapjainak kialakítása R-50-kategóriájú számítógépre" c. OMTB-tanulmányban.
- Faragó T., Kaba M., 1976: Sztatiszticeszkij prognoz meszjacsnih srednyih temperatur i ocenka Bajeszova risska. Időjárás, 6., 313-325.
- Faragó et al., 1976: Távprognosztikai adatbank és kezelőrendszer jelene, az alkalmazói programokkal való kapcsolatai és általános kiépítésének alapvető szempontjai. OMTB-tanulmány, Orsz. Meteorológiai Szolgálat, Budapest.
- Faragó T., 1977: Introduction of characteristic function in stochastic-dynamic prediction. beadva a Tellus-hoz.
- Faragó T., 1977: On the estimation of the probability distribution of maximum values and the statistical analysis of a wind velocity sample. Időjárás, 1., 27-39.

- Fleming, R.J., 1971: On stochastic dynamic prediction. Monthly Weather Rev., 11., 851-872.
- Freibenger, W., Grenander, U., 1965: On the formulation of statistical meteorology. Rev. of Intern. Stat. Institute, 1., 59-68.
- Fritz J., 1971: A Karhunen-Loève sorfejtésről. Szemináriumi Közl., Távközl. Kutató Intézet, Budapest.
- Fritz J., 1975: Distribution free exponential error bound for nearest neighbour pattern classification. IEEE Trans. on Information Theory, 1.
- Fu, K.S., 1968: Sequential methods in pattern recognition and machine learning. Acad. Press, New York.
- Gandin, L.S., 1967.: O primenyeniyi metoda kanonicheskikh korrelatsij v meteorologii. Trudi GGO, 203., Leningrad.
- Gates, G.W., 1972: The reduced NN-rule. IEEE Trans. on Information Theory, 19.
- Glahn, H., 1968: Canonical correlation and its relationship to discriminant analysis and multiple regression. J. Atmos. Sci., 1.
- Gleeson, T.A., 1966: A causal relation for probabilities in synoptic meteorology. J. Appl. Meteorology, 5., 365-368.
- Hart, P.E., 1968: The condensed NN-rule. IEEE Trans. on Information Theory, 15.
- Hershfield, D.M., 1973: On the probability of extreme rainfall events. Bull. Amer. Meteorological Soc., 1013-1022.

- Hess, P., Brazowsky, H., 1952: Katalog der Grosswetterlagen Europas. Berichte des Deutschen Wetterdienstes in der US Zone, 33.
- Ibragimov, I. A., Rozanov, J. A., 1970: Gausszovszkije szlucsajniye processzi. Nauka, Moszkva.
- Jefimov, A. A., 1972: K razscesetu obobszsennik szfericeszkik funkciij. Trudi GGO, 380.
- Judin, M. I., 1967: Fiziko-sztatistsziceszkije metodi prognozov pogodi i vozmozsnosztyi ih vnyedrenyija. Meteorologija i Gidrologija, 11.
- Hald, A., 1967: Statistical theory with engineering applications. Wiley, New York.
- Kapovits R. et al., 1975: A légköri folyamatok szimulációján alapuló közléptávu előrejelzési modell. Reprint. Orsz. Meteorológiai Szolgálat, Budapest.
- Koppány Gy., 1974: Az analógiák extrapolációjának felhasználása a havi és évszakos előrejelzések készítéséhez. Időjárás, 2., 97-108.
- Lehmann, E. L., 1966: Testing Statistical Hypotheses. Wiley, New York.
- Mackenhauer, B., 1974: On the present state of spectral methods in numerical integrations of global atmospheric models. Rep. Intern. Symp. on Spectral Methods in Numerical Weather Prediction; the GARP Programme on Num. Exp., 7.
- Nyberg, A., 1975: An experiment in forecasting monthly mean temperature in Stockholm. Tellus, 1.

- Féczely, Gy., 1957: Grosswetterlagen in Ungarn. Orsz. Meteorológiai Szolgálat Kisebb Kiadványai, 30., Budapest.
- Prókopa, A., 1972: Valószínűségelmélet. Műszaki Kiadó, Budapest.
- Rényi, A., 1970: Valószínűségszámítás. Tankönyvkiadó, Budapest.
- Rosenblatt, M., 1956: A central limit theorem and the strong mixing condition. Proc. Nat. Acad. Sci. USA,
- Smagorinsky, J., 1969: Problems and promises of deterministic extended range forecasting. Bull. Amer. Meteor. Soc., 5., 286-312.
- Szonyecskin, D.M., 1976: Obosznovenyije sztatisztiki v meteorologii. Trudi GMC, 181.
- Tatarskij, V.I., 1969: Iszledovanyije dinamiczeszkij uravnyenij pri verojatnosznom prognoze bariceszkovo polja. Izv. AN SzSzsR, Fizika atmoszferi i okeana, 3.
- Ter-Mkrtcsan, M.G., 1970: O primenyeniji diszkriminantnovo analiza dlja ulucsenyija sztatisticzeszkij prognozov po metoda mnozsosztvennoj regresszii. Trudi GMC, 64.
- Thom, H.C.S., 1966: Some methods of climatological analysis. WMO Techn. Note, 81.
- Vincze, I., 1975: Matematikai statisztika. Műszaki Kiadó, Budapest.
- Wagner, T.J., 1971: Convergence of the NN-rule. IEEE Trans. on Information Theory, 5.
- Watanabe, S., 1972: Karhunen-Loève expansion and factor analysis. Trans. of the Fourth Prague Conf on Inf. Theory.

WMO Techn. Note, 71.,1966: Statistical analysis and prognosis in meteorology.

WMO Techn. Note, 74.,1967: Data processing by machine methods.

Faragó T., Györfi L.,1974: Error distortion function in two - hypothesis decision problem. Proc. of IFAC Conference, Budapest.

Faragó T., Györfi L.,1975: On the continuity of the error distortion function for multiple hypothesis decisions. IEEE Trans. on Information Theory, 6., 458-460.