

Tehnici neparametrice

Bianca Păuna*

Abstract: *Această lucrare prezintă doi estimatori mai puțin cunoscuți pentru vizualizarea funcțiilor de densitate și a relațiilor dintre variabile. Acești estimatori sunt neparametrici în sensul că nu sunt necesare ipoteze a priori privind forma funcțională a dependențelor. În prezentare s-a pus accent pe doi estimatori, estimatorul Kernel al densității și estimatorul Nadaraya – Watson al relațiilor dintre variabile. Aceste tehnici au fost folosite în continuare pentru a evidenția caracteristicile consumului în România.*

Keywords: estimatori neparametrici, estimatorul Kernel, estimatorul Nadaraya-Watson.

Clasificare JEL: C14

1 Introducere

Graficele sunt instrumente potrivite atunci când se dorește vizualizarea anumitor caracteristici ale datelor. Dar, datorită numărului de dimensiuni care pot fi vizualizate odată, folosirea graficelor este restricționată. Grafice/Proiecții cu mai mult de trei dimensiuni sunt dificil de reprezentat și de înțeles. Deși graficele nu pot surprinde modele complicate ele sunt un punct de plecare pentru orice tip de investigație. Tablele sunt un alt mod de a prezenta informațiile, dar devin foarte mari și greu de urmărit atunci când se dorește prezentarea completă a variabilelor și a relațiilor dintre ele.

Estimarea densităților este în general punctul de plecare al oricărei analize. Există mai multe metode care pot fi folosite pentru obținerea unui estimator al funcțiilor de densitate (histograme, estimatorul neparametric al densității, etc.). O discuție a avantajelor și dezavantajelor privind cele două metode va fi prezentată în următoarea secțiune.

În afara densităților, relațiile dintre variabile sunt un alt punct de interes. Reprezentarea punctelor pe un grafic poate fi folositoare pentru descoperirea relațiilor dintre variabile numai în cazul în care observațiile nu sunt foarte disperse sau numeroase. Odată cu creșterea numărului de observații, sau a dispersiei punctelor, ochiului uman îi vine din ce în ce mai greu să identifice relațiile dintre variabile, și de aceea este nevoie de un alt instrument. Alegerea este în acest caz între metodele parametrice și cele neparametrice.

* Drd., cercetător la Institutul Național de Cercetare Economică, Academia Română.

Proprietățile estimatorului parametric OLS, nedepășirea, consistența și eficiența îl fac foarte popular în analizarea datelor. Prețul care trebuie plătit este în materie de ipoteze legate de distribuții și specificarea modelului, dar aceasta este o problemă numai în cazul în care forma funcțională nu este cunoscută. Alternativ pot fi folosite metode neparametrice. Estimatorii sunt polarizați, dar sunt consistenți, și cel mai important, nu este nevoie să se facă nici un fel de ipoteze *a priori* privind forma dependenței, deci sunt un instrument foarte adecvat în analizele datelor.

Estimarea parametrică este instrumentul adecvat în cazurile în care se știe că ipotezele pe care se bazează construcția modelului sunt satisfăcute. În cazul în care modelul nu este bine specificat, estimatorii coeficienților vor fi deplasați și inconsistenți, în timp ce estimatorii neparametrici vor fi consistenți și asimptotic nedepășați.

Un exemplu este cazul estimării dependenței dintre variabila dependentă (Y) și cea explicativă (X). Estimatorul OLS presupune existența unei dependențe liniare și constă într-o linie trasată printre puncte, astfel încât suma pătratului erorilor să fie minimă. Estimând coeficienții prin metoda OLS se obțin cei mai buni estimatori al unei dependențe liniare, sau liniarizate.

În cazul în care nu se cunoaște cu exactitate forma relației dintre variabile, este necesară o alternativă mai puțin rigidă, cum este cazul metodelor neparametrice. Metodele neparametrice au apărut datorită nevoii de a evita aplicarea unei dependențe funcționale rigide modelului, deci termenul de neparametric se referă în acest caz la forma funcțională flexibilă. Exemple tipice de folosire a tehnicilor neparametrice este în construirea modelului, verificarea, inferențe și predicție.

Echivalentul neparametric al OLS-ului este estimatorul neliniar, și ideea este ca în loc de a se impune forme funcționale rigide, datele sunt lăsate să indice forma funcțională a dependenței dintre cele două variabile. Scopul tehnicilor neparametrice este de a contribui la reprezentarea relației adevărate dintre două variabile. Metoda cea mai la îndemână este unirea punctelor, dar graficul care ar rezulta ar fi mult prea eratic pentru a fi de folos, de aceea orice tehnică neparametrică implică și un proces de netezire. Procesul de netezire implică un compromis, graficul trebuie să fie suficient de neted pentru a putea descrie relația dintre cele două variabile, fără a renunța la toată variația locală. Alegerea nivelului de netezime a graficului se poate face subiectiv prin comparație între mai multe grafice cu diverse grade de netezire, sau cu ajutorul metodelor statistice.

Ideea care stă la baza estimatorului neliniar este de a calcula variabila dependentă în fiecare punct folosind informația disponibilă în vecinătatea punctului. Astfel, variabila dependentă este obținută ca o medie a valorilor funcției în vecinătate. Această procedură va produce estimatori polarizați, și ca exemplu se consideră cazul estimării valorii funcției la maxim sau minim. Toate punctele din jurul maximumului (minimumului) au valori mai mici (mari) deci valoarea rezultată va fi mai mică (mare) decât valoarea reală. Este posibil să se reducă valoarea polarizării prin reducerea dimensiunii intervalului, dar odată cu scăderea intervalului, erorile din date vor avea o influență mai mare asupra estimatorului, deci variația acestuia va crește. Deci, procedura de alegere a intervalului optim este și ea rezultatul unui compromis între valoarea polarizării și a variației estimatorului.

2 Estimatorul densității

Cea mai simplă analiză neparametrică este estimarea unei funcții de densitate a unei variabile X_i . Aceasta constă în găsirea probabilităților asociate fiecărei valori din domeniul de definiție a variabilei X_i . Pentru obținerea estimatorului trebuie în primul rând împărțit domeniul funcției în intervale egale, după care se numără datele din fiecare interval.

Pentru un interval de dimensiune h estimatorul histogramei este de forma următoare:

$$\hat{f}(x) = \frac{1}{nh} (\text{numărul de puncte din interval})$$

Mai formal, relația de calcul a histogramei poate fi scrisă astfel:

$$\hat{f}_h(x) = \frac{1}{nh} \sum_{i=1}^n \sum_j I(X_i \in B_j) I(x \in B_j)$$

Factorii care influențează forma histogramei sunt atât punctul de origine a intervalului, cât și dimensiunea intervalului. Controlul netezimii curbei de histogramă este dat de dimensiunea intervalului. Această dependență a formei histogramei de cei doi factori este una din dezavantajele folosirii lor. Pentru corectarea acestei dependențe se poate folosi metoda medierii histogramelor cu origine diferită. Aceasta constă în estimarea mai multor histograme (care au originea intervalului diferită) pentru dimensiunea optimă a intervalului, iar histograma estimată este obținută prin medierea tuturor histogramelor.

Histogramele prezintă unele dezavantaje care nu o fac foarte răspândită în aplicații. Prin definiție, estimatorul este discontinuu: are valori constante în fiecare interval, dar are salturi la capetele intervalului, deci este un instrument nepotrivit mai ales în cazurile în care se dorește obținerea derivatelor de ordinul întâi.

Un alt estimator al densităților este estimatorul naiv. Acesta se obține pornind de la definiția funcției de densitate: probabilitatea ca $x \in [x-\varepsilon, x+\varepsilon]$ atunci când $\varepsilon \rightarrow 0$. Pentru estimarea probabilității se folosește proporția de puncte care se regăsesc în intervalul $[x-\varepsilon, x+\varepsilon]$.

$$\hat{f}(x) = \frac{1}{nh} (\text{număr de } X_i \text{ din } [x-h/2, x+h/2])$$

Atât histograma cât și estimatorul naiv au definiții similare, dar diferă din punct de vedere al calculelor făcute. Histograma are valori constante pe fiecare interval, și salturi la capetele acestuia. Parametrul h controlează numărul de intervale. Estimatorul naiv este calculat pentru fiecare valoare a lui x . Gradul de netezire al graficului este controlat prin dimensiunea lui h . Din punct de vedere computațional, histograma necesită mai puține operații la estimare, în timp ce estimatorul naiv calculează pentru fiecare observație o valoare a densității, fiind mai intensiv din punct de vedere computațional.

Estimatorul densității poate fi scris mai compact dacă se definește o funcție de pondere: $w(u) = \frac{1}{2} \cdot I(|u| \leq 1)$ unde $I(\bullet)$ este funcția indicator care ia valoarea 1 atunci când este expresia este adevărată și zero când este falsă:

$$\hat{f}_h(x) = \frac{1}{n} \cdot \sum_{i=1}^n \frac{1}{h} w\left(\frac{x - X_i}{h}\right)$$

Funcția de pondere dă ponderi egale fiecărei observații din interval, deci estimatorul densității este o funcție discontinuă care are salturi în punctele $X_i \pm h$ și are derivata egală cu zero în rest.

2.1 Estimatorul Kernel al densității

Discontinuitățile estimatorului naiv pot fi corectate cu ușurință prin alegerea unei forme funcționale diferite pentru ponderi. O funcție care să corecteze discontinuitățile trebuie să satisfacă două proprietăți:

$K(u) \rightarrow 0$ as $|u| \rightarrow 1$ and

$$\int_{-\infty}^{\infty} K(x) dx = 1$$

Prima condiție asigură continuitatea funcției de densitate estimată, prin forțarea ponderilor spre zero la sfârșitul intervalului. Aceasta garantează continuitatea funcției de densitate pe întreg domeniul, o proprietate importantă atunci când se dorește estimarea derivatei de ordin întâi.

Cea de-a doua condiție garantează că suma ponderilor este 1, și în consecință satisfacerea condiției pe care toate funcțiile de densitate trebuie să o satisfacă și anume integrala funcției pe întreg domeniul de definiție să fie 1. Cum cea de a doua condiție este similară cu definiția funcției de densitate, de multe ori pentru funcția de ponderi este aleasă o densitate cunoscută.

Ordinul funcției kernel este definit ca primul moment diferit de zero. Deci un kernel este de ordinul p dacă $k_p = \int u^p K(u) du \neq 0$. Un kernel pozitiv poate fi maximum de ordinul 2. Deoarece în estimarea densităților funcția este restricționată la valori pozitive, kernelul de ordinul 2 este cel mai popular în aplicații. Kernelii de ordin mai mare sunt folosiți și ei în anumite cazuri, care vor fi prezentate în această secțiune.

Ținând cont de toate observațiile, un estimator kernel al funcției de densitate este de forma următoare:

$$\hat{f}_h(x) = \frac{1}{nh} \cdot \sum_{i=1}^n K\left(\frac{x - X_i}{h}\right)$$

unde h este dimensiunea intervalului.

O alegere posibilă pentru funcția de ponderi este funcția de densitate normală, care fiind o funcție de densitate, satisface condițiile enumerate mai sus. Totuși distribuția normală nu are suport finit, ceea ce înseamnă că pentru fiecare la calculul densității în fiecare punct se folosește informația de la toate punctele pentru care există din valori ale funcției. Acest lucru crește costurile din punct de vedere computațional, și din acest motiv funcțiile kernel ale căror valori nu sunt într-un interval finit, nu sunt foarte răspândite în practică.

Una din cele mai răspândite funcții de ponderi este kernelul Epanechnikov care ia forma următoare:

$$K(u) = \frac{3}{4} \cdot (1 - u^2) \cdot I(|u| \leq 1)$$

Folosirea funcțiilor de pondere care iau valori pe un domeniu infinit (cum este cazul densității normale) rezultă în creșterea semnificativă numărului de operații necesare la estimare. La fiecare valoare, trebuie calculate ponderi pentru toate observațiile, deci un total de n^2 operații. Când se folosesc funcții de pondere finite numărul de operații necesare pentru estimare scade semnificativ la nh . Chiar în cazul folosirii ponderilor finite, atunci când numărul de date este foarte mare, aspectele computaționale pot fi prohibitive.

2.2 Proprietățile statistice ale estimatorului kernel al densității

Proprietățile statistice sunt criteriul cel mai răspândit pentru clasificarea diversilor estimatori. Estimatorul neparametric este funcție de date, dar forma estimatului este sensibilă la alegerea kernelului și a dimensiunii intervalului. Din acest motiv, proprietățile statistice (deplasarea și variația) sunt folosite ca instrument de decizie pentru alegerea ponderilor și a gradului de netezire.

Proprietățile estimatorului într-un punct sunt judecate prin folosirea MSE^\dagger (media pătratelor erorilor) care este definită în felul următor:

$$MSE_x(\hat{f}) = E\{\hat{f}(x) - f(x)\}^2$$

Dacă se scrie MSE în funcție de deplasare și de variație se obține expresia:

$$MSE_x(\hat{f}) = \{E\hat{f}(x) - f(x)\}^2 + var\hat{f}(x)$$

Integrând MSE pe întregul domeniu de definiție se obține MISE (media integrată a pătratelor erorilor). MISE este o măsură a acurateții globale a estimatorului și are următoarea definiție:

$$MISE(\hat{f}) = E \int \{\hat{f}(x) - f(x)\}^2 dx$$

Prin aplicarea proprietăților integralelor și cele ale operatorului de așteptare, MISE poate fi rescrisă astfel:

[†] MSE – abrevierea lui Mean Squared Error.

$$\begin{aligned} \text{MISE}(\hat{f}) &= \int E\{\hat{f}(x) - f(x)\}^2 dx = \int \text{MSE}(\hat{f}) dx \\ &= \int \{E\hat{f}(x) - f(x)\}^2 dx + \int \text{var } \hat{f}(x) dx \end{aligned}$$

Dimensiunea optimă a intervalului de estimare se obține ca o soluție a maximizării MISE. Pentru aceasta, MISE trebuie scrisă ca o funcție de dimensiunea intervalului. Modul de derivare a dimensiunii intervalului va fi prezentat în secțiunea următoare.

Estimatorul densității kernel are forma următoare[‡]:

$$\hat{f}(x) = n^{-1} \sum K_h((x-y)/h) \text{ where } K_h = n^{-1} K((x-y)/h)$$

$$E(\hat{f}) = n^{-1} E \{ \sum K_h((x-y)/h) \} = n^{-1} \sum E K_h((x-y)/h) = E K_h((x-y)/h)$$

$$\text{var } \hat{f}(x) = \text{var} \{ n^{-1} \sum K_h((x-y)/h) \} = n^{-2} \text{var } \sum K_h((x-y)/h) = n^{-1} \text{var } K_h((x-y)/h).$$

Pentru calculul MISE este necesar să se obțină expresii ale deplasării și ale variației.

$$\text{deplasarea}(x) = E \hat{f}(x) - f(x) = \int h^{-1} K((x-y)/h) f(y) dy - f(x)$$

Se face următoare schimbare de variabile: $y = x - ht$; $dy = h dt$ și pentru că $\int K(u) du = 1$, deplasarea se poate rescrie:

$$\text{deplasarea}(x) = \int K_t(t) \{f(x - ht) - f(x)\} dt$$

Prin folosirea descompunerii Taylor:

$$f(x - ht) = f(x) - ht f'(x) + \frac{1}{2} (ht)^2 f''(x) + \dots$$

și după simplificări deplasarea poate fi rescrisă astfel:

$$\text{deplasarea}(x) = \int K_t(t) \frac{1}{2} (ht)^2 f''(x) dt + O(h^3) \cong \frac{1}{2} h^2 f''(x) \int t^2 K_t(t) dt = \frac{1}{2} h^2 f''(x) k_2$$

Prin definiție variația poate fi scrisă astfel:

$$\begin{aligned} \text{var } \hat{f}(x) &= n^{-1} \text{var } K_h((x-y)/h) = E \{ K_h((x-y)/h) - E K_h((x-y)/h) \}^2 = \\ &= n^{-1} E \{ K_h^2((x-y)/h) + [E K_h((x-y)/h)]^2 - 2 K_h((x-y)/h) E K_h((x-y)/h) \} = \\ &= n^{-1} \{ E K_h^2((x-y)/h) + [E K_h((x-y)/h)]^2 - 2 E K_h((x-y)/h) E K_h((x-y)/h) \} = \\ &= n^{-1} \{ E K_h^2((x-y)/h) - [E K_h((x-y)/h)]^2 \} \\ \text{var } \hat{f}(x) &= n^{-1} \int h^{-2} K^2((x-y)/h) f(y) dy - n^{-1} \{ \int h^{-1} K((x-y)/h) f(y) dy \}^2 \end{aligned}$$

Se efectuează aceeași schimbare de variabilă și se folosește descompunerea Taylor la fel ca mai sus:

[‡] Demonstrațiile urmăresc Silverman, 1986, pg38.

$$\begin{aligned} \text{var } \hat{f}(x) &= n^{-1} \int h^{-1} K^2(t) f(x-h t) dt - n^{-1} \{f(x) + \text{bias}(x)\}^2 = \\ &= n^{-1} h^{-1} \int K^2(t) \left\{ f(x) - h t f'(x) + \frac{1}{2} h^2 t^2 f''(x) \dots \right\} dt - n^{-1} \{f(x) + \text{bias}(x)\}^2 \end{aligned}$$

La descompunerea Taylor, cel de-al doilea termen este zero, iar pentru n mari și intervale h mici, cel de-al treilea este aproximativ egal cu zero. Cea de a doua expresie are ordin de mărime n^{-1} și deci poate fi aproximată și ea cu zero:

$$\text{var } \hat{f}(x) \cong n^{-1} h^{-1} \int K^2(t) f(x) dt - O(n^{-1}) \cong n^{-1} h^{-1} f(x) \int K^2(t) dt$$

Expresia obținută pentru deplasare nu conține n deci deplasarea este influențată numai de dimensiunea intervalului și de forma funcțională a kernelului. Variația estimatorului este funcție de dimensiunea intervalului, forma funcțională a kernelului și numărul de observații. Cu cât numărul de observații este mai mare, ceteris paribus, estimatorul are variația mai mică. Dimensiunea intervalului influențează deplasarea și variația în mod opus. Cu cât dimensiunea intervalului este mai mare, deplasarea crește, dar variația scade, ceea ce sugerează prezența unui compromis între dimensiunea deplasării și a variației la alegerea dimensiunii intervalului.

Proprietățile asimtotice ale estimatorului sunt cele care indică evoluția deplasării și a variației cu creșterea numărului de observații. Consistența estimatorului este definită astfel:

Dacă o funcție kernel satisface următoarele condiții[§]:

- $\int |K(u)| du < \infty$ and $\int K(u) du = 1$
- $\lim_{|u| \rightarrow \infty} uK(u) = 0$
- $EY^2 < \infty$
- $n \rightarrow \infty, h_n \rightarrow 0, nh_n \rightarrow \infty$

dacă f este continuu în x , $\hat{f}(x) \xrightarrow{p} f(x)$ atunci când $n \rightarrow \infty$.

Majoritatea funcțiilor Kernel satisfac condițiile de mai sus, ceea ce înseamnă că estimatorul densității kernel este un estimator consistent.

Deplasarea nu depinde de numărul de observații. Consistența este proprietatea care ne oferă garanția că odată cu creșterea numărului de observații, estimatorii sunt din ce în ce mai apropiați de valoarea reală. Această proprietate nu spune însă nimic de viteza de convergență.

[§] Această proprietate a fost demonstrată de Parzen în 1962.

2.3 Selectarea dimensiunii optime a intervalului și a formei funcției kernel

Dimensiunea optimă a intervalului este, așa cum am arătat deja ca un compromis între deplasare și variație. Pentru a se ajunge la o dimensiune optimă trebuie aleasă o funcție obiectiv de deplasare care să fie minimizată. Un candidat pentru această funcție obiectiv este MISE. Aceasta depinde atât de alegerea intervalului maxim cât și de forma funcțională a funcției de ponderi. Prin introducerea expresiei deplasării și variației în formula lui MISE se obține următoarea expresie:

$$\begin{aligned} MISE(x) &= \left\{ \frac{1}{2} h^2 k_2 \right\}^2 \int f''(x)^2 dx + n^{-1} h^{-1} \int f(x) dx \int K^2(t) dt = \\ &= \frac{1}{4} h^4 k_2^2 \int f''(x)^2 dx + n^{-1} h^{-1} \int K^2(t) dt \end{aligned}$$

Minimizarea MISE conduce la următoarea expresie pentru dimensiunea optimă a intervalului:

$$h_{opt} = k^{-2/5} n^{-1/5} \left\{ \int K^2(t) dt \right\}^{1/5} \left\{ \int f''(x)^2 dx \right\}^{-1/5}$$

Intervalul optim este dependent de forma funcțională a kernelului precum și de densitatea care trebuie estimată prin expresia: $\int f''(x)^2 dx$. Ca să se obțină o formulă analitică pentru calculul intervalului optim este necesar să se facă niște ipoteze legate de distribuția densității.

Derivarea intervalului optim se va face pornind de la ipoteza că densitatea de estimat este normală. Alternativ, există și metode iterative, care nu necesită specificarea formei densității, pentru obținerea dimensiunii optime a intervalului.

Introducând expresia dimensiunii optime a intervalului în formula MISE se ajunge la următoarea expresie:

$$MISE = \frac{5}{4} n^{-4/5} k_2^{2/5} \left\{ \int f''(x)^2 dx \right\}^{1/5} \left\{ \int K^2(t) dt \right\}^{4/5}$$

Pentru un interval de dimensiune optimă, singura modalitate de a scădea MISE este prin alegerea unei forme funcționale a kernelului care minimizează MISE. În formula MISE singura parte care depinde de alegerea kernelului este următoarea:

$$C(K) = k_2^{2/5} \left\{ \int K^2(t) dt \right\}^{4/5}$$

Kernelul optim este cel pentru care se obține valoarea MISE cea mai scăzută, condiție care se reduce la kernelul cu valoarea $C(K)$ cea mai mică. Kernelul Epanechnikov este cel care are valoarea $C(K)$ cea mai redusă.

Comparația dintre diverse forme funcționale de kerneli se face prin calcularea raportului $C(K)$ față de $C(K)$ al kernelului Epanechnikov. Acest raport este destul de aproape de 1, în cazul densității normale este 0.95. Concluzia care se impune este că posibilitățile de îmbunătățire a MISE prin alegerea unei anumite forme de kernel sunt limitate. Din acest

motiv considerentele care predomină la alegerea formei funcționale a ponderilor sunt legate de aspectele computaționale.

Vom deriva în continuare formula pentru dimensiunea intervalului optim în cazul în care densitatea ce se dorește estimată este normală. Densitatea normală are următoarea formă funcțională:

$$f(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-(x-\mu)^2/2\sigma^2}$$

$$\int f''(x)^2 dx = \int \frac{1}{2\pi\sigma^2} \frac{(x-\mu)^4}{\sigma^8} e^{-(x-\mu)^2/\sigma^2} dx$$

Se efectuează schimbarea de variabilă $z=(x-\mu)/\sigma$, și $dz=dx/\sigma$ și se obține următoarea expresie pentru integrală:

$$\int f''(x)^2 dx = \int \frac{1}{2\pi\sigma^5} z^4 e^{-z^2} dz = \frac{1}{\sigma^5} \int \frac{1}{2\pi} z^4 e^{-z^2} dz = \frac{1}{\sigma^5} \int z^4 \phi''(z) dz$$

unde $\phi''(z)$ este funcția de densitate standard normală.

Rezolvarea integralei se face prin părți, ținând cont și de proprietatea că: $\int \phi''(z) dz = 1$.

$$\int f''(x)^2 dx = \frac{3}{8\pi\sigma^5}$$

Prin înlocuirea valorii integralei în formula dimensiunii optime a intervalului se obține expresia intervalului optim a lui Silverman, care are următoarea formă:

$$h_{opt} = 1.364 \sigma \cdot n^{-1/5} \cdot \left(\frac{\int K^2(t) dt}{k_2^2} \right)^{1/5}$$

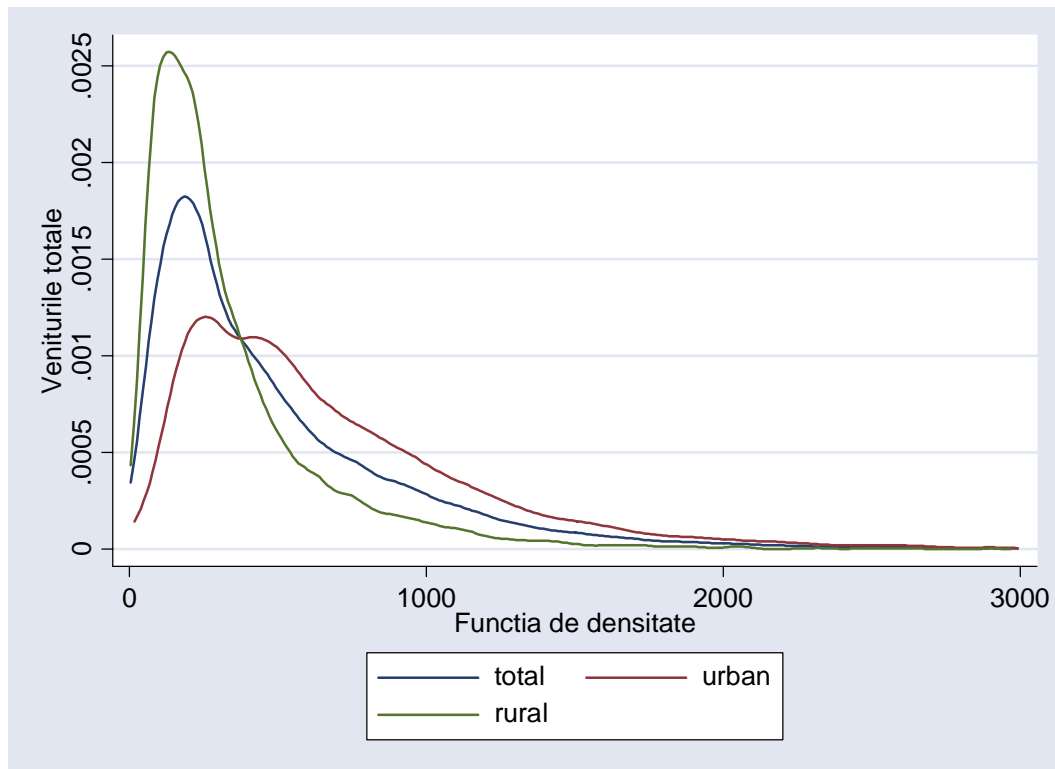
Derivarea expresiei intervalului optim Silverman este făcută în cazurile în care distribuția care se dorește estimată este normală, dar se poate folosi fără erori majore în toate cazurile în care distribuția este uni-modală. Se poate spune deci că într-un număr mare de cazuri aproximarea Silverman este suficientă și dă un punct de plecare bun pentru estimarea densităților.

2.4 Reprezentarea distribuțiilor variabilelor de interes

Estimatorul densității neparametrice a fost folosit pentru a studia consumul gospodăriilor din România. Datele analizate provin din Ancheta Integrată în Gospodării 2003, un sondaj pe care Institutul de Statistică îl face în fiecare an, începând cu 1995.

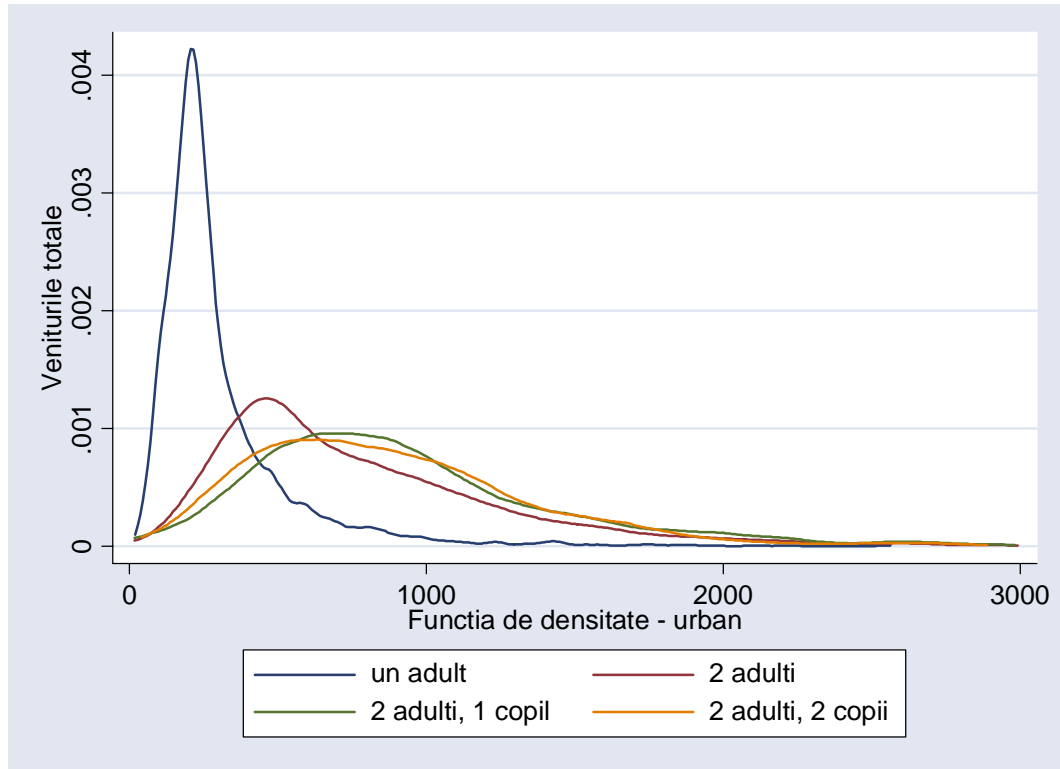
Graficul 2-1 prezintă funcția de densitate a veniturilor totale pentru toate gospodăriile și separat pentru cele din rural și urban. Se poate observa o diferență importantă între cele două distribuții, veniturile din zona urbană sunt, sensibil mai mari decât cele din zona rurală. Diferența poate fi cauzată de mai mulți factori, și nu avem suficiente informații

pentru a identifica cu exactitate care factor este responsabil cu discrepanțele privind veniturile între mediu urban și rural. Dintre factorii responsabili menționăm, caracteristici demografice diferite ale gospodăriilor din cele două medii, care ar duce la o structură diferită a venitului total (mai mulți membri în vârstă de muncă, deci mai multe salarii), caracteristici ale forței de muncă diferite în cele două regiuni, ca de exemplu nivel de educație, calificări, etc. (salarii mai mari pentru salariații din zona urbană datorită prestării unei activități mai productive), dar și din punct de vedere al ocupării și nu în ultimul rând posibile diferențe salariale între cele două zone, chiar pentru munci similare, datorită condițiilor de pe piața muncii.

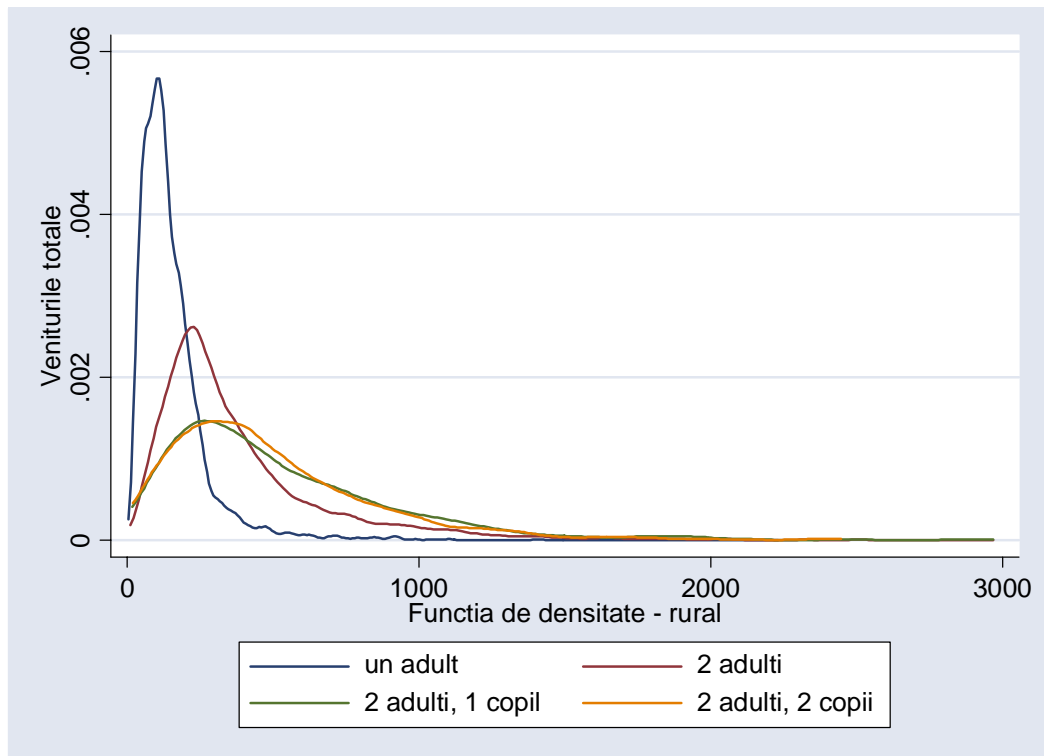


Graficul 2-1 Funcția de densitate a veniturilor totale (RON)

Următoarele două grafice investighează ipoteza că diferențele între veniturile totale ale familiilor din cele două regiuni ar putea fi explicate de diferențe în numărul de persoane care compun familia. Și în acest sens a fost reprezentată grafic funcția de densitate separat pentru zona urbană și rurală și separat pentru gospodăriile formate dintr-un adult, doi adulți, doi adulți și un copil, și doi adulți și doi copii, aceste tipuri de familii fiind cele mai reprezentate în eșantion.

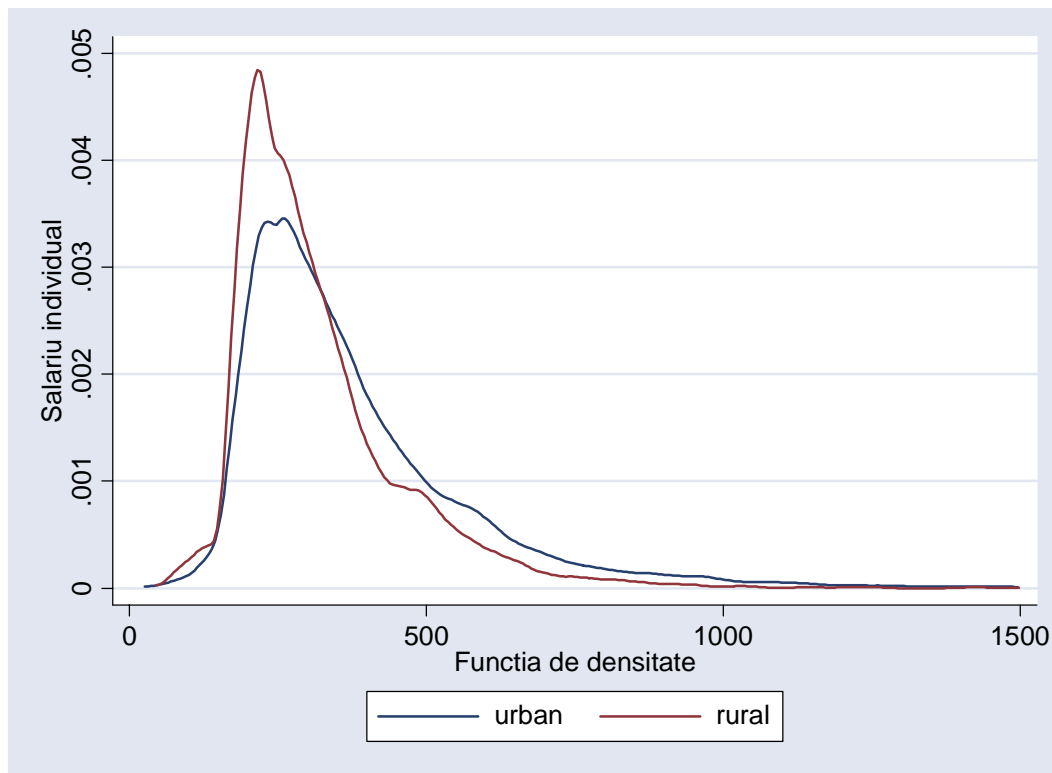


Graficul 2-2 Veniturile totale în funcție de componența familiei, în zona urbană (RON)



Graficul 2-3 Veniturile totale în funcție de componența familiei, în zona rurală (RON)

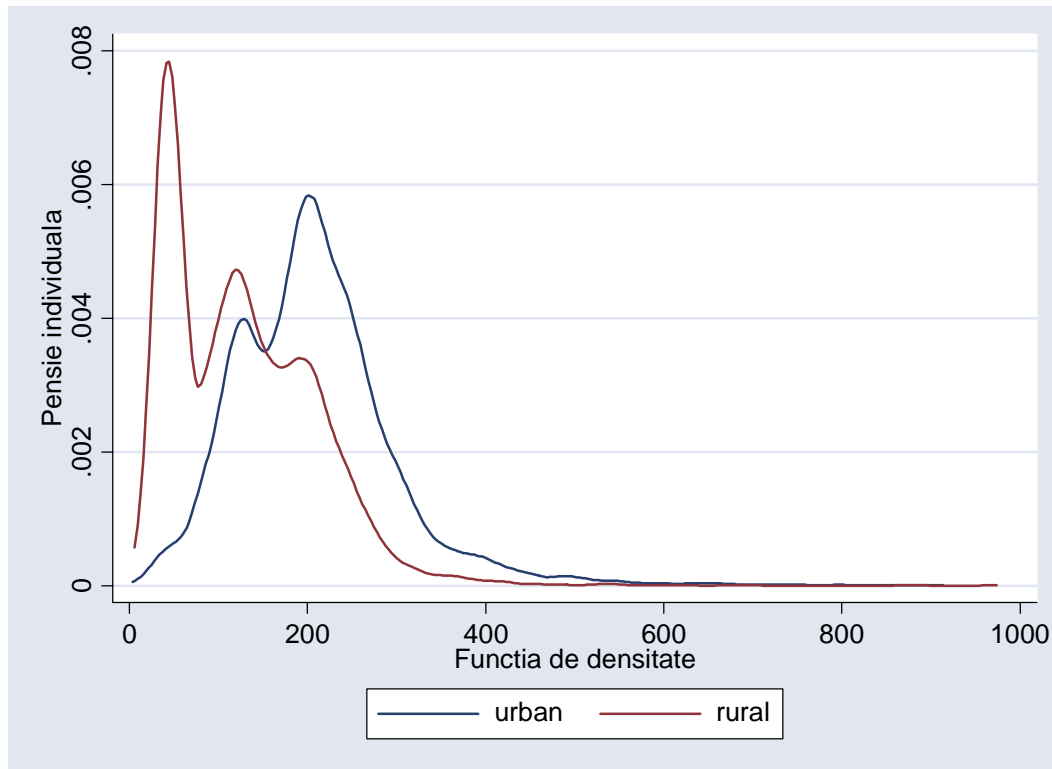
Din Graficul 2-2 și Graficul 2-3 se poate observa că atunci când se controlează pentru structura familiei diferențele privind veniturile totale se mențin, cu alte cuvinte caracteristicile demografice ale populației din zona urbană și rurală din punct de vedere al numărului de membri ai familiei nu au o influență decisivă asupra diferențelor de venituri între cele două zone.



Graficul 2-4 Funcția de densitate a salariilor pentru urban și rural (RON)

Salariile individuale sunt prezentate în Graficul 2-4. Din compararea funcției de distribuție pentru salariilor urbane și rurale se pot observa anumite diferențe între cele două funcții: o dispersie mai mare a salariilor urbane față de cele urbane (distribuția salariilor rurale este mai înaltă ceea ce indică un număr mai mare de indivizi cu salarii apropiate de medie), și de asemenea o medie marginal mai mare a acestora. Totuși distribuțiile salariale sunt relativ apropiate, și nu pot fi la originea diferențelor care au fost observate la veniturile totale.

Graficul 2-5 prezintă distribuția pensiilor urbane și rurale.



Graficul 2-5 Funcția de densitate a pensiilor pentru urban și rural (RON)

În cazul pensiilor situația este alta în sensul că funcțiile de distribuție ale pensiilor pentru zona urbană și cea rurală sunt foarte diferite. În cazul pensiilor rurale distribuția este tri-modală, cu vârful cel mai mare la valori ale pensiilor foarte mici, care corespund pensiilor lucrătorilor în agricultură. distribuția pensiilor urbane are numai două vârfuri, datorită absenței pensiilor agricole în zona urbană. Al doilea vârf al distribuției pentru zona rurală corespunde primului vârf al distribuției pensiilor urbane din punct de vedere al valorii pensiei. La pensia urbană se poate observa că distribuția atinge maximul la valori ale pensiilor superioare pensiilor rurale.

Două concluzii pot fi deduse din cele două distribuții, în primul rând **diferențele dintre venituri pot fi explicate măcar parțial prin diferențele de pensie**, dar motivația existenței diferențelor între pensiile urbane și rurale stă în statutul ocupațional al persoanelor care locuiesc în cele două regiuni. O proporție mai mare de persoane din zona rurală au lucrat și încă lucrează în agricultura de subsistență, activitate care generează pensii și venituri mici. Interesant este că în rest **diferențele între salarii și pensii între cele două zone sunt marginale**. De aceea se pare că în zona urbană o proporție mai mare de persoane active este responsabilă pentru discrepanțele dintre veniturile totale ale gospodăriilor.

3 Regresia neparametrică

La regresia neparametrică se dorește să se modeleze relația dintre variabila dependentă și un set de variabile explicative. Vom începe cu cazul unei singure variabile dependente și explicative. Pentru exemplificare să considerăm modelul:

$$Y_i = m(X_i) + \varepsilon_i$$

Se dorește derivarea curbei de regresie care poate fi scrisă astfel:

$$m(x) = E(Y|X=x)$$

Cel mai simplu estimator poate fi obținut prin calculul mediei variabilei dependente pe diferite intervale. În acest caz, funcția estimată este constantă pe intervalul respectiv, dar are discontinuități la capetele intervalului. Acesta nu este un estimator potrivit, decât în cazul funcțiilor care manifestă asemenea salturi. Acest estimator este o generalizare a histogramei, și este denumit, prin analogie, regresogramă.

Dacă la fiecare valoare a variabilei explicative ar exista mai multe valori pentru variabila dependentă, estimatorul poate fi calculat ca medie a variabilei răspuns în fiecare punct. În cele mai frecvente cazuri, există numai o observație pentru fiecare valoare, în acest caz, valoarea funcției este estimată prin calculul mediei ponderate a variabilei dependente într-un interval mic în jurul punctului. Acest estimator specific este denumit estimatorul mediei mobile.

Deși regresograma și estimatorul mediei mobile sunt calculați similar prin medierea variabilei de răspuns pe un interval, există o diferență majoră între cei doi estimatori. Regresograma are valoare constantă pe un interval, numărul de intervale considerate dă numărul de valori diferite pe care le poate lua funcția. Estimatorul mediei mobile calculează valoarea funcției în fiecare punct al variabilei independente prin medierea unui anumit număr de valori ale variabilei dependente, deci în timp ce regresograma este o funcție care are salturi pe intervale, estimatorul mediei mobile este o funcție continuă. .

3.1 Estimatorul Kernel

Estimatorul kernel pentru dependențele dintre două variabile aplică principiile estimatorului kernel al densității pentru estimarea dependențelor dintre două sau mai multe variabile. Ca și în cazul estimatorului kernel al densității, estimatorul calculează o medie ponderată a variabilei dependente în fiecare punct.

$$\hat{m}(x) = \frac{1}{n} \cdot \sum_{i=1}^n W_{ni}(x) Y_i$$

iar ponderile $\{W_{ni}(x)\}_{i=1}^n$ sunt calculate cu ajutorul funcției kernel astfel:

$$W_{mi}(x) = \frac{K_{hn}(x - X_i)}{\hat{f}_{hn}(x)}$$

unde $\hat{f}_h(x)$ este estimatorul kernel al funcției de densitate și $K(\bullet)$ este funcția kernel:

$$\hat{f}_h(x) = \frac{1}{n} \cdot \sum_{i=1}^n K_h(x - X_i)$$

$$K_u(\bullet) = h^{-1} \cdot K(\bullet/h).$$

Estimarea constă în calculul unor ponderi pentru fiecare observație din intervalul $[x-h, x+h]$ la fiecare valorare a variabilei explicative. Cu excepția cazurilor în care kernelul este o funcție constantă, adică ponderile sunt egale pentru fiecare observație din intervalul $[x-h, x+h]$, ponderile depind de distanța punctelor față de x . Deoarece suma ponderilor trebuie să fie unu, acestea sunt împărțite cu densitatea variabilei explicative în acel punct.

După efectuarea tuturor substituțiilor estimatorul kernel ia forma următoare:

$$\hat{m}_h(x) = \frac{\frac{1}{n} \sum_{i=1}^n K_h(x - X_i) Y_i}{\frac{1}{n} \sum_{i=1}^n K_h(x - X_i)}$$

Acest estimator este cunoscut și sub numele estimatorul Nadaraya – Watson, botezat cu numele primelor persoane care au recomandat folosirea acestuia.

Proprietățile funcțiilor kernel au fost descrise în secțiunea anterioară, dar reamintim aici pe cele importante: funcțiile kernel sunt simetrice, pozitive (pentru kernelii de ordin 2, cei mai uzuali în practică) și integrala lor trebuie să fie unu. Datorită aspectelor computaționale, cei mai folosiți sunt kernelii cu suport finit.

3.2 Proprietățile statistice ale estimatorului Nadaraya – Watson

În această secțiune se vor investiga proprietățile statistice și asimtotice ale estimatorului Nadaraya – Watson. La derivarea formulei deplasării și a variației, se va urmări raționamentul prezentat în Scott (1992) pg. 223. Derivarea formulei deplasării și a variației este mai dificilă în acest caz datorită formei estimatorului definit ca raport a două variabile aleatoare corelate.

În derivații s-a folosit următoarea proprietate:

Dacă numărătorul și numitorul unei fracții converg către o constantă pozitivă, atunci expectațiile asimptotice ale raportului este egal cu raportul expectațiilor asimptotice ale numărătorului și numitorului^{**}.

$$\hat{m}_h(x) = \frac{\frac{1}{n} \sum_{i=1}^n K_h(x - X_i) Y_i}{\frac{1}{n} \sum_{i=1}^n K_h(x - X_i)}$$

Deja se cunoaște valoarea așteptată a numitorului, pentru că acesta este estimatorul kernel a densității, iar formula a fost obținută în secțiunea anterioară.

$$E \hat{f}(x) = f(x) + \frac{1}{2} h^2 f''(x) k_2$$

$$\text{var} \hat{f}(x) = n^{-1} h^{-1} f(x) \int K^2(t) dt$$

În continuare vom calcula valoarea așteptată pentru numărător:

$$\begin{aligned} E \left(\frac{1}{n} \sum_{i=1}^n y_i \frac{1}{h} K \left(\frac{x - x_i}{h} \right) \right) &= \frac{1}{n} \sum_{i=1}^n y_i E \left(\frac{1}{h} K \left(\frac{x - x_i}{h} \right) \right) = \sum_{i=1}^n y_i E \left(\frac{1}{h} K \left(\frac{x - x_i}{h} \right) \right) \\ &= \iint y \frac{1}{h} K \left(\frac{x - z}{h} \right) f(y, z) dy dz \end{aligned}$$

unde $f(y, z)$ este funcția dublă de densitate a lui x și y .

Se efectuează următoarea schimbare de variabilă $s = (x - z)/h$, cu $ds = dz/h$ și funcția poate fi rescrisă:

$$E(\bullet) = \iint y K(s) f(x - hs, y) dy ds$$

Densitatea dublă poate fi scrisă cu ajutorul densității condiționate astfel:

$$f(x - hs, y) = f(y|x - hs) f(x - hs)$$

Cu această substituție variabilele pot fi separate în felul următor:

$$E(\bullet) = \int K(s) f(x - hs) \left(\int y f(y|x - hs) dy \right) ds = \int K(s) f(x - hs) m(x - hs) ds$$

Folosind descompunerea Taylor, integrala poate fi rescrisă:

$$\begin{aligned} E(\bullet) &= \int K(s) [f(x) - f'(x)hs + f''(x)(hs)^2/2] [m(x) - m'(x)hs + m''(x)(hs)^2/2] ds = \\ &= \int K(s) ds - h [f' m' + f'' m] \int K(s) ds + h^2/2 [f''' m + f'' m'' + 2f'' m'] \int s^2 K(s) ds \\ &= f(x) m(x) + h^2/2 [f'''(x) m(x) + f''(x) m''(x) + f'(x) m'(x)] k_2 \end{aligned}$$

Valoare așteptată a estimatorului va fi:

$$E \hat{m}(x) = \frac{f(x) [m(x) + h^2 k_2 (f'''(x) m(x) / 2 f(x) + m''(x) / 2 + f'(x) m'(x) / f(x))]}{f(x) [1 + h^2 k_2 f''(x) / 2 f(x)]}$$

^{**} Scott pg. 222.

Se știe că pentru a suficient de mic $1/(1+a) \cong 1-a$, folosind această aproximare expresia de mai sus poate fi simplificată:

$$E\hat{m}(x) = m(x) + h^2 k_2 / 2 [m''(x) + 2m'(x)f'(x)/f(x)]$$

deci deplasarea poate fi scrisă:

$$\text{deplasarea } \hat{m}(x) = h^2 k_2 / 2 [m''(x) + 2m'(x)f'(x)/f(x)]$$

În cazurile când punctele sunt uniform distribuite pe domeniul funcției, funcția de densitate este o constantă, derivata de ordinul unu este zero, și termenul al doilea din paranteze este zero, iar formula deplasării poate fi simplificată. În realitate, în afara experimentelor din fizică, este foarte dificil de a avea control asupra distribuției punctelor x , mai ales când variabilele sunt obținute în urma unor anchete.

Formula variației poate fi obținută în mod similar cu cea a deplasării:

$$\text{var } \hat{m}(x) = \frac{\sigma_\varepsilon^2 \int K^2(u) du}{nhf(x)}$$

În acest moment se poate calcula MSE pentru estimatorul Nadaraya – Watson:

$$\text{MSE } \hat{m}(x) = \frac{\sigma_\varepsilon^2 \int K^2(u) du}{nhf(x)} + h^4 k_2^2 [m''(x) + 2m'(x)f'(x)/f(x)]^2$$

Pe lângă proprietățile statistice ne interesează și proprietățile asimptotice ale estimatorului. Estimatorul este deplasat, și în formula deplasării influența creșterii numărului de observații nu este vizibilă. De aici ar părea că creșterea numărului de observații nu are nici o influență asupra scăderii deplasării. Acesta este motivul pentru care proprietatea de consistență a estimatorului este atât de importantă.

Condițiile^{††} care asigură consistența estimatorului sunt enunțate în următoarea propoziție:

Dacă condițiile următoare sunt satisfăcute:

$$- \int |K(u)| du < \infty$$

$$- \lim_{|u| \rightarrow \infty} uK(u) = 0$$

$$- EY^2 < \infty$$

$$- n \rightarrow \infty, h_n \rightarrow 0, nh_n \rightarrow \infty$$

atunci în fiecare punct de continuitate a lui $m(x)$, $f(x)$ și $\sigma^2(x)$ unde $f(x) > 0$:

$$n^{-1} \sum_{i=1}^n W_{hi}(x) Y_i \xrightarrow{p} m(x)$$

^{††} Vezi Hardle pag. 29.

Această propoziție stipulează că odată cu creșterea numărului de observații, estimatorul Nadaraya – Watson converge în probabilitate către estimatorul real. Această proprietate este foarte importantă în special datorită deplasării estimatorului, recomandând folosirea acestuia atunci când există un număr suficient de observații. Demonstrația^{††} propoziției este complicată datorită exprimării estimatorului ca un raport a două variabile aleatoare corelate.

3.3 Alegerea gradului de netezire – metoda validării

Pentru a putea face o alegere cât mai obiectivă, este nevoie de un procedeu de selecție care să poată ordona după anumite criterii variantele posibile. Instrumentul folosit în general este minimizarea unei funcții obiectiv. În contextul stabilirii dimensiunii intervalului optim la estimarea densității, funcția de minimizat era integrala mediei pătratelor erorilor (MISE). Datorită dificultății obținerii unei expresii explicite pentru intervalul optim, una care să nu depindă de densitatea estimată, și datorită faptului că expresia lui MISE este mult prea complicată este necesară găsirea altor criterii pentru selectarea intervalului optim.

Procedura alternativă pentru alegerea intervalului optim, care va fi descrisă în continuare, poartă numele de metoda validării (cross validation). Aceasta constă în minimizarea mediei *MSE*.

$$MSE(\lambda) = 1/n \sum E\{\hat{m}(x) - m(x)\}^2$$

Deoarece funcția care se dorește estimată nu este cunoscută, în locul lui *MSE* se folosește o funcție care o aproximează. Cel mai des folosit estimator al lui *MSE* este media pătratelor reziduurilor *ASR* definit sub forma următoare:

$$ASR(\lambda) = 1/n \sum \{y_i - \hat{m}(x_i)\}^2$$

În mod intuitiv, *ASR* nu este un estimator bun pentru *MSE*, pentru că minimizând *ASR* se obține un interval optim care duce la estimarea unei funcții care trece prin toate punctele observate, în care caz *ASR* este zero. Metoda de validare constă în calcularea valorii funcției în punctul x_i eliminând observația i , folosind restul de puncte. După care, pentru diferite intervale optime λ se calculează funcția *CV*. Se alege ca interval optim, valoarea λ care minimizează funcția *CV*.

Expresia pentru *CV* este următoarea:

$$CV(\lambda) = 1/n \sum \{y_i - \hat{m}^{-i}(x_i)\}^2$$

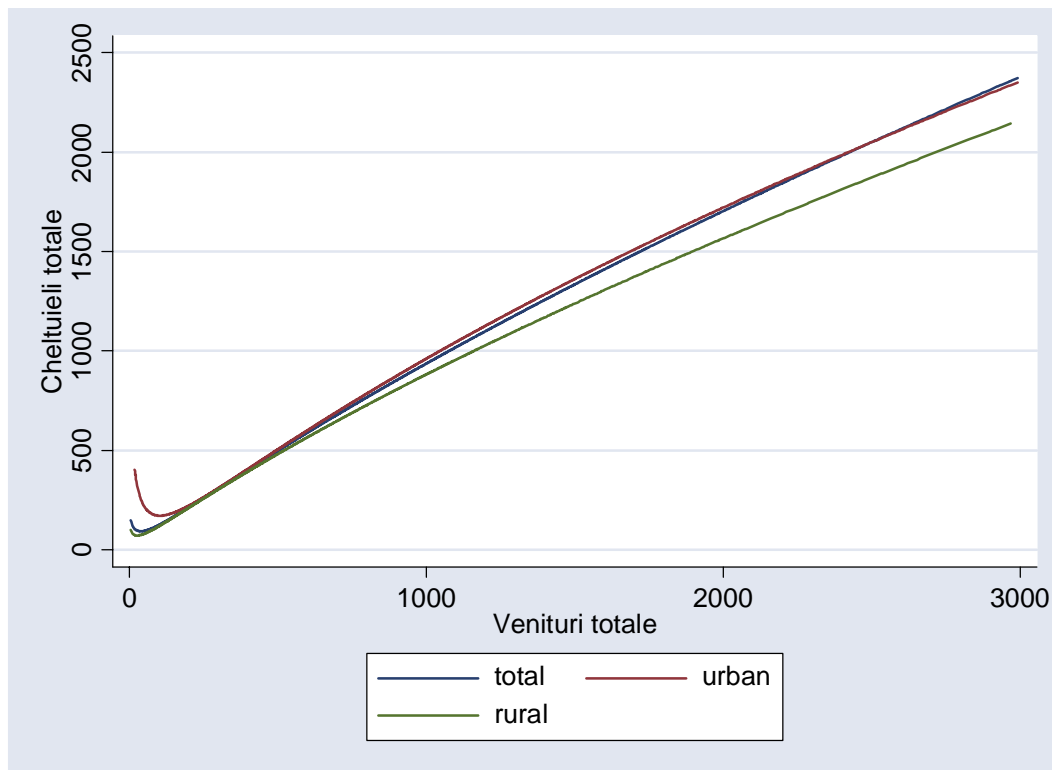
Există și alte metode pentru alegerea dimensiunii intervalului optim (statistica C_p , etc.) dar din compararea celor două curbe folosind metodele diferite s-a constatat că rezultatele obținute sunt similare.

^{††} Pentru demonstrația propoziției vezi Hardle pag. 39.

3.4 Dependența cheltuielilor de veniturile gospodăriilor

În această secțiune, estimatorii neparametrici sunt aplicați pentru a studia dependențele dintre variabile, ne interesează în principal modul de variație al diverselor cheltuieli cu veniturile gospodăriilor.

Înainte de a începe prezentarea graficelor, trebuie menționat că metodele neparametrice funcționează cel mai bine acolo unde există un număr suficient de mare de observații. În zonele în care numărul de observații este mai redus, așa cum este la extremitățile domeniului de referință a funcției (cozile distribuției de repartiție) valoarea estimatorului este mai instabilă, și mai sensibilă la prezența valorilor extreme. Din acest motiv, atunci când analizăm graficele trebuie să fim conștienți de limitările estimatorului neparametric.

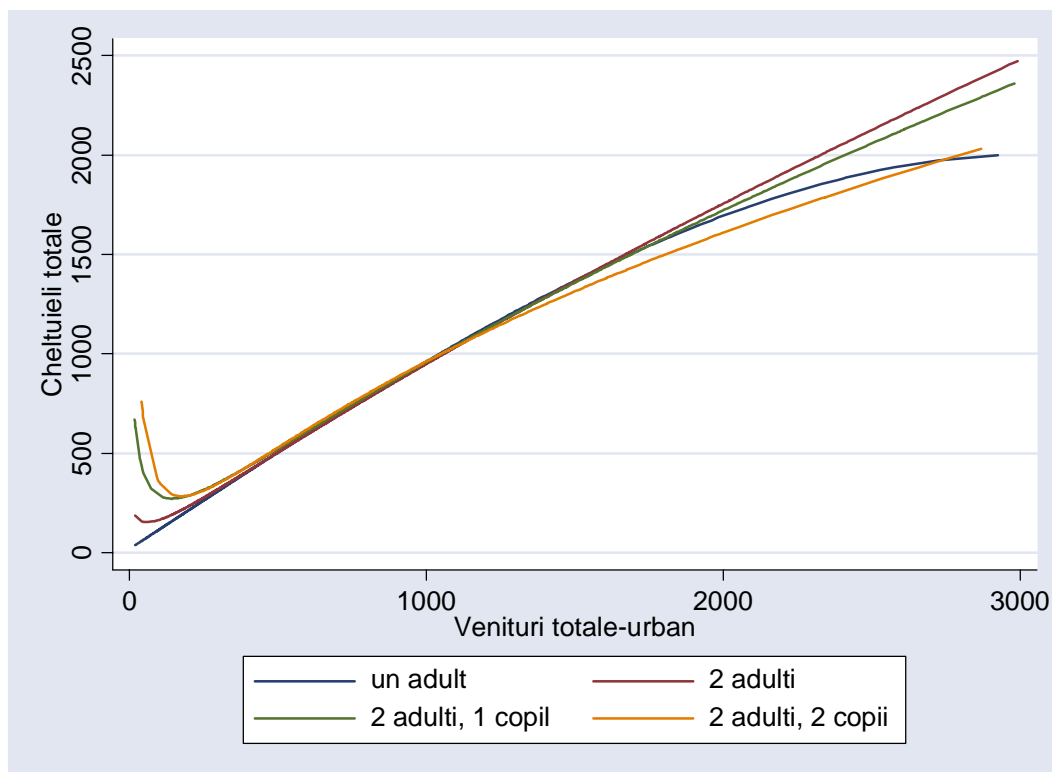


Graficul 3-1 Dependența cheltuielii totale- venituri totale pentru total populație, pentru populația urbană și populația rurală

Prima dependență ce va fi analizată este relația dintre cheltuielile totale și veniturile totale pentru toată populația, și separat pentru populația urbană și rurală. Ca o primă observație se poate remarca o dependență aproape liniară a cheltuielilor de venituri, cel puțin pentru venituri mai mici, cu propensitatea de consum mare, în jurul unității. Aceasta înseamnă că o proporție constantă din venituri este folosită pentru cheltuieli de către cea mai mare parte a gospodăriilor cu venituri mai mici. Însă odată cu creșterea veniturilor, cheltuielile dau semne că se aplatizează, sugerând că odată cu creșterea veniturilor proporția care se

cheltuie scade, deci propensitatea de consum a gospodăriilor se micșorează, un obicei de consum care este conform teoriilor economice.

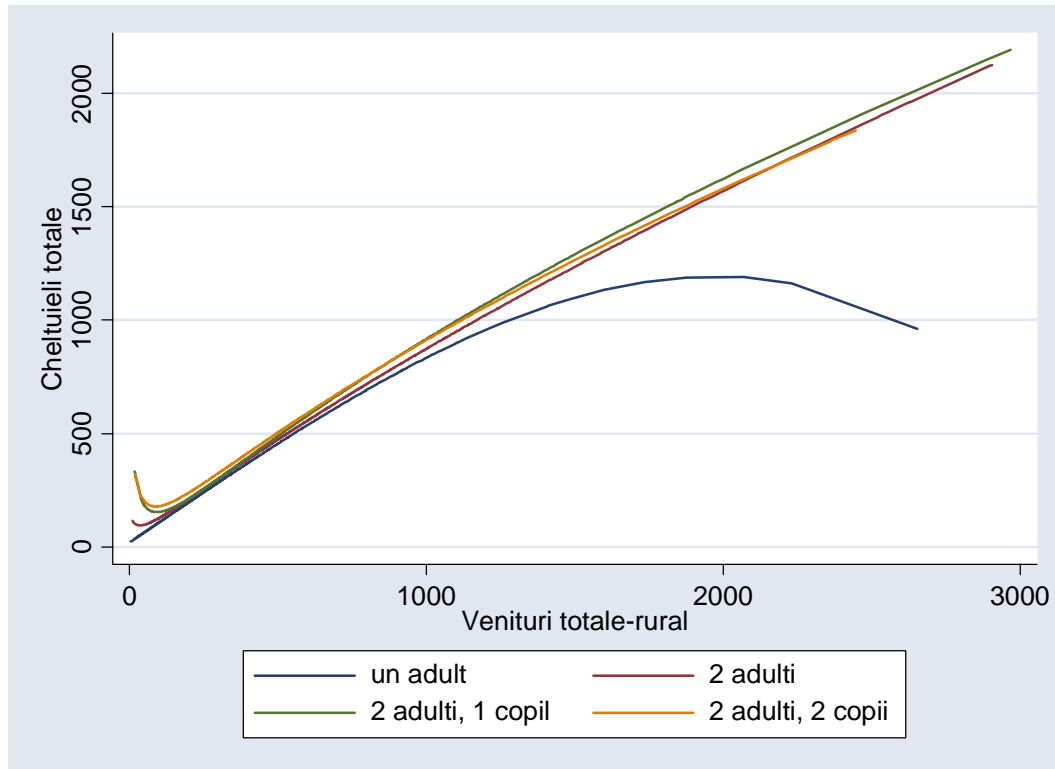
La analiza relației cheltuieli totale – venituri totale pentru cele două medii de rezidență se poate observa la venituri mici aproape o suprapunere a celor trei grafice. Odată cu creșterea veniturilor, cheltuielile totale din zona rurală nu în pasul cu cele din zona urbană, ecartul dintre cele două grafice mărindu-se cu creșterea veniturilor. Explicațiile privind discrepanțele dintre consumul gospodăriilor urbane și al celor rurale pot fi multiple, pornind de la un număr de membri mai reduși, prezența auto-consumului care contribuie la reducerea cheltuielilor alimentare în zona rurală, dar și o preponderență mai mare în mediu rural a grupelor de vârstă care au o propensitate de economisire mai mare.



Graficul 3-2 Relația cheltuieli totale – venituri totale în funcție de componența familiei în urban

Graficul 3-2 și Graficul 3-3 prezintă relația dintre cheltuieli totale și venituri totale în funcție de componența familiei. Mai ales în mediu urban se poate observa o suprapunere a celor patru grafice pentru venituri totale până la 1.500 RON. După acest prag, panta graficelor se modifică, propensitatea de consum scăzând masiv mai ales pentru familiile formate din doi adulți și doi copii. Gospodăriile formate din doi adulți și mențin constantă proporția din venituri destinată consumului, care este totuși subunitară. Gospodăriile formate din doi adulți sunt fie tinere cupluri sau pensionari,

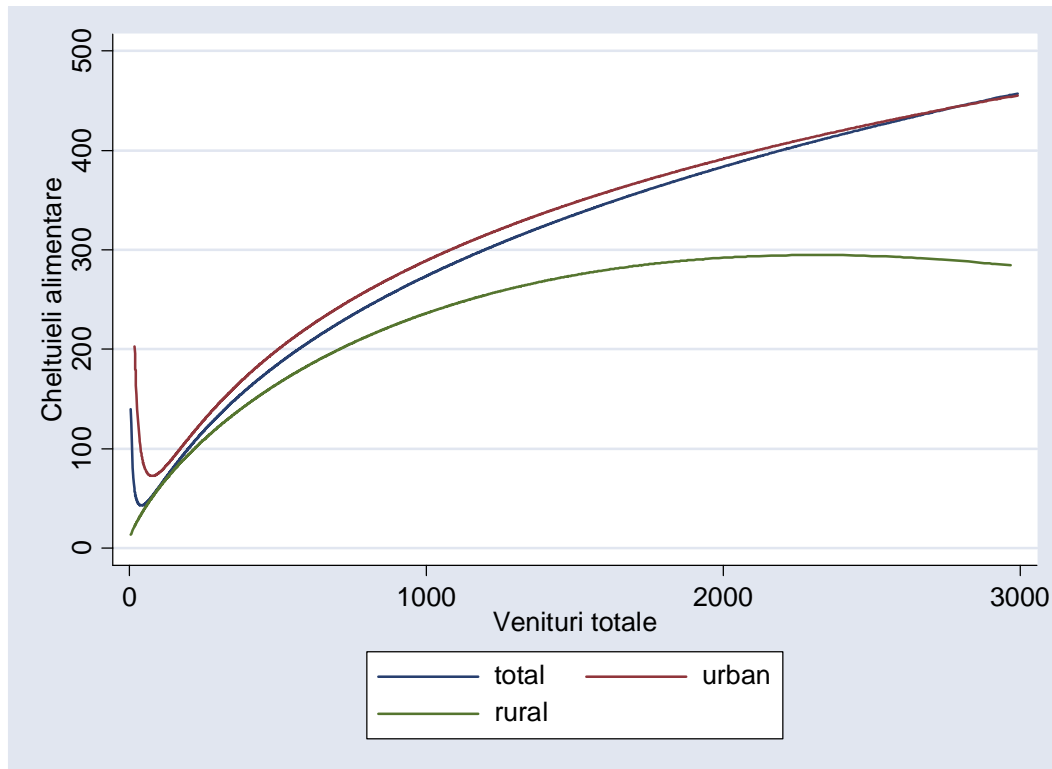
două categorii de populație care nu au motivația de economisire foarte puternică. Prezența copiilor pare să modifice comportamentul adulților, inducând dorința de economisire. Acest lucru este evident însă numai atunci când veniturile sunt suficient de mari să permită atingerea unui standard minim de viață.



Graficul 3-3 Relația cheltuieli totale – venituri totale în funcție de componența familiei în rural

În zona rurală, lucrurile sunt un pic diferite. Se poate remarca o apropiere destul de mare a celor patru grafice în zona cu venituri mici și medii, grupare care se menține și la venituri mari (cu excepția gospodăriilor formate dintr-un adult^{§§}) sugerând că la venituri egale se cheltuie sume similare indiferent de componența familiei. Trebuie remarcat însă o propensitate de consum mai mică decât în mediu urban, adică un comportament similar cu gospodăriile urbane formate din adulți și copii pentru toate tipologiile de gospodării din mediu rural.

^{§§}Dependenței în cazul gospodăriilor formate dintr-un adult la venituri mari nu este relevantă, pentru că în acea zonă graficul suferă de influențe de la niște observații extreme, ceea ce influențează disproporționat forma relației.

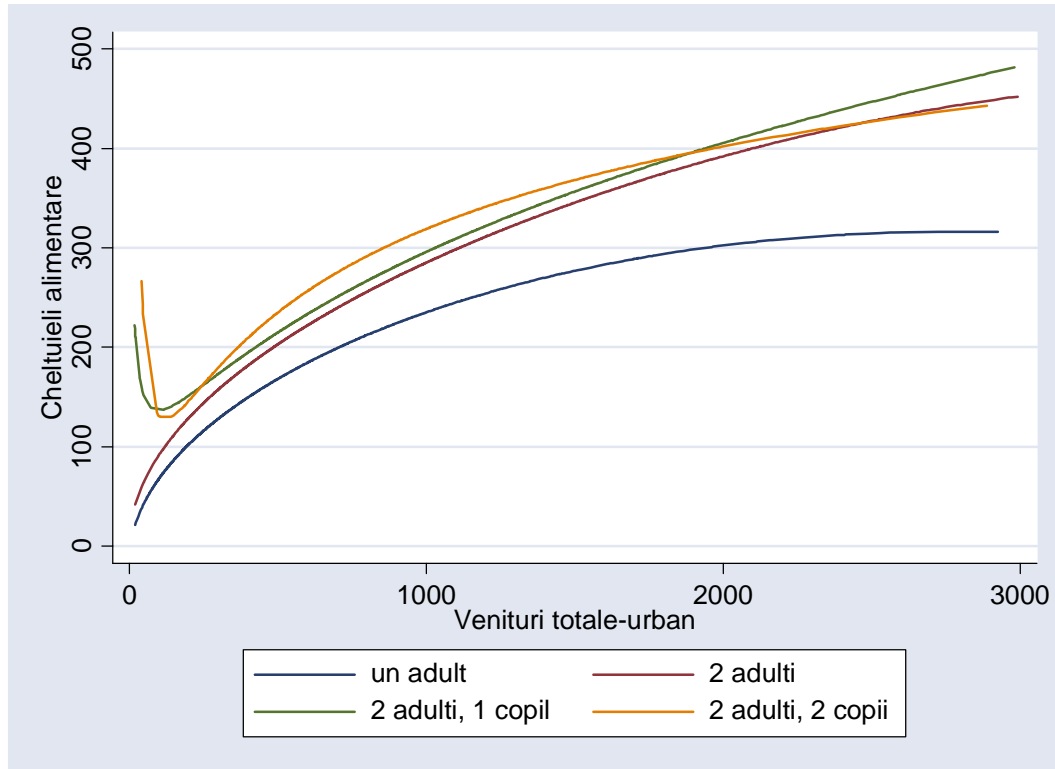


Graficul 3-4 Dependenta cheltuieli alimentare – venituri totale pentru total populație, pentru populația urbană și populația rurală

O primă observație este relația neliniară dintre cheltuielile alimentare și veniturile totale. Deci ipoteza de liniaritate a dependențelor dintre variabile necesară în cazul aplicării estimatorului metodei celor mai mici pătrate nu este îndeplinită. În acest caz este necesar să se liniarizeze dependențele dintre variabile prin aplicarea unor transformări – ca de exemplu logaritizarea – sau prin introducerea de variabile pătratice în venituri.

În cazul cheltuielilor alimentare se pot observa deosebiri importante între comportamentele gospodăriilor din mediu urban față de cele din mediu rural. Deși cheltuielile alimentare dau semne de saturație atât în mediu urban cât și în mediu rural, cele din mediu rural se stabilizează în jurul valorii de 300 RON, în timp ce cele din mediu urban continuă să crească dar cu o pantă mult mai redusă. Prezența în mediu rural a auto-consumului este o parte parțial responsabilă pentru saturația consumului de alimente în mediu rural.

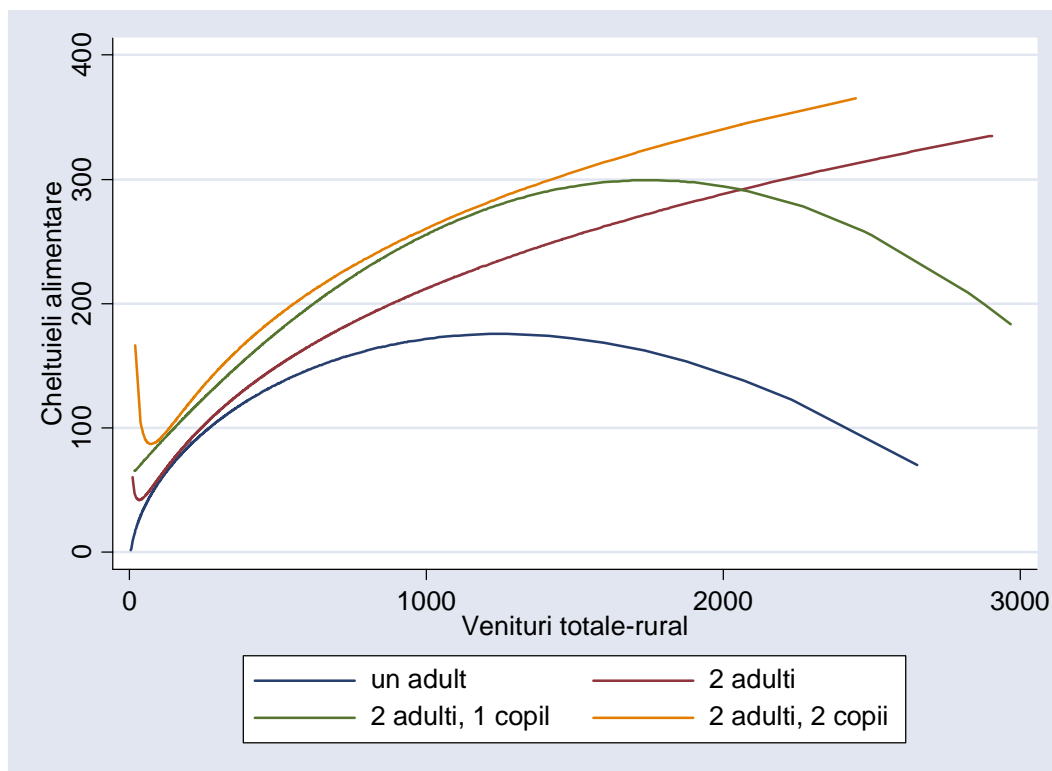
În Graficul 3-5 sunt prezentate cheltuielile alimentare în funcție de numărul de membri pentru zona urbană. Caracteristica de neliniaritate a consumului alimentar se menține și la gospodăriile urbane. Cum este de așteptat, atingerea pragului de saturație se face la venituri mai mici și la cheltuieli mai mici în cazul familiilor formate dintr-un singur adult. Proporția cea mai mare a cheltuielilor alimentare în cazul veniturilor mici și medii o au gospodăriile formate din doi adulți și doi copii, însă acestea se saturează mai rapid decât cheltuielile celorlalte tipuri de familii.



Graficul 3-5 Dependența cheltuieli alimentare – venituri totale în funcție de componența familiei în mediu urban

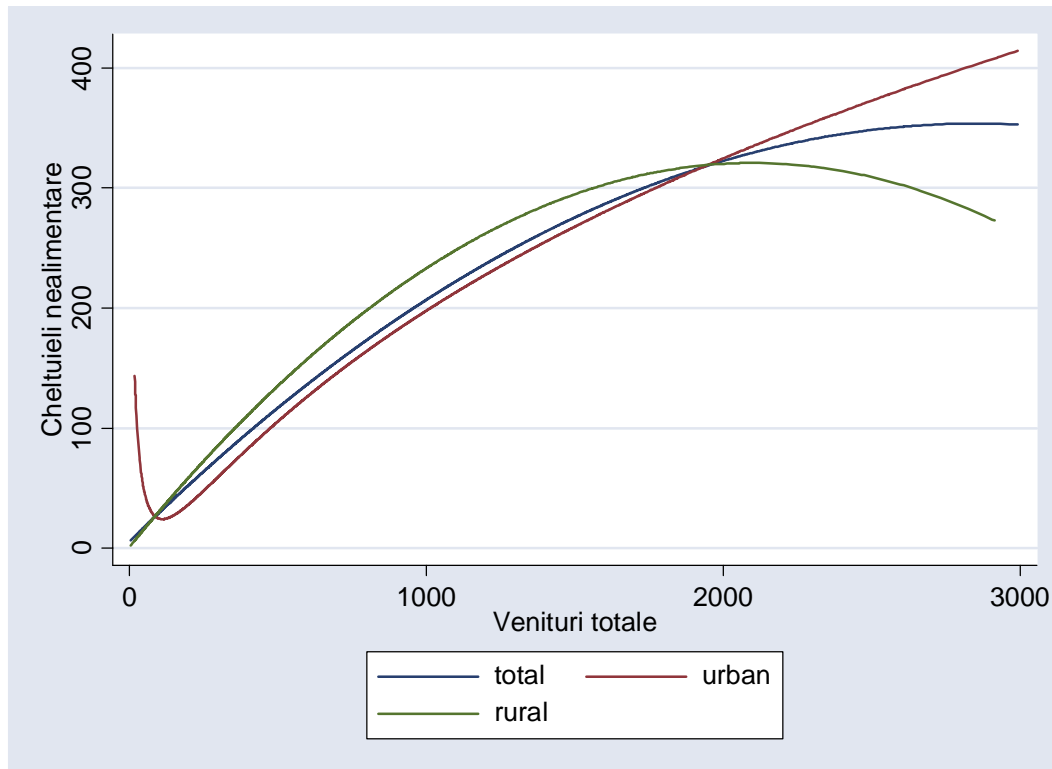
Dacă comparăm relația dintre cheltuieli alimentare – venituri pentru gospodăriile din mediu urban față de cele din mediu rural se poate observa diferența de cheltuieli alimentare la venituri similare pentru același tip de familie, lucru evident încă din Graficul 3-4.

Dacă în zona urbană există o apropiere destul de mare a relației dintre cheltuieli alimentare și venituri la gospodăriile formate din mai mulți adulți, acest lucru nu mai este valabil în cazul familiilor din zona rurală. Cheltuielile cele mai mari alimentare le au, așa cum este de așteptat familiile formate din doi adulți și doi copii, iar pentru veniturile mici și medii, familiile formate din doi adulți și un copil au cheltuieli foarte apropiate, aproape identice. Cheltuielile adulților care trăiesc singuri sunt foarte scăzute, probabil că o mare proporție din aceștia participă în agricultura de subzistență. Forma parabolică a relației pentru un adult și poate și cea pentru doi adulți cu un copil poate fi dată de absența suficientelor observații în zona de venituri mari.



Graficul 3-6 Dependența cheltuieli alimentare – venituri totale în funcție de componența familiei în mediu rural

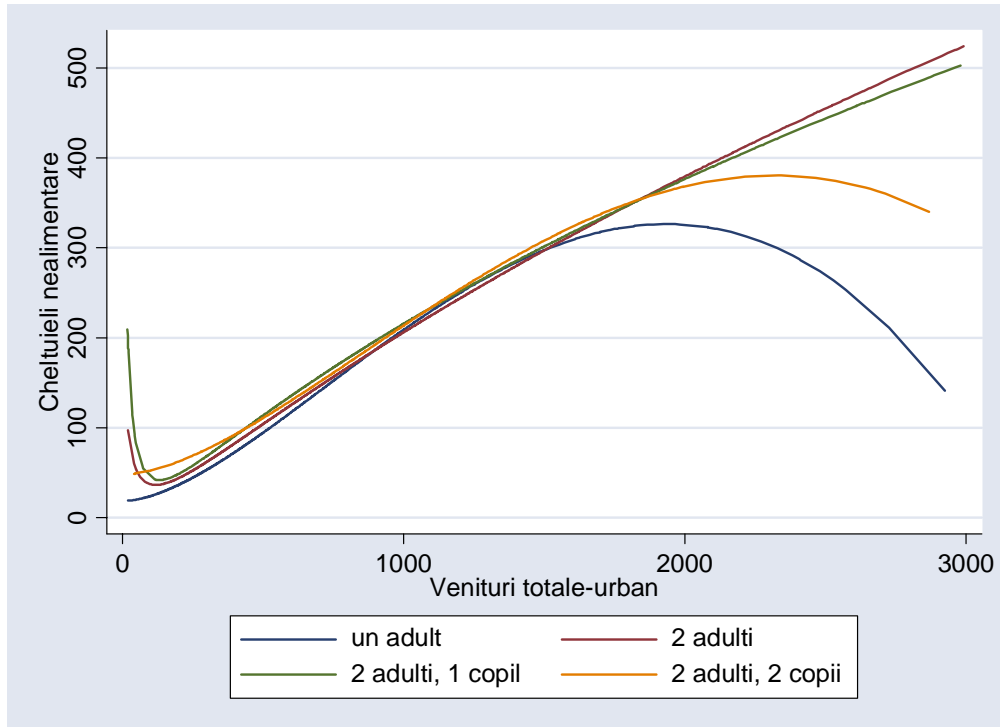
Graficul 3-7 descrie dependența cheltuielilor nealimentare de venituri. O primă observație este forma aproape liniară a dependenței la venituri mici și medii. De asemenea, **nu se mai observă diferențe foarte importante între consumul gospodăriilor urbane față de cele rurale**. La venituri mici cheltuielile nealimentare sunt chiar mai mari pentru gospodăriile rurale, dar la venituri peste 2000 RON situația se inversează.



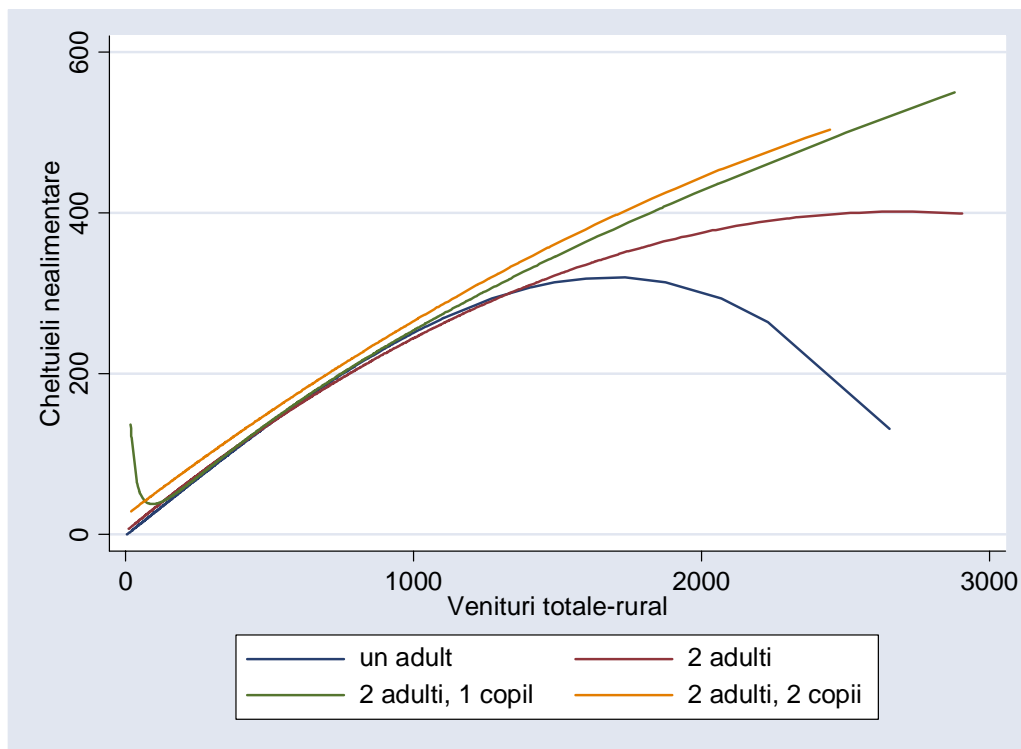
Graficul 3-7 Relația cheltuieli nealimentare – venituri totale pentru total populație, populația urbană și populația rurală

Dependența cheltuieli nealimentare – venituri pentru diverse tipuri de familii din zona urbană este similară pentru venituri mici și medii, independent de componența familiei, graficele practic se suprapun unul peste altul. În zona de variație comună se poate observa o dependență liniară a cheltuielilor de venituri, după care cheltuielile pentru o familie formată dintr-un adult și cele ale familiei formate din doi adulți și doi copii dau semne că suferă de efectele date de influența valorilor extreme în absența datelor suficiente.

Aceleași caracteristici ale dependențelor cheltuielilor nealimentare – venituri sunt valabile și pentru diversele tipuri de familii din zona rurală. Graficele sunt aproape suprapuse unul peste celălalt, la venituri mici și medii, dependența fiind aproape liniară pe secțiunea respectivă, după care cheltuielile familiilor formate din doi adulți dau semne de saturație în timp ce pentru doi adulți cu unu sau doi copii ele continuă să crească în continuare.



Graficul 3-8 Dependența cheltuieli nealimentare – venituri totale în funcție de componența familiei în mediul urban



Graficul 3-9 Dependența cheltuieli nealimentare – venituri totale în funcție de componența familiei în mediul rural.

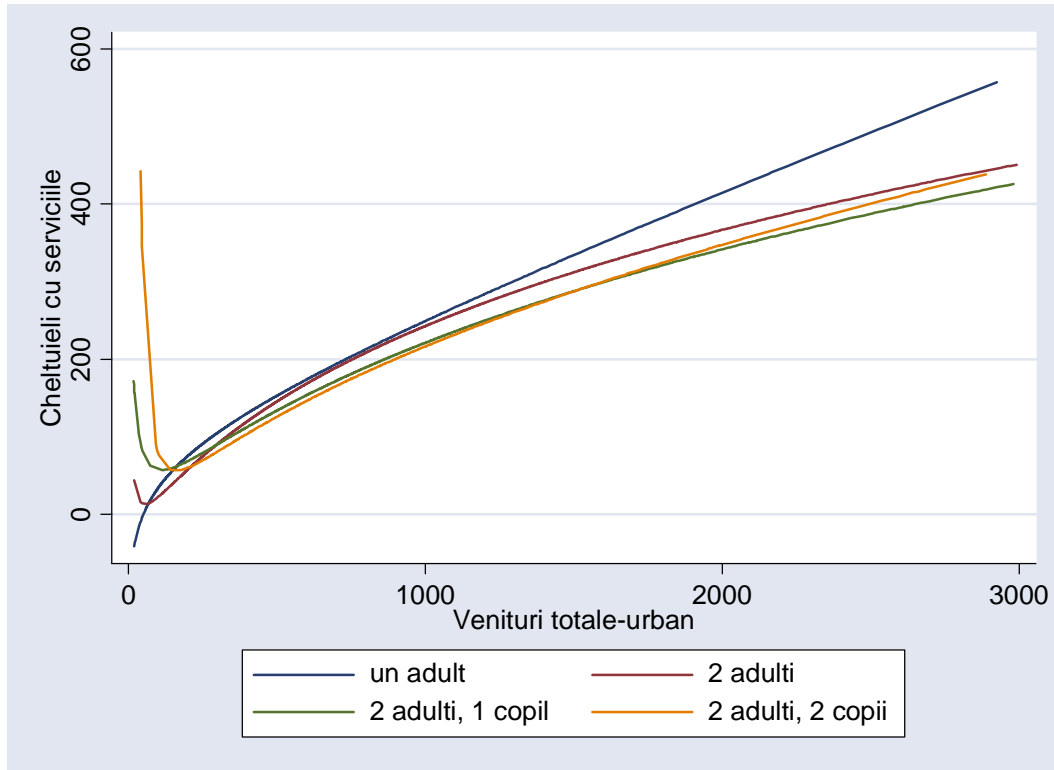


Graficul 3-10 Relația dintre cheltuielile cu serviciile – veniturile totale pentru total populație, populația urbană și populația rurală

O ultimă dependența pe care o vom analiza este relația dintre cheltuielile cu serviciile și veniturile totale. Este interesant de observat **forma dependenței**, care în cazul **consumatorilor rurali este linară, iar în cazul consumatorilor urbani are o formă logaritmică**. Cheltuielile cu serviciile sunt mai mari în mediu urban, în mare parte datorită existenței unui pachet mai divers de servicii în mediu urban. Nu discutăm aici numai de utilități, despre care am mai discutat, dar și de servicii către populație pentru petrecerea timpului liber și nu numai.

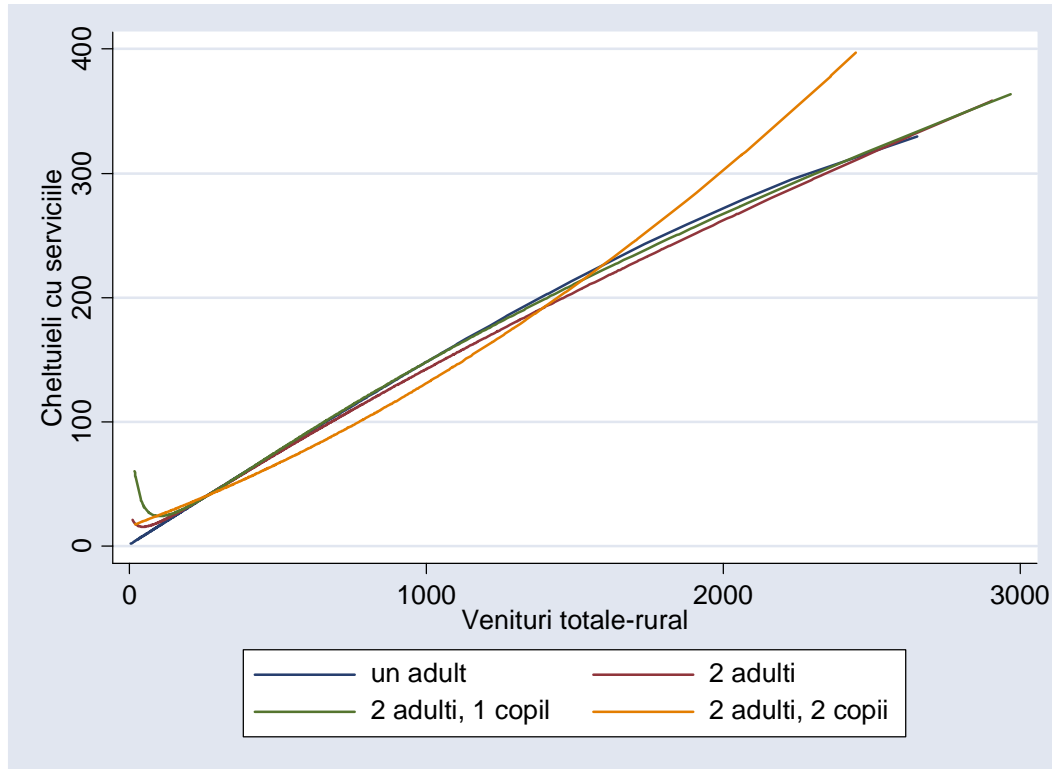
Graficul 3-11 prezintă relația dintre cheltuielile cu serviciile față de venituri pentru tipuri de gospodării, și se poate observa o **variație foarte apropiată a cheltuielilor cu serviciile pentru toate gospodăriile care conțin doi adulți**. S-ar părea că aceste cheltuieli nu depind de numărul de copii, și reprezintă mai mult niște cheltuieli ale adulților. Bineînțeles că în servicii intră și utilitățile, dar există o corelație puternică între venituri și consumul de utilități, și mai puțin puternică între consumul de utilități și numărul de persoane din gospodărie. Spre exemplu, considerați cheltuielile cu încălzirea, acestea depind de volumul de încălzit, de tipul de locuință; de tipul de încălzire folosit, temperatura exterioară, etc. și mai puțin de numărul de persoane care locuiesc în camera respectivă. Veniturile, pe de altă parte, pot fi decisive în

stabilirea tipului de încălzire folosit^{***}, și a numărului de camere încălzite. În cazul unor venituri mici, o familie poate decide să încălzească o cameră în timp ce la venituri mari aceeași familie poate decide să încălzească mai multe camere.



Graficul 3-11 Dependența cheltuieli cu serviciile – venituri totale în funcție de componența familiei în mediul urban

^{***} În cazul unor venituri mai mari, o gospodărie poate decide să-și instaleze o centrală de apartament în loc să se încălzească cu sobe



Graficul 3-12 Dependența cheltuielilor cu serviciile – venituri totale în funcție de componența familiei în mediul rural

Interesant este însă dependența cheltuielii cu serviciile – veniturile totale pentru gospodăriile urbane formate dintr-un singur adult. Acestea depășesc, la venituri similare, cheltuielile cu serviciile pentru gospodăriile formate din mai multe persoane. Cum am observat deja, acest tip de cheltuieli par să fie specific adulților și adulții singuri cheltuie mai mult la venituri similare, probabil datorită cheltuielilor mai mari pe care aceste persoane le au în petrecerea timpului liber. Adulții singuri au cheltuieli mai reduse atât la bunurile alimentare cât și la cele nealimentare, ce le rămân din venituri cheltuiesc pentru servicii.

În Graficul 3-12 se prezintă dependența cheltuielii cu serviciile – venituri totale pentru gospodăriile din mediul rural. În acest caz, cheltuielile cu serviciile sunt destul de apropiate independent de tipologia familiei, ele depinzând numai de venituri. Față de dependențele urbane cheltuielile adulților singuri nu mai ies în evidență, în sensul că în pasul cu restul cheltuielilor la venituri similare, însă dacă le raportăm la numărul de adulți sunt practic duble în comparație cu restul familiilor. În schimb, familiile cu doi adulți și doi copii la venituri mari au cheltuieli cu serviciile mult mai mari decât restul.

4 Concluzii

În această lucrare au fost folosite metodele neparametrice pentru a studia caracteristicile consumului gospodăriilor. Au fost construite funcții de densitate pentru venituri, dar și estimatori ai relațiilor dintre cheltuieli și venituri.

Din punctul de vedere al veniturilor, s-a constatat că veniturile urbane sunt sensibil superioare veniturilor rurale. Iar acest lucru nu a fost dat de tipologia diferită a familiilor în cele două medii, pentru că diferențele între venituri s-au menținut și atunci când s-a studiat veniturile medii în funcție de tipul familiei, și nici de nivele de salarizare sensibil diferite între cele două medii, pentru că deși distribuțiile salariilor nu sunt identice, nici nu sunt atât de diferite încât să fie responsabile de diferențele înregistrate la venituri.

Cele mai mari discrepante s-au înregistrat în pensii, distribuția pensiilor rurale este trimodală cu cel mai mare procent de populație la pensiile mici – corespunzătoare pensiilor agricultorilor. În zona urbană pensiile agricole lipsesc, distribuția având doar două vârfuri, iar cel mai mare procent de pensionari sunt la nivelul mare de pensie. Astfel discrepanțele veniturilor pot fi explicate măcar parțial de diferențele privind statutul ocupațional al persoanelor din urban și rural, în zona rurală se află mai puțini salariați și mai mulți lucrători/pensionari din agricultură.

La analiza dependențelor diverselor cheltuieli față de venituri rezultatele au fost mai variate. Cheltuielile alimentare sunt foarte diferite între rural și urban, iar la venituri similare sunt mai mari în urban, în timp ce cheltuielile nealimentare sunt mai mari în rural pentru venituri medii și mici. Cheltuielile cu serviciile în zonele rurale sunt mult mai scăzute decât în zonele urbane, iar cel puțin parțial responsabil pentru asta este lipsa se ofertă în rural.

Cheltuielile cu alimentele sunt tipul de cheltuială cel mai sensibil la tipologia familiei, în timp ce cheltuielile cu serviciile sunt aproape identice la venituri similare indiferent de caracteristicile gospodăriei, chiar și în cazul familiilor formate dintr-un singur adult în mediu urban, ceea ce practic înseamnă că un adult singur cheltuie de două ori mai mult cu serviciile decât un adult căsătorit.

Din punctul de vedere al relației dintre variabile s-a observat că există cazuri când dependența veniturilor de cheltuieli nu este liniară, mai ales în cazul cheltuielilor alimentare. Acest lucru nu este surprinzător, dar poate ridica probleme privind folosirea estimatorului celor mai mici pătrate. În cazul în care neliniaritatea dependențelor nu poate fi corectată prin introducerea unor variabile suplimentare, atunci probabil că ar trebui să se liniarizeze relația prin logaritmare, sau poate este necesară folosirea unor metode neparametrice.

Bibliografie

1. Ahmad, I.A., Lin, P. E. (1984): “*Fitting a multiple regression*”, Journal of Statistics Planning and Inference, , vol. 2, pp. 163-176.
2. Benedetti, J. K. (1977): “*On the Nonparametric Estimates of Regression Function*” *Journal of the Royal Statistical Society, Series B*, vol. B, vol 39, pp. 248-253.
3. Bierns, H. J., Pott-Bitter, H. A. (1990): “*Specification of Household Engel Curves by Nonparametric Regression*” *Econometric Reviews*, vol.9, pp.123-184.
4. Chu, C. K., Marron, J. S. (1991): “*Choosing a Kernel Regression Estimator*” *Statistical Science*, vol. 6, pp. 404-436.
5. Devroye, L. P., Gyorfi, L. (1985): *Nonparametric Density Estimation*, New York, Wiley.
6. Dodge, Y. (1986): “*Some Difficulties Involving Nonparametric Estimation of a Density Function*” *Journal of Official Statistics*, vol. 2, pp. 193-202
7. Epanechnikov (1969) *Nonparametric estimation of a multidimensional probability density function* *Theory Probability Appl.*, vol. 14, pp. 153-158.
8. Greene, W.H. (1993): *Econometric Analysis*, New Jersey: Prentice Hall International Editions.
9. Greblick, W., Krzyak (1980): “*Asymptotic Properties of Kernel Estimates of a Regression Function*” *Journal of Statistical Planning and Inference*, vol. 4, pp. 81-90.
10. Hardle, W. (1990): *Applied Nonparametric Regression*, New Zork, Cambridge University Press.

11. Hardle, W. (1991): *Smoothing techniques with implementation in S, ser. Springer Series in Statistics. New York: Springer-Verlag*
12. Hardle, W., Oliver, L. (1994): *Applied Nonparametric Methods*, in Handbook of Econometrics vol. 4.
13. Jordan, M., Regep, M., Chilian, M. N. (2001): *"Household Consumption in the Central and East - European Countries Aspiring to Joint the EU"* Romanian Journal of Economic, vol. 1-2.
14. Maddala, G. S. (1992): *Introduction to econometrics* Maxwell MacMillan International Editions.
15. Pagan, A., Ullah, A. (1999): *Nonparametric Econometrics*, Cambridge, Cambridge University Press.
16. Pindyck, R. S., Rubinfeld, D.L. (1991): *Econometric Models and Economic Forecasts*, New York: McGraw – Hill International Edition.
17. Scott, D. W. (1992): *Multivariate Density Estimation: theory practice and visualization*, Wiley.
18. Silverman, B.W. (1986): *Density Estimation for Statistics and Data Analysis* London: Chapman and Hall.
19. Stanciu, Mariana. 2006. Metode de cercetare a modelelor de consum. București: CIDE
20. Stanciu, Mariana. 2001. Structuri moderne ale consumului european. București: Editura Genicod