

Problema stabilității estimărilor econometrice

Date anormale(outliers), metode robuste de regresie ortogonală

Corina Sâman¹

Introducere

Datele anormale(outliers), prin definiție, pun sub semnul întrebării orice analiză și inferență pe care dorim s-o obținem din structura statistică a oricărui set de date. Sursa acestei anormalități (erori de măsurare, erori datorate abaterii față de structura statistică considerată a datelor) poate fi de interes, dar mai important este ca aceste date să fie ponderate corespunzător sau excluse din setul de date pentru a preveni perturbarea parametrilor estimați și a permite inferențe statistice.

Acestea vor fi prezente mai ales în seturi de date cu multe variabile și/sau multe observații. Cu privire la o distribuție multivariată a datelor, distanța statistică a unei observații prezintă proprietatea că o creștere a distanței de la medie reflectă o descreștere în probabilitate. Datele anormale (outliers) vor fi deci acele observații care sunt anormale în raport cu distribuția lor de probabilitate, având o probabilitate semnificativ mică de apariție. Datele anormale sunt deci un tip de contaminare. Pot exista astfel de date rătăcite în setul de date sau pot exista clusteri cu astfel de puncte. Prezența anomaliilor în date (outliers) face ca încrederea în inferențele statistice bazate pe modelul asociat datelor să fie mică.

Tehnicile de analiză multidimensională (regresii multivariate, componente principale, analiza factorilor, clasificarea, analiza clusterilor, etc.) se bazează pe statistici empirice (medie, covarianță, corelație) și minimizarea unor funcții de reziduri care pot fi afectate sever de existența unor date anormale, fie ele cât de puține. Doua lucruri se întâmplă:

- estimările diferă substanțial de cele obținute dacă nu ar exista date anormale;
- modelul rezultat nu permite detectarea datelor anormale pe baza rezidualelor, distanțelor Mahalanobis sau a altor teste diagnostic.

Două proceduri pot fi folosite:

- estimari robuste: găsirea unui model robust care este similar cu cel ce ar rezulta dintr-un set de date fără anomalii;
- detectarea datelor anormale.

¹ Institutul de Prognoză Economică, Centrul de Modelare Macroeconomică, Academia Română

Ambele răspund problemei, dar pornesc din direcții diferite. În primul caz se găsește modelul potrivit cu datele fără anomalii, ceea ce permite identificarea datelor anormale pornind de la reziduale, iar în cel de-al doilea caz detectarea anomaliilor permite să se trateze acestea prin înlăturare sau orice altă procedură de diminuare a efectului lor asupra estimării.

Ne vom ocupa în primul rând de detectarea anomaliilor în date (puncte numite outliers).

Media unui eșantion și matricea de covarianță empirică stau la baza analizei clasice a datelor, fiind estimatori optimi ai parametrilor de localizare și împrăștiere, dar foarte sensibili la date anormale. Acești estimatori sunt necesari pentru multe tehnici de analiză ca regresia, analiza componentelor principale, analiza factorilor, clasificarea, etc.

În prima parte a studiului am discutat tehnici robuste de estimare a localizării și împrăștierii datelor, iar în partea a doua am aplicat aceste metode pe un caz real cu trei variabile.

În partea a treia am descris o metodă robustă pentru estimarea unui model de regresie ortogonală, iar în secțiunea a patra am făcut simulări Monte-Carlo pentru a compara performanța diferiților estimatori ai modelului de regresie ortogonală robustă.

1. Tehnici de estimare robustă

1.1 Distanțe statistice

Dacă avem un set de date de dimensiune d și presupunem că ele reprezintă observații ale unui vector de variabile aleatoare $\mathbf{X}^T = (\mathbf{X}_1, \dots, \mathbf{X}_d)$, iar o realizare a variabilei aleatoare este $\mathbf{x}^T = (x_1, \dots, x_d)$.

Pentru a analiza acest set de date și a vedea dacă există date anormale, adică date care nu aparțin modelului statistic considerat, este necesar să avem definită o distanță. Doi parametri esențiali pentru definirea distanței sunt media μ și matricea de covarianță Σ a variabilei aleatoare \mathbf{X} .

Putem porni în analiza datelor de la distanța Euclidiană, D , față de un estimator al mediei $\hat{\mu}$:

$D^2 = ((x_i - T)'(x_i - T))$, dar atunci nu am luat în considerare variația setului de date de-a lungul axelor care este dată de un estimator al împrăștierii datelor. O distanță mare față de centru este mai semnificativă pentru un set de date care prezintă o variație mică decât pentru un altul cu variație mare. Distanța trebuie deci să fie determinată în mod invers proporțional de o măsură a împrăștierii datelor. Astfel de distanțe sunt numite distanțe statistice.

Cu privire la o distribuție multivariată a datelor, distanța statistică a unei observații prezintă proprietatea că o creștere a distanței de la medie reflectă o descreștere în probabilitate.

Datele anormale (outliers) vor fi deci acele observații care sunt anormale în raport cu distribuția lor de probabilitate, având o probabilitate semnificativ mică de apariție.

O distribuție multivariat normală de dimensiune d este definită de densitatea de probabilitate:

$$f(x) = \frac{1}{(2\pi)^{d/2}} \frac{1}{|\Sigma|^{1/2}} \exp(-1/2(x - \mu)' \Sigma^{-1}(x - \mu)) \quad (1)$$

unde Σ este matricea de covarianță și $\frac{1}{|\Sigma|^{1/2}}$ este Iacobianul transformării.

Exponentul din această expresie corespunde cu unul dintre exemplele de distanță statistică foarte cunoscute, distanța Mahalanobis definită de acesta în 1930 numită și distanța T^2 a lui Hotelling.

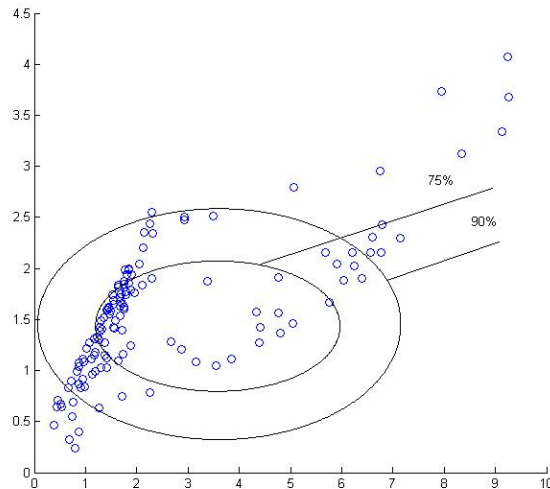
Dacă \mathbf{X} este d -multivariat normală atunci orice punct cu distanță echivalentă față de centrul datelor se găsește pe un elipsoid care descrie mulțimile de nivel ale funcției de densitate de probabilitate.

Pentru date astfel distribuite este cunoscut că pătratul distanței Mahalanobis urmează o distribuție chi-pătrat cu d grade de libertate:

$$M_i = ((x_i - \mu)' \Sigma^{-1}(x_i - \mu))^{1/2} \sim \chi_d^2, \quad (2)$$

unde $P(M_i \leq \chi_{0.90,d}^2) = 0.90$. Astfel încât, dacă $d = 2$, $c^2 = \chi_{0.90,2}^2 = 4.605$ și este 90% șansă ca variabila aleatoare x_i să se găsească în interiorul elipsei descrise de curba $M^2 = c^2$.

O ilustrare a unei astfel de situații pentru un set de date particular este dat în figura următoare. O valoare mai mică pentru $c > 0$ înseamnă că mai puține date vor fi conținute în elipsa corespunzătoare.



Când identificăm datele anormale (outliers) folosind metode bazate pe distanțe localizăm acele observații care sunt în exteriorul unei astfel de elipse. Determinarea frontierei acestei elipse care separă cele două regiuni se face pornind de la estimarea robustă a localizării (centrului) și împrăștierii (covarianței).

Seturile de date pot conține date anormale (outliers), puncte care deviază de la modelul sugerat de majoritatea datelor, așa că o analiză statistică a acestora trebuie să înceapă în mod necesar cu estimarea lor.

1.2 Invarianța Afină și Estimarea de Verosimilitate Maximă

În cazul particular când \mathbf{X} are o densitate de probabilitate multivariată normală (1), estimatorii de verosimilitate maximă pentru eșantionul x_1, \dots, x_n pentru medie și împrăștiere sunt dați de acei parametri $\hat{\mu}$, $\hat{\Sigma}$ care satisfac condiția:

$$L_n(\hat{\mu}, \hat{\Sigma}) = \max_{\mu, \Sigma} L_n(\mu, \Sigma), \quad (3)$$

unde $L_n(\mu, \Sigma)$ este probabilitatea dată de modelul asociat datelor

$$L_n(\mu, \Sigma) = \frac{1}{|\Sigma|^{1/2}} \exp\left(-\frac{1}{2} \sum_{i=1}^n (x_i - \mu)' \Sigma^{-1} (x_i - \mu)\right). \quad (4)$$

Estimatorii sunt:

$$\hat{\mu} = \bar{x} = \frac{1}{n} \sum_{i=1}^n x_i, \quad \hat{\Sigma} = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})'$$

Invarianța estimatorilor la transformările afine este o proprietate naturală făcând analiza independentă de unitatea de măsură și orice translație sau rotație a datelor. Spunem că un estimator $T(\mathbf{X})$ al localizării și un estimator $C(\mathbf{X})$ al covarianței este invariant în sensul precizat dacă îndeplinesc condițiile:

$$\begin{aligned} T(\mathbf{X}\mathbf{A} + \mathbf{1}_n \mathbf{v}') &= T(\mathbf{X})\mathbf{A} + \mathbf{v} \\ C(\mathbf{X}\mathbf{A} + \mathbf{1}_n \mathbf{v}') &= \mathbf{A}'C(\mathbf{X})\mathbf{A} \end{aligned} \quad (5)$$

pentru orice vector \mathbf{v} de dimensiune $d \times 1$ și orice matrice nesingulară \mathbf{A} de dimensiune $d \times d$, unde d este dimensiunea spațiului de date.

Estimatorii de verosimilitate maximă pentru modelul multivariat normal al datelor sunt afin-invarianți.

Totuși acești estimatori pot fi afectați de prezența datelor anormale pe care vrem să le detectăm, diminuând șansele de a le identifica. De aceea este necesar să considerăm variante robuste ale localizării și împrăștierii datelor.

1.3 Funcția de influență

O altă discuție legată de estimatorii robuști privește funcția de influență (Hampel, 1974).

O metodă de a măsura robustețea unui estimator $T(X)$, dacă se presupune că distribuția lui X este F , este de a evalua efectul asupra estimatorului când distribuția reală a datelor este perturbată într-un punct x din \mathbb{R}^d . Fie δ_x măsura de probabilitate exprimată de funcția Dirac care atribuie 1 punctului x . Hampel (1971) a definit funcția de influență pentru $T(X)$ dat de distribuția F într-un punct x :

$$IF(x; T, F) = \lim_{0 < \varepsilon \rightarrow 0} \frac{T[(1 - \varepsilon)F + \varepsilon\delta_x] - T(F)}{\varepsilon} \quad (6)$$

dacă limita există.

Funcția de influență măsoară efectul asupra lui T a unei conatminări infinitezimale a distribuției F .

Estimatorul T care are o funcție de influență mică sau cel puțin mărginită este considerat robust.

Dacă înlocuim ε cu $1/n$ în expresia funcției de influență putem măsura efectul unui outlier pe poziția x .

Funcția de influență are scăderea că măsoară efectul unei singure date anormale (Lopuhaa, 1989).

Pentru media și covarianța clasică avem:

$$IF(x; \mu, F) = x - \mu(F)$$

$$IF(x; \Sigma, F) = (x - \mu)(x - \mu)' - \Sigma(F) \text{ } ^?$$

ceea ce dovedește că sunt nemărginite și deci că estimatorii sunt nerobuști.

1.4 Estimatori robuști pentru localizare și împrăștiere

Bickel (1964) pare să fie primul care a considerat alternative robuste pentru media unui vector: mediana pe fiecare direcție și estimatorul Hodge-Lehman. Acești estimatori fiind mai robuști decât media empirică, totuși nu respectă proprietate de invarianță afină.

Hampel (1973) a fost primul care a propus o procedură iterativă pentru un estimator afin-invariant pentru matricea de împrăștiere (scatter), care s-a dovedit a fi un M-estimator.

Inspirat de Huber (1964), Maronna (1976) a definit și studiat pentru prima oară M-estimator pentru localizare și împrăștiere.

O metoda clasică pentru detectarea datelor anormale (outliers) este să se calculeze distanța Mahalanobis față de centrul spațiului datelor (norul de puncte) pentru fiecare punct:

$$MD_i = ((x_i - T)' C^{-1} (x_i - T))^{1/2},$$

unde T și C sunt respective media și matricea de covarianță a eșantionului.

MD ar trebui să măsoare depărtarea de centrul norului de date, ținându-se cont în același timp de forma acestuia.

Numai că aceasta distanță este influențată de datele anormale (outliers) care pot exista care denaturează măsurile T și C făcând să apară efectul de mascare prin care datele anormale nu primesc neapărat o măsură MD mare.

Robustificarea acestor măsuri se poate face folosind alte marimi pentru T și C , care să aprecieze corect centrul (locația) și împrăștierea. Dacă se consideră estimatori robuști pentru T și C atunci distanța Mahalanobis va fi mare pentru datele anormale și va deveni:

$RD_i = ((x_i - T_R)' C_R^{-1} (x_i - T_R))^{1/2}$, unde indicele R semnaleză o măsură robustă pentru estimator.

Estimatorii robuști pentru localizare și împrăștiere pot fi calculați deodată ca în cazul M -estimatorilor, sau separat.

Invarianța estimatorilor la transformările afine este o proprietate naturală făcând analiza independentă de unitatea de măsură și orice translație sau rotație a datelor:

$$T(XA + 1_n v') = T(X)A + v$$

$$C(XA + 1_n v') = A' C(X) A$$

pentru orice vector v de dimensiune $dx1$ și orice matrice nesingulară A de dimensiune $dx d$, unde d este dimensiunea spațiului de date.

O măsură globală de robustețe a estimatorilor este punctul de rupere (breakdown point).

Donoho și Huber (1983) au considerat definiția acesteia pentru eșantioane finite, iar Hampel (1971) sensul asimptotic.

Acest punct de rupere se poate defini informal ca proporția minimă de contaminare a datelor din mulțimea X care să facă estimatorul T să devieze oricât de mult de la adevărata valoare:

$$BP(T, X) = \min \left\{ \frac{m}{n} : \sup_{x_m} |T(X_m) - T(X)| = \infty \right\} . \quad (6)$$

Unde X_m este o mulțime de date contaminată rezultând din înlocuirea a m puncte din X cu puncte arbitrare din R^d .

Pentru matricea de covarianță se înlocuiește în definiția T cu vectorul valorilor proprii logaritmice ale lui C . Estimatorul lui C ar deveni nefolositor dacă valorile proprii ar fi 0 sau infinit așa cum pentru localizare ar fi infinit de mare.

1.4.1 M-estimatori

M -estimatorii nu sunt influențați de perturbații mici în date și au o eficiență relativ bună pentru multe modele ale unei populații, dar nu sunt destul de robuști ca măsură globală (în termeni de punct de rupere = breakdown point) pentru spațiile cu multe dimensiuni.

M-estimatorii pentru localizare și împrăștierte sunt soluțiile T (în R^d) și C (o matrice simetrică pozitiv-definită) ale ecuațiilor următoare:

$$\frac{1}{n} \sum_{i=1}^n u_1 (((x_i - T)' C^{-1} (x_i - T))^{1/2}) (x_i - T) = 0 \quad (7)$$

$$\frac{1}{n} \sum_{i=1}^n u_2 (((x_i - T)' C^{-1} (x_i - T))^{1/2}) (x_i - T) (x_i - T)' = C \quad (8)$$

unde u_i , $i=1,2$, sunt funcții care satisfac niște condiții:

$u: R \rightarrow [0, \infty)$ este simetrică, $u(0)=0$ și $u(y) \rightarrow \infty$ când $y \rightarrow \infty$.

Sunt generalizări ai estimatorilor de verosimilitudine maximă și pot fi priviți ca media și matricea de covarianță ponderate.

Verosimilitatea maximă pentru distribuțiile normale este dată de alegerea:

$$u_1(s) = -\frac{1}{s} \frac{d(\log f(s))}{ds} \text{ și } u_2(s^2) = u_1(s) \text{ pentru } s > 0.$$

Ei au funcție de influență mărginită pentru anumiți u_i , dar un punct de rupere relativ mic ($< 1/(d+1)$) (Maronna, 1976) și de aceea nu sunt global robuști în spațiile de date cu multe dimensiuni.

1.4.2 Estimatori Stahel-Donoho

Stahel (1981) și Donoho (1982) au introdus ca măsuri pentru localizare și împrăștiere o medie și matrice de covarianță ponderate bazate pe proiecții care au în acelaș timp proprietatea de afin-invarianță și o robustețe globală înaltă dată de punctul de rupere (breakdown point).

Detectarea punctelor anormale se face după o măsură care se bazează pe proiecția unui punct x din mulțimea datelor $X \subset R^d$ ($d \geq 1$ este dimensiunea datelor) în R :

$$O(x, X) = \sup_{\{u \in R^d, \|u\|=1\}} \frac{|u'x - \mu(uX)|}{\sigma(uX)} \quad (9)$$

unde $X = \{u'X_1, \dots, u'X_n\}$, $X = \{X_1, \dots, X_n\}$.

Atunci locația și împrăștierea se calculează astfel:

$$T_{SD}(X) = \sum_{i=1}^n w_i x_i / \sum_{i=1}^n w_i, \quad (10)$$

$$C_{SD}(X) = \sum_{i=1}^n w_i (x_i - T_{SD}(x))(x_i - T_{SD}(x))' / \sum_{i=1}^n w_i, \quad (11)$$

unde $w_i = w(O(x_i, x))$ și w este o funcție pondere care penalizează punctele anormale.

Donoho a demonstrat robustețea acestor estimatori pentru $(\mu, \sigma) = (Med, MAD)$, mediana și median absolute deviation și w funcția pondere cu proprietăți asemănătoare funcțiilor u de la M-estimatori.

1.4.3 Estimatori MVE și MCD

Rousseeuw (1985) a introdus ca estimatori afin-invarianți pentru localizare și împrăștiere cei ce se bazează pe MVE (minimum volume elipsoid) elipsoidul de volum minim și cei ce se bazează pe MCD (minimum covariance determinant).

Estimatorii MVE pentru localizare și împrăștiere sunt respectiv centrul și elipsoidul de volum minim care cuprind cel puțin h puncte din X . S-a dovedit de către Davies (1987) că acești estimatori au un punct de rupere foarte înalt pentru $h = \lfloor (n+d+1)/2 \rfloor$, dar nu sunt asimptotic normali nici \sqrt{n} consistenți (Davies, 1992)

Estimatorii MCD sunt media empirică și un multiplu al matricea de covarianță a h puncte din X pentru care determinantul matricei de covarianță este minim.

TMCD și CMCD au punctul de rupere aproximativ $(n-h)/n$, adică $1/2$ pentru $h = \lfloor (n+d+1)/2 \rfloor$ (Lopuhaa and Rousseeuw, 1991) valoarea maximă pe care o poate atinge orice estimator invariant la transformări afine.

Acești estimatori MCD sunt \sqrt{n} consistenți (Davies and Jhun, 1993) și de asemenea este demonstrată Rousseeuw normalitatea asimptotică pentru localizare dar nu și pentru matricea de covarianță și nu au o eficiență mare pentru modelele multivariat normale (Croux and Haesbroeck, 1999), dar sunt foarte des folosite datorită algoritmilor rapizi de calcul existenți.

Un algoritm eficient de calcul pentru estimatorilor MCD pentru localizare și covarianță este FAST-MCD dezvoltat de Rousseeuw și Van Driessen (1999) care calculează estimatori ponderați (one step reweighted):

$$T_{RMCD}(X) = \sum_{i=1}^n w_i x_i / \sum_{i=1}^n w_i ,$$

$$C_{RMCD}(X) = D_{h,n} \sum_{i=1}^n w_i (x_i - T_{RMCD}(x))(x_i - T_{RMCD}(x))' / \sum_{i=1}^n w_i ,$$

unde $w_i = \begin{cases} 1 & \text{daca} \\ 0 & \text{altfel} \end{cases} RD_{MCD}(X_i) \leq \sqrt{\chi_{d,0.975}^2}$, iar $D_{h,n}$ este un factor de corecție-

1.4.4 S-estimatori

O metodă de a îmbunătăți eficiența estimatorilor MVE și MCD este să se considere o funcție obiectiv netedă. Astfel o clasă importantă de esimatori sunt S-estimatorii

(Rousseeuw and Leroy , 1987, Davies, 1987) pentru localizare și covarianță care minimizează funcția:

$$\min \det(C)$$

$$\frac{1}{n} \sum_{i=1}^n \rho(\sqrt{(x_i - T)' C^{-1} (x_i - T)}) \leq b, \quad (12)$$

pentru o funcție ρ care de multe ori se alege a fi funcția lui Tukey:

$$\rho(x) = \begin{cases} \frac{x^2}{2} - \frac{x^4}{2c^2} + \frac{x^6}{6c^4}, & \text{daca } |x| \geq c \\ \frac{c^2}{6}, & \text{daca } |x| < c \end{cases} \quad (13)$$

iar punctul de rupere este dat de valoarea lui c , $BP=6b/c^2$.

Dacă considerăm $\rho(x) = x^2$ și $b=1$ obținem cazul metodei celor mi mici pătrate.

S-estimatorii au un punct de rupere înalt și proprietatea de normalitate asimptotică (Rousseeuw și Zohai 1984) dar metodele de calcul sunt intensive, iar când datele sunt foarte contaminate nu sunt unici (Woodruff și Rocke, 1994).

1.4.5 Detectarea datelor anormale (outliers)

Considerăm un model standard localizare-covarianță, adică un set de date d -dimensionale x_1, x_2, \dots, x_n care sunt realizări ai unor variabile aleatoare X_1, X_2, \dots, X_n cu o distribuție de probabilitate eliptică $P_{\mu, \Sigma}$ cu densitatea de probabilitate:

$$f(x) = |\Sigma|^{-1/2} g(\|x - \mu\|_{\Sigma^{-1}}^2) (*)$$

cu $\mu \in R^d, \Sigma \in PDS(d)$, clasa matricilor de dimensiune $d \times d$ simetrice și pozitiv definite și $g : [0, \infty) \rightarrow [0, \infty)$.

O dată ce am obținut estimări robuste ale localizării și împrăștierii putem aborda problema căutării datelor care au o deviație anormală față de media estimată și conform cu matricea de covarianță.

Există trei tipuri de contaminare care cer diferite tehnici de detecție.

Prima categorie sunt acele date anormale pentru care distanțele robuste RD sunt mai mari decât un prag fix (Rousseeuw și van Zomeren 1990, Rocke și Woodruff 1996) sau decât un prag adaptativ care depinde nu numai de structura datelor ci și de mărimea eșantionului (Gervini 2003).

Dacă f este multivariat normală atunci $RD(X_i) > \sqrt{\chi_{d,0.975}^2}$ determină aceste anormalități cu prag de semnificație 0.975.

Al doilea tip este cel al datelor anormale (shift outliers) care corespund unor clusteri care urmează aceiași structură ca a majorității datelor, dar sunt deplasați față de media

majorității populației. De exemplu pentru o populație distribuită $N(\mu, \Sigma)$ un nor de date anormale deplasate pot fi distribuite $N(\mu + s, \Sigma)$ și astfel corup metrica bazată pe distanța Mahalanobis (Rocke și Woodruff 1999) făcând ca diferența dintre datele normale și acestea să dispară dacă se consideră această distanță.

Al treilea tip de anormalitate este cel al datelor anormale localizate (point mass outliers) dat de cel rezultat dintr-o concentrare a contaminării într-o regiune redusă (Pena și Prieto 2001)

2. Un exemplu de estimare a datelor anormale folosind estimatorii MCD

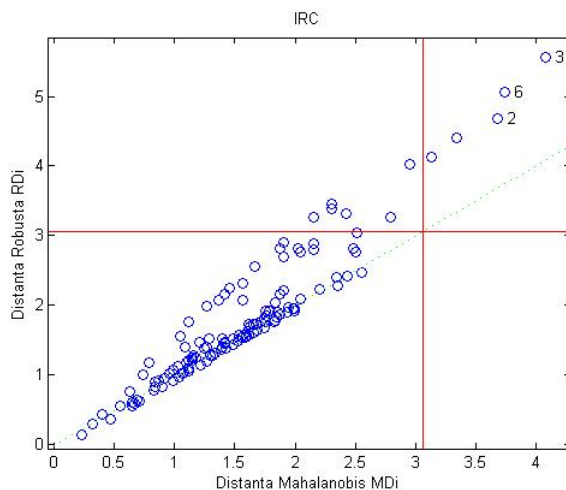
Considerăm în spațiul R^3 , variabilele Inflație, Rata dobânzii active a clienților nebancari, Cursul Euro-Ron în perioada februarie 2000 - iunie 2010 (n=125 de puncte) care urmează un model multivariat normal și calculăm distanțele Mahalanobis MD_i și distanțele robuste RD_i cu pragul de semnificație 97,5% și $h=0.9n$ ceea ce asigură un punct de rupere BP=10%.

Din graficul următor (figura 1) se observă efectul de mascare al distanțelor nerobuste MD_i , care în comparație cu RD_i sunt mai mici.

Punctele desemnate ca outliers sunt cu indexii în intervalul 1-12 și 26. Iată că deviația față de modelul normal se întâmplă în primul interval (an), apoi se poate considera că datele urmează acest model (graficul 3 în care se raportează RD_i la rădăcină pătrată a cuantilelor distribuției chi-pătrat). De fiecare dată, pentru orice h în intervalul $(n/2, n)$, datele anormale se găsesc în prima perioadă. Pentru $h=n$ obținem puncte anormale 1-4, 6, 12, 26, care apar ca anormale în toate estimările așa că putem conchide că aceste puncte sunt outliers.

În figura 1, se poate observa efectul de mascare asupra distanțelor Mahalanobis, care sunt subevaluate datorită măsurilor nerobuste ale mediei și covarianței.

Figura 1



În figura 2 se pot observa datele anormale comparându-se distanțele robuste cu valoarea chi-pătrat, iar în figura 3 se raportează cuantile distribuției empirice a distanțelor robuste cu cuantile distribuției chi-pătrat cu d grade de libertate și prag de semnificație 0.975.

Figura 2

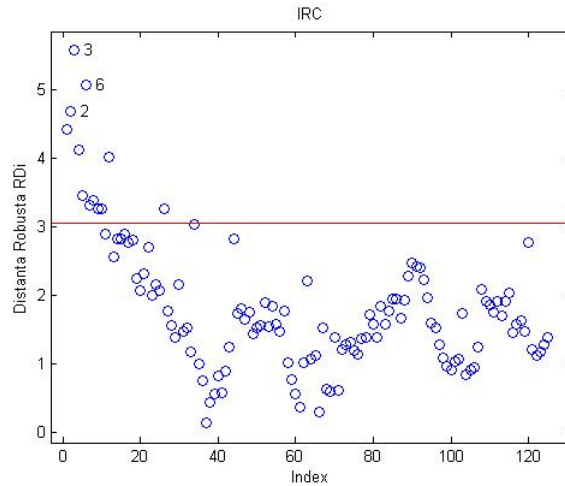
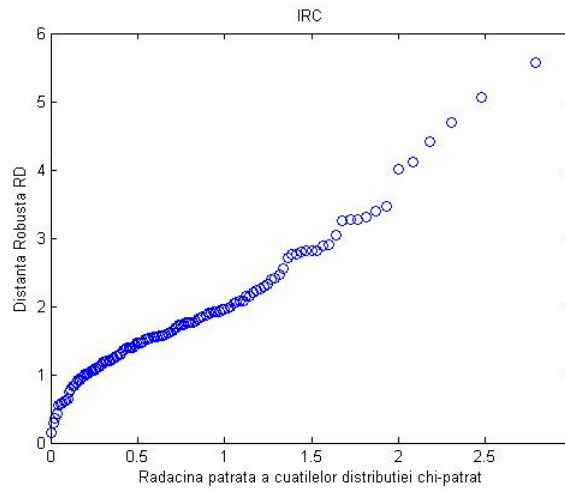


Figura 3



3. Regresie ortogonală robustă

Fekri și Ruiz-Gazen (2004) au definit regresie ortogonală ponderată robustă.

Fie $X \in R^{d \times n}$ matricea conținând pe coloane cele n date observate $X_i \in R^d$, $i=1, \dots, n$ și fie x_i aparținând direcției afine $D_{A,b} = \{x \in R^d, A^1 x = b; A \in R^d, b \in R\}$, $X_i = x_i + e_i$ și e_i distribuit normal cu medie 0 și matrice de covarianță diagonală Σ .

Modelul regresiei ortogonale estimează A și b minimizând criteriul:

$$\sum_{i=1}^n \|X_i - x_i\|^2 = \sum_{i=1}^n \|X_i - \Pi_{D_{A,b}}(X_i)\|^2 + \sum_{i=1}^n \|x_i - \Pi_{D_{A,b}}(X_i)\|^2$$

Minimizarea acestui criteriu duce la $x_i = \Pi_{D_{A,b}}(X_i)$ și minimizarea :

$$\sum_{i=1}^n \|X_i - \Pi_{D_{A,b}}(X_i)\|^2 = \sum_{i=1}^n \|A'X_i - b\|^2 = \|(A'X - bu)\|_F^2$$

cu $u \in R^n = [1,1,\dots,1]$ și $\|M\|_F^2 = Tr(M'M) = Tr(MM')$ reprezentând norma Frobenius.

Regresia ortogonală robustă schimbă criteriul de minimizare într-unul ponderat:

$$\sum_{i=1}^n w(\|X_i - \mu\|_{\Sigma^{-1}}^2) \|X_i - x_i\|^2$$

$$\text{care conduce la } \sum_{i=1}^n w(\|X_i - \mu\|_{\Sigma^{-1}}^2) \|A'X_i - b\|^2 = \|(A'X - bu')W\|_F^2,$$

unde $W = [w_1, \dots, w_n]$, $w_i = w(\|X_i - \mu\|_{\Sigma^{-1}}^2)$.

Soluția optimă (A*,b*) este dată de:

$$A^* = \text{vectorul propriu de normă unitate asociat matricii pozitiv semidefinite} \\ H = X(W^2 - \frac{1}{v}W^2uu'W^2)X' \quad \text{și} \quad b^* = \frac{1}{v}A^*XW^2u.$$

Fie :

$$\mu_{R,n} = \frac{\sum_{i=1}^n w(\|X_i - \mu_n\|_{\Sigma_n^{-1}}^2) X_i}{\sum_{i=1}^n w(\|X_i - \mu_n\|_{\Sigma_n^{-1}}^2)} \quad \text{și} \quad \Sigma_{R,n} = \frac{\sum_{i=1}^n w(\|X_i - \mu_n\|_{\Sigma_n^{-1}}^2) (X_i - \mu_{R,n})(X_i - \mu_{R,n})'}{\sum_{i=1}^n w(\|X_i - \mu_n\|_{\Sigma_n^{-1}}^2)},$$

unde μ_n și Σ_n sunt estimatori robusți ai localizării și covarianței.

Minimul este atins de direcția afină care conține $\mu_{R,n}$ și este paralelă direcției dată de vectorul propriu asociat celei mai mici valori proprii a matricii $\Sigma_{R,n}$.

Funcția w este o funcție descrescătoare care penalizează observațiile cu distanța $\|X_i - \mu_n\|_{\Sigma_n^{-1}}^2$ mare. Vom folosi în exemple funcția pondere $1I_{[0,c]}$ unde c este cuantila 97,5% a distribuției chi-pătrat.

Tabelul 1

	Coeficienții hiperplanului		
	TLS	RS	RMCD
Inflație	-0.4239	-0.4875	-0.4870
Rata dobânzii	0.9057	0.8732	0.8734

Curs	0.0041	-0.0028	0.0027
Termen liber	-0.0249	-0.0195	-0.0191

4. Rezultate de simulare

În această secțiune am făcut simulări pentru a compara estimatorii coeficienților hiperplanului obținuți prin metoda robustă și cea standard. Parametrii populației sunt cei obținuți cu metoda RMCD pentru exemplul prezentat.

Am considerat datele necontaminate care sunt conforme cu hiperplanul generat de ecuația ortogonală robustă obținută folosind covarianța RMCD și urmează modelul regresiei ortogonale cu erori de tip multivariat normal. Contaminarea s-a realizat considerând că erorile urmează alta distribuție de tip eliptic, sau combinații de distribuții normale.

Astfel, rezultatele obținute arată ce s-ar întâmpla dacă ar exista date anormale modelului considerat și am aplica una dintre metodele robuste sau nerobuste. Aceste simulări constituie o măsurare a robusteții globale a metodelor.

Simulările constau în 1000 de eșantioane cu dimensiunea 125 pentru mai multe distribuții:

1. distribuția normală $N(\mu, \Sigma)$ (NOR)
2. distribuția normală simetric contaminată (SCN), care este un mixaj de 90% $N(\mu, \Sigma)$ și 10% $N(\mu, 10\Sigma)$
3. distribuția multivariată Cauchy (CAU)
4. distribuția normală asimetric contaminată (ACN), care este o mixtură de 90% $N(\mu, \Sigma)$ și 10% $N(\mu_1, 10\Sigma)$ cu $\mu_1 = 2 * \mu$

Tabelul 2

	Media empirică			MSE(media pătratelor erorilor)		
	TLS	RS	RMCD	TLS	RS	RMCD
NOR	-0.4870	-0.4870	-0.4870	4.36e-13	0.25e-13	4.36e-13
	0.8734	0.8734	0.8734	1.41e-11	0.54e-11	1.45e-11
	0.0027	0.0027	0.0027	1.37e-16	0.212e-14	2.75e-14
	-0.0190	-0.0191	-0.0191	0.14e-15	0.3e-15	0.1e-15
SCN	-0.4874	-0.4870	-0.4870	2.58e-10	0.31e-11	2.6e-13
	0.8744	0.8734	0.8734	2.35e-12	2.2e-12	2.3e-12
	0.0026	0.0027	0.0027	2.8e-14	3.2e-17	2.7e-17
	-0.0191	-0.0191	-0.0191	1.4e-3	2.3e-4	3.7e-4
CAU	0.7135	-0.4934	-0.4865	0.0514	8.9e-4	2.26e-4

	0.7004	0.8761	0.8735	0.0299	8.4e-4	7.02e-5
	0.0133	0.0023	0.0027	1.132e-4	4.3e-6	8.31e-8
	0.0647	-0.0798	-0.0989	0.0021	0.8260	1.8918
ACN	0.4224	-0.4794	-0.4824	0.0265	5.6e-6	0.4e-6
	-0.49118	0.8700	0.8703	0.0129	3.4e-4	1.69e-4
	6.7188e-04	0.0025	0.0026	2.9e-6	4.6e-7	2.07e-7
	0.0019	-0.0189	-0.0204	0.0004	0.1e-5	0.0003

În cazul contaminării simetrice rezultatele pentru regresia TLS sunt foarte bune deoarece datele (inflația, rata dobânzii, cursul de schimb) sunt puternic corelate așa modificarea considerată matricei de covarianță înseamnă o corelare mai puternică a datelor.

În ce privește modelele de regresie ortogonală robuste, ele par echivalente și deci poate fi folosită oricare dintre ele și cum algoritmul pentru metoda RMCD este mai eficient este de preferat utilizarea acestei metode.

Ca măsură globală performanței estimatorilor ecuației regresiei ortogonale am considerat valoarea absolută a cosinusului unghiului dintre parametrii estimați și cei ai populației. Această măsură poate fi interpretată ca mărime a influenței datelor anormale (outliers) asupra estimatorilor.

Figura 4 și 5 prezintă funcția de distribuție cumulată empirică pentru cele 1000 de eșantioane ale simulărilor, pentru distribuțiile Cauchy și distribuția normală asimetric contaminată. Distribuția cumulativă empirică ar trebui să fie masată în apropierea lui 1 dacă parametrii estimați au valori aproape de valorile reale deoarece atunci cosinusul unghiului dintre parametrii ar fi aproape de 1.

În cazul distribuției normale toate cele trei curbe au vârful de maxim către 1. Însă nu aceeași este situația pentru estimarea TLS (regresia ortogonală nerobustă) în timp ce metodele robuste dau rezultate bune.

Figura 4

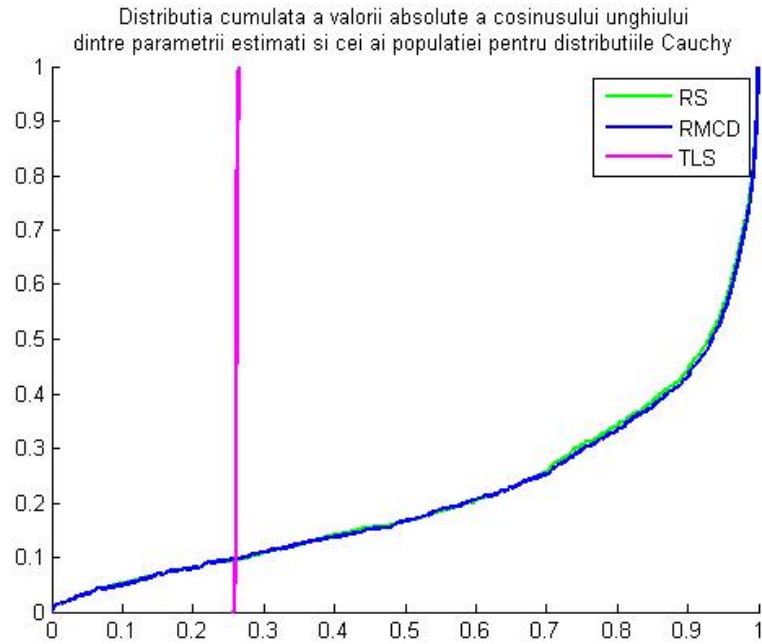
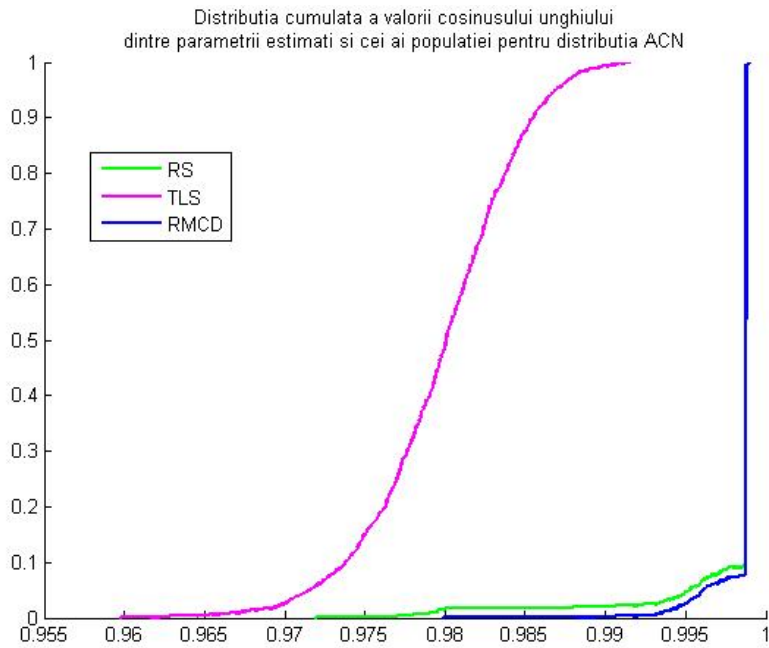


Figura 5



Concluzii

În acest studiu am analizat influența datelor anormale asupra estimatorilor în cazul modelului Eroare-in-Variabile cu distribuția erorilor multivariat eliptice pentru regresia ortogonală.

Estimatorii robuști sunt definiți ca estimatori ponderați de tip cele-mai-mici-pătrate fiind derivați din estimatori robuști ai localizării (mediei) și împrăștierii (covarianței) datelor. În acest studiu am folosit estimatori S și RMCD, dar pot fi utilizați orice alți estimatori robuști ai primelor două momente ale distribuției empirice a datelor cum sunt estimatorii de tip proiecție Stahel-Donoho sau M-estimatori.

Am văzut că estimatorii clasici ai hiperplanului care se află la distanță minimă de date (distanța Euclidiană deoarece este analizat cazul regresiei ortogonale) se îndepărtează de valorile reale dacă datele sunt contaminate de erori care urmează alte distribuții decât cea gaussiană. Acest caz este des întâlnit în situațiile reale când datele deviază față de distribuția normală prezentând de multe ori distribuții asimetrice sau cu valori extreme.

Folosirea unor metode nerobuste în aceste cazuri duce la estimări incorecte și implicit la inferențe greșite.

Rezultatele obținute din simulări și prezentate în tabel sunt similare celor din Fekri și Ruiz-Gazen (2004) dovedind că pentru modelul necontaminat cei mai eficienți estimatori ai hiperplanului ce reprezintă regresia liniară ortogonală fără ca celelalte metode robuste să fie ineficiente, iar pentru modelele contaminate regresia robustă este eficientă dar regresia ortogonală clasică este complet inadecvată în cazurile distribuției Cauchy și celei normal asimetric contaminate.

Programele utilizate pentru calculul estimatorilor și pentru simulări au fost scrise în limbaj MathLab. Estimatorii RMCD au fost calculați folosind rutinele LIBRA.

Referințe

- Bickel, P.J. (1964), On some alternative estimates for shift in the p-variate one sample problem, *Ann. Math. Statist.*, 35: 1079-1090.
- Brown, M. (1982), Robust line estimation with errors in both variables, *Journal of the American Statistical Association*, 77: 71-79.
- Croux, C. and Haesbroeck, G. (1999), Influence and Efficiency of the Minimum Covariance Determinant Scatter Matrix Estimator, *Journal of Multivariate Analysis*, 71: 161-190.
- Davies. P.L. (1987), Asymptotic behaviour of S-estimators of multivariate location parameters and dispersion matrices, *Annals of Statistics*, 15: 1269-1292.

- Fekri, M. and Ruiz-Gazen, A. (2004), Robust weighted orthogonal regression in the errors-in-variables model, *Journal of Multivariate Analysis*, 88: 89–108.
- Gervini, D. (2002), A robust and efficient adaptive reweighted estimator of multivariate location and scatter, *Journal of Multivariate Analysis*, 84:116-144.
- Hampel, F.R. (1974), The influence curve and its role in robust estimation, *Journal of the American Statistical Association*, 69: 383-393.
- Hampel, F.R. (1971), A general qualitative definition of robustness, *Ann. Math. Statist.*, 42:1887-1896.
- Huber, P.J. (1972), Robust statistics: A review, *Ann. Math. Statist.*, 35:73-101.
- Hubert, M. Rousseeuw, P. J. And Van Aelst, S. (2008), High-Breakdown Robust Multivariate Methods, *Statistical Science*, 23(1): 92-119.
- Lopuhaä, H.P. (1989), On the relation between S-estimators and M-estimators of multivariate location and covariance, *Ann. Statist.*, 17(4):1662-1683.
- Lopuhaä, H.P. (1999), Asymptotics of reweighted estimators of multivariate location and scatter, *Ann. Statist.*, 27: 1638–1665.
- Lopuhaä, H.P. and Rousseeuw, P.J. (1991), Breakdown points of affine equivariant estimators of multivariate location and covariance matrices, *Ann. Statist.*, 19: 229–248.
- Maronna, R.A. (1976), Robust M-estimates of multivariate location and scatter, *Ann. Statist.*, 4: 51–67.
- R.A. Maronna and S. Morgenthaler, (1986), Robust regression through robust covariances, *Comm. Statist.—Theory Methods*, 15: 1347–1365.
- Maronna, R.A. Stahel, W.A. and Zohai, V.J. (1992), Bias-robust estimates of multivariate scatter based projections, *Journal of Multivariate Analysis*, 42:141-161.
- Rocke, D.M. and Woodruff, D.L. (1999),
- Rousseeuw, P. J. and Van Zomeren, B.C. (1990), Unmasking multivariate outliers and leverage points, *J. Amer. Statist. Assoc.*, 85: 633–651.
- Rousseeuw, P. J. And Van Driessen, K. (1999), A Fast Algorithm for the Minimum Covariance Determinant Estimator, *Technometrics*, 41: 212-223.
- Rubinstein, R.Z, and Kroese, D.P. (2008), Simulation and the Monte Carlo Method, John Wiley / Sons, Hoboken, New Jersey.
- Woodruff, D.L. and Rocke, D.M. (1994), Computable robust estimation of multivariate location and shape in high dimension using compound estimators, *J. Amer. Statist. Assoc.*, 89: 888-896.

Anexa.

In această anexă se prezintă programul de simulare pentru datele contaminate cu distribuția Cauchy.

```
function res=Multivar_tLSSim(IRCplan,Mer,Ser,trial,p,multisigma,multm)
% Simulare Monte Carlo pentru setul de date IRCdata contaminat cu distributie Cauchy
n = size(IRCplan(:,1)); % The number of function evaluations
% --- Generate vectors of random inputs
% r ~ distributie Normala N(medie=m,sd=SIGMA)
% rc ~ distributie Normala N(medie=multm*m,sd=10*SIGMA)
% re ~ distributie Cauchy CAU(sd=Ser)
% rce ~ distributie Normala N(medie=Mer,sd=Ser)

res.TLS=[];
res.TLSm=[];
res.TLSmse=[];
res.RMCDmse=[];
res.sz=[];
res.multisigma=multisigma;
res.multm=multm;
for i=1:4
res.TLSmse(i)=0;
res.RMCDmse(i)=0;
end
j=0;
SIGMA=cov(IRCplan);
m=mean(IRCplan);
while j<trial
r=mvnrnd(m,SIGMA,n(1)-floor(n(1)*multisigma));
re=mvnrnd(multm*m,10*SIGMA,floor(n(1)*multisigma));
IRCrand=[r; rc];
re=mvtrnd(Ser,1,25); rce=mvnrnd(Mer,Ser,100);
IRCrand=IRCrand+[re;rce];
m=mean(IRCplan);
COVrand=cov(IRCrand); mrand=mean(IRCrand);
j=j+1;
[v,d]=eig(COVrand);
k=1;
for i=2:3
if d(i,i) < d(k,k)
k=i;
end
end
res.TLS(4,j)=0;
for i=1:3
res.TLS(i,j)=v(i,k);
res.TLSmse(i)=res.TLSmse(i)+(abs(v(i,k))-abs(p(i)))*(abs(v(i,k))-abs(p(i)));
res.TLS(4,j)=res.TLS(4,j)+v(i,k)*mrand(i);
end
```

```

res.TLSmse(4)=res.TLSmse(4)+(res.TLS(4,j)-abs(p(4)))*(res.TLS(4,j)-abs(p(4)));
[rew,raw]=mcdcov(IRCrnd,'plots',0, 'alpha',0.75) ; plan=size(rew.plane);
res.sz(j)=plan(1);
[v,d]=eig(rew.cov);
k=1;
for i=2:3
    if d(i,i) < d(k,k)
        k=i;
    end
end
res.RMCD(4,j)=0;
for i=1:3
    res.RMCD(i,j)=v(i,k);
    res.RMCDmse(i)=res.RMCDmse(i)+(abs(v(i,k))-abs(p(i)))*(abs(v(i,k))-abs(p(i)));
    res.RMCD(4,j)=res.RMCD(4,j)+v(i,k)*mrand(i);
end
res.RMCDmse(4)=res.RMCDmse(4)+(res.RMCD(4,j)-abs(p(4)))*(res.RMCD(4,j)-abs(p(4)));
end
res.TLSm(1)=mean(res.TLS(1,:)); res.TLSm(2)=mean(res.TLS(2,:)); res.TLSm(3)=mean(res.TLS(3,:));
res.TLSm(4)=mean(res.TLS(4,:)); res.RMCDm(1)=mean(abs(res.RMCD(1,:)));
res.RMCDm(2)=mean(abs(res.RMCD(2,:))); res.RMCDm(3)=mean(abs(res.RMCD(3,:)));
res.RMCDm(4)=mean(abs(res.RMCD(4,:)));
end

```