

NĚKTERÉ NOVÉ PŘÍSTUPY K ANALÝZE ČASOVÝCH ŘAD

T. Cipra, MFF UK

1. Úvod

Přestože se v názvu tohoto příspěvku mluví o nových přístupech k analýze časových řad, je toto označení velice relativní, neboť daná disciplína se velice rychle rozvíjí a její postupy jsou neustále vylepšovány, modifikovány či překonávány novými lepšími metodami. Střízlivý odhad počtu původních prací věnovaných této problematice je 200 až 300 ročně, přičemž sem nejsou zahrnuty početné výzkumné zprávy a některé zajímavé práce čistě aplikačního charakteru. Od r. 1979 vychází časopis Journal of Time Series Analysis věnovaný výlučně časovým řadám. Byla také založena společnost TSA & F (Time Series Analysis & Forecasting), která pravidelně pořádá různé konference a kurzy z oboru a vydává informativní brožuru TSA & F News. Na druhé straně tento rychlý rozvoj má také některé negativní stránky, jako je propagace nových metod bez jejich dostatečného numerického ověření nebo chaotické používání nejrůznějších zkratek a akronymů typu DARM, EWMA, FFT, TAR aj. (viz Granger (1982)).

Z předchozího by mělo být zřejmé, že v příspěvku tohoto typu lze postihnout (a to jen velice schematicky) pouze úzký výsek z daného oboru. Konkrétně se zde zaměříme na Boxovu-Jenkinsovu metodologii (pro jednorozměrné časové řady), která zůstává stále ve středu kritického zájmu statistiků zabývajících se analýzou časových řad, takže v jejím rámci se stále intenzivně pracuje jak po stránce teoretické, tak po stránce softwarové.

Uspořádání článku je následující: v 1. kapitole jsou připomenuty základní principy klasické Boxovy-Jenkinsovy metodologie, zatímco v kapitolách 2-8 jsou uvedeny některé nové směry, v nichž se dnes tato metodologie rozvíjí. Obsah těchto kapitol detailně rozpracován a konkrétními numerickými příklady bude možné nalézt v učebnici Cipra (1986). Zároveň 9. kapitola je věnována robustním metodám v časových řadách.

2. Klasická Boxova-Jenkinsova metodologie

Hlavní myšlenky této metodologie byly zformulovány ve známé monografii Boxe a Jenkinse (1970). Na rozdíl např. od většiny dekompozičních metod zpracovává Boxova-Jenkinsova metodologie řady se závislými pozorováními a dokonce těžiště jejich postupu spočívá právě ve vyšetřování těchto závislostí neboli v tzv. korelační či kovarianční analýze (COVA analýze). Na první pohled by se mohlo zdát, že posornost věnovaná v B.-J. metodologii náhodné složce je přehnaná a že ztrácíme možnost modelovat systematické složky řady, jako je trend nebo sezónní složka. Ale i tyto složky je schopna B.-J. metodologie zvládnout pomocí tzv. modelů ARIMA nebo sezónních ARMA modelů, a to daleko flexibilněji než při dekompozičním přístupu založeném např. na regresní analýze. K hlavním výhodám B.-J. metodologie patří

- flexibilita a adaptivnost na změny v charakteru modelovaného procesu;
- systematický způsob výstavby modelu vhodný pro algoritmické zpracování;
- existence již bohatého a obecně rozšířeného softwaru (např. součást statistických programových systémů SAS, IMSL, BMDP, PACK);
- přibývající pozitivní praktické výsledky.

Na druhé straně nelze zastírat, že tato metodologie má také závažné nedostatky, jako je např.

- pořadavek dostatečné délky řady (obvykle minimálně 50 pozorování);

- značná časová, finanční a softwarová náročnost praktické implementace včetně nemáloch nároků na zkušenosť statistika, který provádí analýzu;
- často nemožnost jednoduché interpretace výsledků, přestože formálně se dosahuje velmi dobrá shoda se skutečností.

Ve stacionárním případě, kdy se momenty procesu do druhého řádu nemění a průběhem času, se provádí výstavba modelu podle B.-J. metodologie ve třech fázích, které nyní stručně připomeneme; na závěr této kapitoly se také zmíníme o nestacionárním případě.

2.1. Identifikace modelu

V rámci identifikační fáze je nutné rozhodnout, zda pro analyzovaný úsek řady x_1, \dots, x_n je vhodné použít autoregresní model AR(p), model klouzavých součtů MA(q) nebo smíšený model ARMA(p,q) včetně stanovení jejich řádu p a q; tyto modely mají tvar postupně

$$(2.1) \quad AR(p) : \quad x_t = \varphi_1 x_{t-1} + \dots + \varphi_p x_{t-p} + \varepsilon_t,$$

$$(2.2) \quad MA(q) : \quad x_t = \varepsilon_t + \theta_1 \varepsilon_{t-1} + \dots + \theta_q \varepsilon_{t-q},$$

$$(2.3) \quad ARMA(p,q) : \quad x_t = \varphi_1 x_{t-1} + \dots + \varphi_p x_{t-p} + \varepsilon_t + \theta_1 \varepsilon_{t-1} + \dots + \theta_q \varepsilon_{t-q}.$$

S výhodou se používá zápis pomocí operátora zpětného posunu B definovaného jako

$$(2.4) \quad B x_t = x_{t-1} \quad (B^k x_t = x_{t-k}),$$

kdy např. model (2.3) lze zapsat jako

$$(2.5) \quad \varphi(B)x_t = \theta(B)\varepsilon_t,$$

kde

$$(2.6) \quad \varphi(B) = 1 - \varphi_1 B - \dots - \varphi_p B^p, \quad \theta(B) = 1 + \theta_1 B + \dots + \theta_q B^q.$$

Je vidět, že základní stavební jednotkou modelů B.-J. metodologie je tzv. bílý šum ε_t , což je posloupnost nekorelovaných náhodných veličin s nulovou střední hodnotou a konstantním kladným rozptylem σ_ε^2 (většinou se také předpokládá, že má normální rozdělení). Modely jsou bez střední hodnoty μ procesu x_t , tj. předpokládá se, že řada je nejprve centrována odečtením aritmetického průměru $\bar{x} = \sum x_t / n$.

Identifikační fáze je plně založena na COVA analýze, která jako hlavní nástroje používá autokorelační funkci

$$(2.7) \quad \varphi_k = \text{corr}(x_t, x_{t+k}) = \frac{\gamma_k}{\gamma_0} = \frac{E(x_t - \mu)(x_{t+k} - \mu)}{\sigma_x^2}$$

a parciální autokorelační funkci φ_{kk} , která je definována jako parciální korelace x_t a x_{t+k} při pevných $x_{t+1}, \dots, x_{t+k-1}$. Pro praktickou identifikaci však máme k dispozici jen výběrové verze předechozích veličin, totiž

$$(2.8) \quad r_k = \frac{c_k}{c_0} \quad (c_k = \sum_{t=1}^{n-k} (x_t - \bar{x})(x_{t+k} - \bar{x}) / n),$$

$$(2.9) \quad r_{kk} = (r_k - \sum_{j=1}^{k-1} r_{k-j} r_{k+j}) / (1 - \sum_{j=1}^{k-1} r_{k-j} r_j) \quad (r_{kj} = r_{k-j} r_{k+j}, r_{11} = r_1).$$

Doporučuje se, aby $n > 50$ a $k < n/4$. COVA analýza umožňuje identifikovat příslušný model díky tomu, že jednotlivé modely (2.1)-(2.3) mají zcela specifické průběhy φ_k a φ_{kk} ; např. pro AR(p) je $\varphi_{kk} = 0$ při $k > p$ a pro MA(q) je $\varphi_k = 0$ při $k > q$ (o případné nulovosti výběrových hodnot r_k či r_{kk} je nutné rozhodnout na základě jednoduchého statistického testu využívajícího páry spolehlivosti pro odhady r_k a r_{kk}).

Zároveň se v identifikační fázi obvykle konstruuje počáteční odhad parametrů $\gamma_1, \dots, \gamma_p, \theta_1, \dots, \theta_q, \sigma_\epsilon^2$, které jsou často dost nepřesné. Např. pro AR(1) lze ze tyto počáteční odhadů vzít

$$\hat{\gamma}_1 = r_1, \quad \hat{\sigma}_\epsilon^2 = \hat{\sigma}_x^2(1 - \hat{\gamma}_1 r_1)$$

($\hat{\sigma}_x^2$ je výběrový rozptyl řady x_1, \dots, x_n).

2.2. Odhad parametrů modelu

Klasický B.-J. přístup k odhadu parametrů je založen na principu maximální věrohodnosti (předpokládá se normalita bílého šumu). V praxi se však většinou používá tzv. metoda nejmenších nelineárních čtverců, která je approximací přesné metody maximální věrohodnosti. Při této metodě se minimalizuje výraz

$$(2.10) \quad S(\varphi, \Theta) = \sum \epsilon_t^2(\varphi, \Theta),$$

kde $\epsilon_t(\varphi, \Theta)$ se počítají rekurentně jako

$$(2.11) \quad \epsilon_t(\varphi, \Theta) = x_t - \varphi_1 x_{t-1} - \dots - \varphi_p x_{t-p} - \theta_1 \epsilon_{t-1}(\varphi, \Theta) - \dots - \theta_q \epsilon_{t-q}(\varphi, \Theta).$$

Konkrétní odhadové metody se pak liší tím, jak je zahájen rekurentní výpočet (lze např. položit $x_0 = x_{-1} = \dots = \epsilon_0(\varphi, \Theta) = \epsilon_{-1}(\varphi, \Theta) = \dots = 0$ a součet v (2.10) brát pro $t=1, \dots, n$) a jaký minimalizační algoritmus je použit.

2.3. Ověření modelu

V této fázi je nutné statisticky ověřit správnost zkonstruovaného modelu (při negativním výsledku je nutné zopakovat identifikační a odhadovou fázi). Bylo navrženo velké množství různých ověřovacích testů. Jeden z nich je tzv. portmanteau test, který pracuje s testovou statistikou

$$(2.12) \quad Q = n \sum_{k=1}^K r_k^2(\hat{\epsilon}),$$

kde K je vhodně zvolené číslo (doporučuje se $K \sim \sqrt{n}$) a $r_k(\hat{\epsilon})$ je výběrová auto-korelační funkce řady vypočtených reziduí $\hat{\epsilon}_t$, které vznikne dosazením odhadnutých parametrů $\hat{\varphi}$ a $\hat{\Theta}$ do (2.11). Pokud se analyzovaná řada řídí modelem ARMA(p,q), pak má Q asymptotické rozdělení χ^2_{K-p-q} (viz také Cipra (1982a)).

2.4. Nestacionární případ

V nestacionárním případě se doporučuje na začátku posoudit, zda není prospěšné řadu nejprve vhodně transformovat. I když Box a Cox (1964) navrhli obecný přístup k takovým transformacím, v praxi se vystačí s transformacemi typu

$$(2.13) \quad \begin{aligned} x_t^{(\lambda)} &= x_t, \quad \lambda \neq 0, \\ &= \ln x_t, \quad \lambda = 0, \end{aligned}$$

ude λ lze stanovit jednoduchou grafickou metodou (viz Jenkins (1979)).

Pokud řada zůstává nestacionární i po této transformaci, navrhli B.-J. použít model ARIMA(p,d,q) (integrated ARMA). Zjednodušeně řečeno, konstruuje se model ARMA, ale už pro řadu, která byla d-krát zdiferencována

$$(2.14) \quad \begin{aligned} \nabla x_t &= x_t - x_{t-1} \quad (\nabla = 1 - B), \\ \nabla^2 x_t &= \nabla(x_t - x_{t-1}) = x_t - 2x_{t-1} + x_{t-2}, \\ &\vdots \\ \nabla^d x_t &= x_t - ({}^d_1)x_{t-1} + ({}^d_2)x_{t-2} - \dots + (-1)^d x_{t-d}. \end{aligned}$$

Konečně pro nestacionaritu ve tvaru sezónního kolísání lze použít v B.-J. metodologii tzv. sezónní modely. Např. při měsíčních pozorováních, kdy "sezóna" obsahuje 12 měření, lze použít model SARIMA(p,d,q) $x(P,D,Q)_{12}$. Např. model SARIMA(0,1,1)x(0,1,1),₁₂ má tvar

$$(1-B)(1-B^{12})x_t = (1+\theta_1 B)(1+\theta_1 B^{12}) \varepsilon_t,$$

$$\text{tj. } x_t - x_{t-1} - x_{t-12} + x_{t-13} = \varepsilon_t + \theta_1 \varepsilon_{t-1} + \theta_1 \varepsilon_{t-12} + \theta_1 \theta_1 \varepsilon_{t-13}$$

(viz také Cipra (1984a)).

3. Sčítavání řádu modelu

Identifikační postup pomocí $\hat{\sigma}_k^2$ a $\hat{\sigma}_{kk}^2$ nebyl zatím zautomatizován natolik, aby nebyla nutná interakce s lidským činitelem. Byly proto hledány identifikační procedury, při nichž by hlavní těža rozhodnutí byla přenesena na počítači. Nejslibnější v tomto směru se jeví kriteria, jejichž minimizaci získáme přímo odhad řádu modelu (viz také Anděl, Cipra (1981)).

Označme $\hat{\sigma}_{k,l}^2$ odhad rozptylu bílého šumu vypočtený za předpokladu, že se řada řídí modelem ARMA(k,l) (např. $\hat{\sigma}_{k,l}^2 = S(\hat{\phi}, \hat{\theta})/n$). Metoda předpokládá, že k dispozici jsou horní hranice pro čísla k a l, tj.

$$(3.1) \quad k \leq K, \quad l \leq L$$

(často se volí K=L=10, ale v literatuře jsou také příklady s K=L=50 nebo K=90, L=0). Přestože s rostoucím k a l má odhad $\hat{\sigma}_{k,l}^2$ tendenci klesat, až při překročení skutečných hodnot p a q kolísá kolem σ_e^2 , nelze doporučit přímou minimizaci $\hat{\sigma}_{k,l}^2$, neboť tak by byly neuděrně preferovány velké hodnoty k a l. Uspěšnější je postup minimalizující

$$(3.2) \quad A_{k,l} = \ln \hat{\sigma}_{k,l}^2 + \ln [1 + w_n(k,l)].$$

Funkce $w_n(k,l)$ je tzv. penalizační funkce, která (i) při pevném n je rostoucí funkce k a l (tj. velké hodnoty k a l opravdu penalizuje), zatímco (ii) při pevném k a l konverguje pro rostoucí n k nule (takže při dostatečně velkém n je pro k < p a l < q s velkou pravděpodobností $A_{k,l} > A_{p,q}$). Konkrétní volba funkce $w_n(k,l)$ pak odlišuje různá kriteria:

1) AIC (název se dnes většinou interpretuje jako Akaike's Information Criterion)

$$(3.3) \quad AIC(k,l) = \ln \hat{\sigma}_{k,l}^2 + 2(k+l)/n.$$

Kriterium odvodil Akaike (1974) na základě jistých principů z teorie informace.

Odpovídající odhad řádu modelu se označuje jako MAICE (Minimal AIC Estimator).

Ozaki (1977) dokonce zobecnil toto kriterium i pro modely ARIMA

$$(3.4) \quad AIC(k,d,l) = \ln \hat{\sigma}_{k,l}^2 + 2(k+l+1+\delta_{d0})/(n-d),$$

kde δ_{d0} je Kroneckerovo delta. Před kriteriem AIC navrhl Akaike (1969) výslovně pro AR modely kriterium FPE (Final Prediction Error), při němž se hledá řád modelu s minimální chybou předpovědi o jeden krok dopředu a které je speciálním případem AIC. Později bylo zjištěno, že AIC není slabě konzistentní a může vést k přečítání řádu modelu.

2) BIC (Bayesian Information Criterion)

$$(3.5) \quad BIC(k,l) = \ln \hat{\sigma}_{k,l}^2 + (k+l) \ln n/n.$$

Toto kriterium nezávisle na sobě odvodili Schwarz (1978) na základě bayesovského

přístupu u Rissanen (1978) na základě informačního přístupu, který souvisejí s hospodárným uložením informace o analyzovaném modelu v paměti počítače. Odhad řádu modelu pomocí BIC je dokonce silně konzistentní.

3) Hannanovo-Quinnovo kriterium

$$(3.6) \quad HQ(k,1) = \ln \hat{\sigma}_{k,1}^2 + c(k+1) \ln \ln n/n.$$

Navrhli je Hannan a Quinn (1979) a v případě AR modelů dokázali jeho silnou konzistence při volbě konstanty $c > 1$. Později byla dokázána jeho silná konzistence i v modelech ARMA při $c > 2$.

4) CAT (Criterion Autoregressive Transfer) - Autorem je Parzen (1974), který toto kriterium odvodil na základě jistých úvah ve spektrální doméně.

5) Metoda maxima χ^2

Navrhl ji McClave (1978a) pro vyšetřování AR modelů, v nichž je jen malý počet parametrů nenulový (tzv. subset autoregression), např. pro model tvaru

$$x_t = 0.6 x_{t-1} - 0.1 x_{t-4} + 0.3 x_{t-12} + \epsilon_t$$

označovaný jako AR (3 max 12) se zpožděními 1, 4 a 12. Později McClave (1978b) tuto metodu rozšířil také na modely MA tak, že obratně použil tzv. inverzní autokorelační funkci.

4. Inverzní autokorelační funkce

Tato funkce představuje nový důležitý nástroj pro identifikaci modelu. Původně ji definoval pomocí pojmu spektrální analýzy Cleveland (1972), později názornější definici v časové doméně předložil Chatfield (1979). Nechť γ_k je autokovarianční funkce uvažované časové řady. Pak inverzní autokovarianční funkce γ_{ik} je definována pomocí vztahu

$$(4.1) \quad \sum_{k=-\infty}^{\infty} \gamma_{ik} z^k = 1 / \sum_{k=-\infty}^{\infty} \gamma_k z^k.$$

Inverzní autokorelační funkce se pak dodefinuje přirozeným způsobem jako

$$(4.2) \quad \phi_{ik} = \gamma_{ik} / \gamma_{i0}.$$

Atraktivnost ϕ_{ik} spočívá v následujícím pravidlu duality: hodnoty ϕ_{ik} v modelu ARMA(p,q) tvaru $\varphi(B)x_t = \Theta(B)\epsilon_t$ odpovídají hodnotám γ_k v modelu ARMA(q,p) tvaru $\Theta(B)x_t = \varphi(B)\epsilon_t$ a naopak. Proto ϕ_{ik} při identifikaci modelu supluje funkci γ_{kk} a v jistých směrech ji ještě předčí (např. při vyšetřování AR modelů s malým počtem nenulových parametrů). Má také význam při určování počátečních hodnot parametrů v modelu ARMA(p,q), neboť známému vztahu

$$(4.3) \quad \gamma_k - \varphi_1 \gamma_{k-1} - \dots - \varphi_p \gamma_{k-p} = 0, \quad k > p$$

odpovídá vztah

$$(4.4) \quad \phi_{ik} + \varphi_1 \phi_{i,k-1} + \dots + \varphi_p \phi_{i,k-p} = 0, \quad k > p.$$

5. Rekurentní metody odhadu parametrů

Rekurentní odhady v časových řadách mají mnoho výhod: především umožňují snadno korigovat již zkonstruované odhady parametrů, jakmile dostaneme k dispozici nová pozorování, a pro svůj algoritmický charakter se snadno implementují na samočinný počítač. Nejpoužívanější z těchto odhadů jsou založeny na principu Kalmanových filtrů (viz Kalman (1960)). Omezme se pro jednoduchost na AR model, který lze napsat ve tvaru

(5.1)

$$x_t = x'_t \varphi + \varepsilon_t,$$

kde $x_t = (x_{t-1}, \dots, x_{t-p})'$, $\varphi = (\varphi_1, \dots, \varphi_p)'$. Pak rekurentní soustava vzorců pro odhady $\hat{\varphi}_t$ a $\hat{\sigma}_{\varepsilon t}^2$ konstruované na základě pozorování x_1, \dots, x_t má tvar (viz Fahrmeier 1981))

$$(5.2) \quad \hat{\varphi}_t = \hat{\varphi}_{t-1} + P_t x_t (x'_t \hat{\varphi}_{t-1}),$$

$$\hat{\sigma}_{\varepsilon t}^2 = \frac{1}{t-p} \left[(t-p-1) \hat{\sigma}_{\varepsilon, t-1}^2 + (x'_t P_{t-1} x_t + 1)^{-1} (x_t - x'_t \hat{\varphi}_{t-1})^2 \right],$$

$$P_t = P_{t-1} - P_{t-1} x_t x'_t P_{t-1} (x'_t P_{t-1} x_t + 1)^{-1},$$

kde P_t je pomocná matice typu $p \times p$. Co se týče počáteční volby parametrů, lze např. volit $\hat{\sigma}_{\varepsilon 0}^2 > 0$ libovolně a $\hat{\varphi}_0 = 0$, $P_0 = c I$ ($c \gg 0$).

6. Ověřovací testy založené na Lagrangeových multiplikátozech

V souvislosti s ověřováním modelu se jen zmíníme o tzv. testech založených na Lagrangeových multiplikátozech (souvisejí s hledáním vázaného extrému), které mají na rozdíl od klasických ověřovacích testů zformulovanou alternativní hypotézu (provádíme např. test AR(1) proti AR(2)). Při tom sympathetic na těchto testezech je fakt, že vystačíme s (maximálně věrohodnými) odhady pořízenými za platnosti nulové hypotézy a není třeba odhadovat model za platnosti alternativní hypotézy (viz např. Godfrey (1979), Poskitt, Tremayne (1980)).

7. Modelování změn v časových řadách

Velká pozornost je v poslední době věnována modelování řad, jejichž charakter se v čase mění. Zde jsou uvedeny tři příspěvky k této problematice :

7.1. Intervenční analýza

Intervenční analýza (viz Box, Tiao (1975), Cipra (1983)) se uplatní v situacích, kdy je průběh řady porušen nějakým jednorázovým zásahem zvenčí (tzv. intervencí, jako je např. výpadek energie, úspěšná reklamní kampaň, změna zákona aj.). Příslušný model má pak dát odpovědi na otázky typu : jak velká je intenzita intervence?, je intenzita intervence konstantní nebo se mění? aj. Princip intervenční analýzy dosudatečně objasní následující příklad :

Analyzuje se měsíční údaje o zněčištění ovzduší výfukovými plyny v jistém velkoměstě od ledna 1955 do prosince 1972. Redukce tohoto znečištění se očekávala od následujících dvou intervenčních zásahů:

I_1 : od ledna 1960 otevření vnějšího dálničního okruhu;

I_2 : od ledna 1966 technologická úprava motorů nově vyráběných aut.

Situaci lze pak zachytit následujícím modelem

$$(7.1) \quad x_t = \gamma_1 S_{t1} + \frac{\gamma_2}{1-B^{1/2}} S_{t2} + \frac{\gamma_3}{1-B^{1/2}} S_{t3} + \varepsilon_t,$$

kde $\gamma_1, \gamma_2 = \gamma_3$ jsou parametry, S_{t1} je skoková proměnná vztahující se k I_1 , a S_{t2}, S_{t3} jsou skokové proměnné vztahující se k I_2 , přičemž

$$S_{t1} = 0, \quad t < leden 1960, \\ = 1, \quad t \geq " ;$$

$$S_{t2} = 1, \quad červen, \dots, říjen od 1966, \\ = 0, \quad jinak;$$

$$S_{t3} = 1, \text{ listopad}, \dots, \text{květen od 1966}, \\ = 0, \text{ jinak.}$$

V případě I_2 je v (7.1) také zachycen fakt, že motory znečišťují ovzduší více v letních měsících než v zimních. Přitom např. pro

$$(7.2) \quad U_{t2} = \frac{\gamma_2}{1-B} S_{t2}$$

je $U_{t2} - U_{t-12,2} = \gamma_2 S_{t2}$, což při $\gamma_2 < 0$ vyjadřuje skutečnost, že podíl nových vozů se zvětšuje každý rok o konstantní hodnotu, takže znečištění klesá o konstantní hodnotu.

1.2. Detekce časových změn v modelu

Zmíníme se zde o metodě CUSUM (Cumulative Sums), což je sekvenční metoda, která ve své původní verzi (viz Page (1955)) je určena pro testování změn v nulové úrovni řady v neznámých okamžicích (je to modifikace Waldova sekvenčního testu). Používá se zde testová statistika

$$(7.3) \quad T_t = \max_{r=1, \dots, t} \{S_r - \min_{i < r} S_i\}$$

kde S_r je kumulativní součet tvaru

$$(7.4) \quad S_r = \sum_{i=1}^r x_i.$$

Hypotéza o nulové úrovni se zamítá, jakmile T_t překročí kritickou hodnotu zjištěnou na základě jistých asymptotických úvah. Později Bagshaw a Johnson (1977) metodu rozšířili na testování změn parametrů v modelech ARMA.

1.3. Modely s časově závislými parametry

Jedná se např. o AR modely obecného tvaru

$$(7.5) \quad x_t = \varphi_1(t)x_{t-1} + \dots + \varphi_p(t)x_{t-p} + \varepsilon_t,$$

kde $\varphi_1(t), \dots, \varphi_p(t)$ jsou funkce času. Jeden z úspěšných přístupů k takovým modelům je založen na teorii evolučních spekter (viz Hussain, Subba Rao (1976)). V poslední době se pozornost také věnuje speciálnímu případu (7.5), kdy $\varphi_1(t), \dots, \varphi_p(t)$ jsou periodické funkce času s periodou d , neboť v tomto tvaru lze vyjádřit každý d -rozměrný autoregresní model a vlastně tak řešit problém odhadu parametrů ve více-rozměrném modelu (pro AR modely viz např. Pagano (1978), Anděl (1983) a pro MA modely viz Cipra (1984b)).

8. Nelineární modely

Další silná tendence v současné teorii časových řad spočívá v přechodu k nelineárním modelům. I když je, samozřejmě, výstavba takových modelů daleko komplikovanější než v lineárním případě, přináší nelineární modely řadu výhod (viz Cipra (1982b)), jako je např. konstrukce přesnějších předpovědí nebo možnost modelovat v diskrétním čase jevy typické pro fyzikální vibrace (např. závislost frekvence na amplitudě). Uvedeme nyní některé důležité typy nelineárních modelů časových řad:

1) Bilineární modely

$$(8.1) \quad x_t = \sum_{n=1}^p \sum_{m=1}^q \beta_{mn} x_{t-n} \varepsilon_{t-m} + \varepsilon_t.$$

Např. tzv. diagonální bilineární model tvaru

$$(8.2) \quad x_t = \beta x_{t-1} \varepsilon_{t-1} + \varepsilon_t$$

je stacionární, právě když $\lambda^2 < 1$. ($\lambda = \beta \sigma_\epsilon$), jeho autokorelační funkce ρ_k má COVA strukturu jako v případě MA(1), zatímco autokorelační funkce $\gamma_k^{(2)}$ procesu x_t^2 má COVA strukturu jako v případě ARMA(1,1) (viz Granger, Anderson(1978)).

2) Exponenciální autoregresní modely. Příkladem je model (viz Haggan, Ozaki(1981))

$$(8.3) \quad x_t = (\gamma_1 + \pi_1 \exp\{-\gamma x_{t-1}^2\})x_{t-1} + \dots + (\gamma_p + \pi_p \exp\{-\gamma x_{t-1}^2\})x_{t-p} + \epsilon_t.$$

Konkrétně např. model

$$(8.4) \quad x_t = (1.5 + 0.28 \exp\{-x_{t-1}^2\})x_{t-1} - 0.96 x_{t-2} + \epsilon_t$$

generuje řadu, jejíž frekvence se s rostoucí amplitudou zvětšuje.

3) Prahové modely. Jedná se o procesy, které mění svůj charakter při překročení jisté úrovni, např.

$$(8.5) \quad \begin{aligned} x_t - a_1 x_{t-1} - \dots - a_p x_{t-p} &= \epsilon_t^{(1)}, & x_{t-d} \leq p_1, \\ x_t - b_1 x_{t-1} - \dots - b_r x_{t-r} &= \epsilon_t^{(2)}, & p_1 < x_{t-d} \leq p_2, \\ x_t - c_1 x_{t-1} - \dots - c_s x_{t-s} &= \epsilon_t^{(3)}, & p_2 < x_{t-d}. \end{aligned}$$

Tento model se označuje jako SETAR($j; p, r, s$)d (Self Exciting Threshold Autoregressive). Pro výstavbu těchto modelů byla navržena metodologie založená na mnohonásobném použití některého kriteria pro odhad řady modelu z 3.kapitoly (viz Tong, Lim (1980), Anděl, Cipra(1982)).

4) Asymetrické modely. Např. Wecker (1981) definuje asymetrický MA(q) model jako

$$(8.6) \quad x_t = \epsilon_t^+ + \theta_1^+ \epsilon_{t-1}^+ + \dots + \theta_q^+ \epsilon_{t-q}^+ + \epsilon_t^- + \theta_1^- \epsilon_{t-1}^- + \dots + \theta_q^- \epsilon_{t-q}^-$$

kde $\epsilon_t^+ = \max\{0, \epsilon_t\}$, $\epsilon_t^- = \min\{0, \epsilon_t\}$ (při $\theta_i^+ = \theta_i^-$, $i=1, \dots, q$, zřejmě došlo klasický MA(q) model).

9. Robustní metody v časových řadách

9.1. Odlehlá pozorování

V konkrétních časových řadách z praxe jsou velice často přítomna tzv. odlehlá pozorování (outliers), která nezypadají do charakteristického průběhu řady a mohou často až drasticky narušit klasické statistické procedury, které s jejich výskytem nepočítají.

Jako jednoduchý příklad lze uvést model AR(1)

$$(9.1) \quad x_t = \gamma_1 x_{t-1} + \epsilon_t, \quad \sigma_x^2 = 1,$$

kde však místo x_t pozorujeme "znečištěné" hodnoty

$$(9.2) \quad y_t = x_t + v_t,$$

přičemž

$$(9.3) \quad v_t \sim \text{iid } 0.9 \delta_0 + 0.1 N(0, 10), \quad P(\delta_0 = 0) = 1$$

(tj. v jedné desetině celkového počtu pozorování stoupne směrodatná odchylka pozorování, která máme k dispozici, přibližně třikrát). Není těžké ukázat, že pak klasický odhad parametru γ_1 metodou nejménších čtverců má přibližně 50% asymptotické vychýlení (takže např. při teoretické hodnotě $\gamma_1 = 0.8$ můžeme obdržet odhad kolem 0.4).

Taková odlehlá pozorování se zpracovávají v rámci statistické analýzy časových řad nejrůznějšími způsoby :

- pomocí subjektivních metod, jako je např. optická prohlídka řady a nahrazení odlehlých pozorování interpolovanými hodnotami;
- nalezení odlehlých pozorování pomocí různých statistických testů, jejich vynechání a pak použití tzv. modelů s chybějícími pozorovánimi (viz např. Kaukenas(1983));
- ponechání původních pozorování řady, ale použití robustních odhadových metod, které respektují rušivý vliv odlehlých pozorování.

V této kapitole se omezíme na poslední z uvedených případů. K detailnějšímu studiu v tomto směru lze doporučit některé práce R.D. Martina (velice instruktivní je např. článek Denby, Martin (1979) o robustních metodách odhadu v AR(1)) a přehledovou zprávu Stockinger, Dutter (1983).

V oblasti časových řad se vystačí s následujícími typy robustnosti:

- kvalitativní robustnost, kdy malé změny v rozdělení procesu vyvolají malé změny v rozdělení odhadu stejnoměrně pro různé délky analyzovaného úseku řady;
- robustnost_v_eficienci, kdy eficiency (vydatnost) odhadu zůstává vysoká v určitém okolí výchozího rozdělení procesu.

Vyšetřují se přitom především následující dva typy modelů časových řad s odlehlými pozorováními:

- modely_s_vnitřními_odlehlými_pozorovánimi_I0 (Innovation Outliers) tvaru

$$(9.4) \quad \begin{aligned} \varphi(B)(x_t - \hat{x}_t) &= \theta(B) \varepsilon_t, \\ \varepsilon_t \sim \text{iid } F &= (1-\gamma)N(0, \sigma^2) + \gamma G, \\ \text{var } G &> \sigma^2 \end{aligned}$$

(odlehlá pozorování zde tedy vznikají tak, že ve zlomku pozorování γ se náhle zvětší rozptyl bílého šumu; nepravidelnosti jsou tedy generovány jakoby uvnitř modelu a proces x_t stále vyhovuje základní dynamické rovnici modelu ARMA);

- modely_s_vnějšími_odlehlými_pozorovánimi_AO (Additive Effects Outliers) tvaru

$$(9.5) \quad \begin{aligned} y_t &= x_t + v_t, \\ x_t \sim \text{ARMA}(\varepsilon_t \sim N(0, \sigma^2)), \\ v_t \sim \text{iid } H &= (1-\gamma) \delta_0 + \gamma G, \quad P(\delta_0 = 0) = 1 \end{aligned}$$

(obvykle se předpokládá, že rozdělení G je absolutně spojité, takže $P(v_t = 0) = 1 - \gamma$; při tomto typu modelu je tedy s pravděpodobností $1 - \gamma$ pozorován vlastní ARMA proces x_t a s pravděpodobností γ proces x_t zatížený chybou s rozdělením G).

2.2. Odhady metodou nejmenších čtverců

První otázka, kterou je nutné v rámci robustního přístupu k časovým řadám zodpovědět, zní, jak se s předchozími dvěma typy modelů vyrovnávají klasické odhady v časových řadách, tj. především odhady metodou nejmenších nelineárních čtverců (viz(2.10)), které se často označují jako LS odhady (Least Squares).

V běžném modelu ARMA, v němž $\varepsilon_t \sim \text{iid } N(0, \sigma^2)$, jsou všechny odhady $\hat{\varphi}_{LS}$, $\hat{\theta}_{LS}$, $\hat{\mu}_{LS}$ a $\hat{\sigma}_{LS}^2$ konzistentní, asymptoticky normální a asymptoticky eficientní. Odhady $\hat{\varphi}_{LS}$ a $\hat{\theta}_{LS}$ však zůstanou konzistentní i v některých případech s bílým šumem bez konečného rozptylu, např. když ε_t mají symetrické stabilní rozdělení F s indexem $0 < \alpha \leq 2$ (tj. $\int_{\mathbb{R}} |x|^{\alpha} dF(x) = \exp\{-c|\alpha|\}, c > 0$).

V modelu I0 jsou $\hat{\varphi}_{LS}$ a $\hat{\theta}_{LS}$ kvalitativně robustní (jejich asymptotická rozptylová matici nezávisí na F), zatímco $\hat{\mu}_{LS}$ a $\hat{\sigma}_{LS}^2$ již tuto vlastnost

nemají (neboť jejich asymptotický rozptyl závisí podstatně na σ_ϵ^2). Žádny z odhadů $\hat{\gamma}_{LS}$, $\hat{\theta}_{LS}$, $\hat{\mu}_{LS} = \hat{\sigma}_{LS}^2$ není robustní v eficienci: např. pro AR(1) je

$$(9.6) \quad \text{eff}(\hat{\gamma}_1, LS) = \frac{\text{var}_{CR} \hat{\gamma}_1}{\text{var} \hat{\gamma}_1, LS} = \frac{(1 - \hat{\gamma}_{1, LS}^2) / (\sigma_\epsilon^2 i(f))}{1 - \hat{\gamma}_{1, LS}^2} = \frac{1}{\sigma_\epsilon^2 i(f)},$$

kde $\text{var}_{CR} \hat{\gamma}_1$ je dolní Cramerova-Reedová hranice pro rozptyl odhadu $\hat{\gamma}_1 = \hat{\sigma}_\epsilon^2$ a $i(f)$ je rozptyl a Fisherova informace rozdělení F s hustotou f . Přitom σ_ϵ^2 může být libovolně velké v libovolně malém okolí $N(0, \sigma_\epsilon^2)$, zatímco $i(f)$ zůstává relativně stabilní.

Konkrétně v AO modelu nejenomže LS odhady nejsou robustní, ale jsou dokonce silně asymptoticky vychýlené. LS odhady jsou proto pro AC modely velice nevhodné.

2.3. M-odhad

Většina robustních odhadů navržených pro modely časových řad je typu M (včetně různých modifikací). Důvodem je patrně fakt, že M-odhad lze považovat za zobecněné maximálně věrohodné odhady, které se v klasických odhadových procedurách pro ARMA modely nejvíce používají (buď přímo nebo po sproximaci jako LS odhady). Např. pro IO model AR(p) přepsaný do tvaru

$$(9.7) \quad x_t = \lambda + \varphi_1 x_{t-1} + \dots + \varphi_p x_{t-p} + \epsilon_t = \beta' z_t + \epsilon_t,$$

kde $\lambda = \mu(1 - \varphi_1 - \dots - \varphi_p)$, $\beta = (\lambda, \varphi_1, \dots, \varphi_p)'$, $z_t = (1, x_{t-1}, \dots, x_{t-p})'$, lze použít Huberovu metodu M-odhadu pro regresní modely. Speciálně je nutné řešit soustavu rovnic (v neznámých $\hat{\beta}_M = \hat{\sigma}_M$)

$$(9.8) \quad \sum_{t=p+1}^n \psi\left(\frac{x_t - \hat{\beta}'_M z_t}{\hat{\sigma}_M}\right) z_t = 0$$

$$(9.9) \quad \frac{1}{B} \sum_{t=p+1}^n \psi^2\left(\frac{x_t - \hat{\beta}'_M z_t}{\hat{\sigma}_M}\right) = B,$$

kde $B = \int \psi^2(x) d\Phi(x)$ (Φ je distribuční funkce $N(0,1)$) a ψ je funkce známá z obecné teorie M-odhadů.

Pro vlastní výpočet M-odhadu je doporučován tzv. IWLS algoritmus (Iterated Weighted Least Squares).

V IO modelech jsou M-odhady konzistentní, asymptoticky normální, robustní v eficienci, ale nejsou kvalitativně robustní (to však lze považovat v jistém smyslu za jejich klad, neboť lze snadno ukázat, že pro rozdělení F s těžkými konci jsou M-odhady přehánějící než LS odhady). Naproti tomu v AO modelech mají M-odhady nepříznivé vlastnosti LS odhadů v těchto modelech.

2.4. Modifikace M-odhadů

GM odhad (Generalized M Estimator) se snaží napravit nepříznivé vlastnosti M-odhadů v AO modelech. Např. ve svém nejjednodušším tvaru pro model AR(1) vznikají takto takzvané rovnice

$$(9.10) \quad \sum_{t=2}^n w(x_{t-1}) \psi(x_t - \hat{\gamma}_1, GM x_{t-1}) = 0,$$

kde w je omezená funkce. V IO modelech jsou GM odhady konzistentní, mají však horší robustnost v eficienci než M-odhady. V AO modelech jsou však již poměrně robustní a mají zde menší asymptotické vychýlení než LS a M-odhady.

AML odhady (Approximate Maximum Likelihood) se konstruuje tak, že se provede jistá approximace přesné věrohodnostní funkce, aby se podobala funkci pro M odhad.

Zajímavý odhad funkcionální metodou nejmenších čtverců navrhl Heathcote a Welsh (1983)..

2.5. Další robustní analýza modelů časových řad

Je nutné zdůraznit, že robustní přístup k modelům časových řad se neomezuje jen na odhady parametrů. Byly už také např. navrženy

- robustní intervaly spolehlivosti pro odhadnuté parametry;
- testy pro rozlišení IO a AO modelů, v nichž se např. používá testová statistika $V_n(\hat{\gamma}_M - \hat{\gamma}_{GM})$ (neboť GM odhady mají přijatelné vlastnosti v obou typech modelů, zatímco M odhady dominují jen v IO modelech) nebo princip věrohodnostních poměrů (viz Fox (1972));
- robustní kriteria pro stanovení řádu modelu.

Literatura

- [1] Akaike,H. (1969). Fitting autoregressive models for prediction. Ann.Inst.Statist. Math. 21, 243-247.
- [2] Akaike,H.(1974). A new look at the statistical model identification. IEEE Trans.Automat. Contr. AC-19, 716-723.
- [3] Anděl,J. (1983). Statistical analysis of periodic autoregression. Aplikace matematiky 28, 364-385.
- [4] Anděl,J., Cipra,T. (1981). Stanovení řádu autoregresního modelu a jeho využití při odhadu spektrální hustoty. Výzkumná zpráva IHE, Praha.
- [5] Anděl,J., CiPrá,T.(1982). Přehové modely. Výzkumná zpráva IHE, Praha.
- [6] Bagshaw,M., Johnson, R.A. (1977). Sequential procedures for detecting parameter changes in a time series model. JASA 72, 593-597.
- [7] Box,G.E.P., Cox,D.R.(1964). An analysis of transformations. J.R.Statist.Soc. B 26, 221-252.
- [8] Box,G.E.P., Jenkins, G.M.(1970). Time series analysis, forecasting and control. Holden Day, San Francisco (2.vydání (1976), ruský překlad (1974)).
- [9] Box,G.E.P., Tiao,G.C.(1975). Intervention analysis with application to economic and environmental problems. JASA 70, 70-79.
- [10] Cipra,T. (1982a). Behaviour of the portmanteau statistic for white noise with periodical components. Research Report 82-3, Dept.of Statistics, University of Uppsala.
- [11] Cipra,T.(1982b). Nelineární modely časových řad. Ekonomicko-matematický obzor 19, 164-178.
- [12] Cipra,T.(1983). Intervenční analýza .Výzkumná zpráva IHE, Praha.
- [13] Cipra,T:(1984a). Investigation of periodicity for dependent observations. Aplikace matematiky 29, 134-142.
- [14] Cipra,T.(1984b). Periodic moving average process (zasláno k otištění).
- [15] Cipra,T. (1986). Analýza časových řad s aplikacemi v ekonomii. SNTL, Praha.
- [16] Cleveland, W.S.(1972). The inverse autocorrelations of a time series and their applications. Technometrics 14, 277-293.
- [17] Denby,L., Martin,R.D. (1979). Robust estimation of the first-order autoregressive parameter. JASA 74, 140-146.
- [18] Fahrmeier,L. (1981). Rekursive Algorithmen für Zeitreihenmodelle. Vandenhoeck und Ruprecht, Göttingen.
- [19] Fox,A.J.(1972). Outliers in time series. J.R.Statist.Soc. B 34, 350-363.
- [20] Godfrey,L.G. (1979). Testing the adequacy of a time series model. Biometrika 66, 67-72.
- [21] Granger, C.W.J. (1982). Acronyms in time series analysis (ATSA). J.Time Series Analysis 3, 103-107.
- [22] Granger, C.W.J., Anderson, A.P.(1978). An introduction to bilinear time series models. Vandenhoeck und Ruprecht. Göttingen.

- [23] Haggan, V., Ozaki, T. (1981). Modelling nonlinear vibrations using an amplitude-dependent autoregressive time series models. *Biometrika* 68, 189-196.
- [24] Hannan, E.J., Quinn, B.G. (1979). The determination of the order of an autoregression. *J.R.Statist.Soc. B* 41, 190-195.
- [25] Heathcote, C.R., Welsh, A.G. (1981). The robust estimation of autoregressive process by functional least squares. *J.Appl.Prob.* 20, 737-753.
- [26] Hussain, M.J., Subba Rao, T. (1976). The estimation of autoregressive, moving average and mixed autoregressive moving average systems with time dependent parameters of nonstationary time series. *Intern.J.Control* 23, 647-656.
- [27] Chatfield, C. (1979). Inverse autocorrelations. *J.R.Statist.Soc. A* 142, 363-377.
- [28] Jenkins, G.M. (1979). Practical experiences with modelling and forecasting time series. GPF Publication, Jersey.
- [29] Kalman, R.E. (1960). A new approach to linear filtering and prediction problems. *Trans. ASME J.Basic Eng.* 82, 34-45.
- [30] Kaukenes, J. (1983). Uznávání lokálních něstacionarností v interpolaciích v realizacích stacionárného procesu autoregressi. *Statistické problémy upravení* 61, 9-23.
- [31] McCleave (1978a). Estimating the order of autoregressive models: the max χ^2 method. *JASA* 73, 122-128.
- [32] McCleave (1978b). Estimating the order of moving average models: the max χ^2 method. *Commun.Statist-Theor.Meth.* A7, 259-276.
- [33] Ozaki, T. (1977). On the order determination of ARIMA models. *Appl.Stat.* 26, 290-301.
- [34] Pagan, M. (1978). On periodic and multiple autoregressions. *Ann.Stat.* 6, 1310-1317.
- [35] Page, E.S. (1955). A test for a change in a parameter occurring at an unknown point. *Biometrika* 42, 523-527.
- [36] Parzen, E. (1974). Some recent advances in time series modelling. *IEEE Trans. Automat.Contr.* AC-19, 723-730.
- [37] Poskitt, D.S., Tremayne, A.R. (1980). Testing the specification of a fitted autoregressive-moving average model. *Biometrika* 67, 359-363.
- [38] Rissanen, J. (1978). Modelling by shortest data description. *Automatica* 14, 465-471.
- [39] Schwarz, G. (1978). Estimating the dimension of a model. *Ann.Stat.* 6, 461-464.
- [40] Stockinger, N., Dutter, R. (1983). Robust time series analysis : an overview. Research Report No.9, Technische Universität Graz.
- [41] Tong, H., Lim, K.S. (1980). Threshold autoregression, limit cycles and cyclical data. *J.R.Statist.Soc. B* 42, 245-292.
- [42] Wecker, W.E. (1981). Asymmetric time series. *JASA* 76, 16-21.