

Ans 2

În cele ce urmează vom considera un alt model unapăr de rețea de neuroni în care activarea unei percepții asums este determinată de distanță dintre vectorul de intrare și un vector prototip.

O proprietate interesantă și importantă a acestor tipuri de rețele sunt rețele de funcții radiale, denumite cu acronimul RBF, este aceea că unifică un număr de concepte importante cum ar fi: aproximarea funcțiilor, regularizare, interpolare din zgomot - noisy interpolation -, estimarea densităților de probabilitate, clasații optimale și funcții de potențial - potențial funcțional.

O consecință a acestui punct de vedere unificator este că aceste proceduri de antrenare substanțial mai rapide decât cele utilizate pentru învățarea RBF. Această rezultă din interpretarea care poate fi dată inter reprezentării formate din stratul funcțional esens și care condus la o procedură de învățare în doi pași. În primul pas parametrii care guvernează funcțiile radiale (corespunzătoare percepțiilor de pe stratul asums) sunt determinate relativ rapid folosind metode nesupervizate (adică metode în care sunt utilizate doar datele de intrare nu și valorile țintă corespuizătoare). Al doilea pas al învățării presupune determinarea ponderilor de pe ultimul strat, soluție unei probleme liniare care se rezolvă rapid.

Def 9 O rețea de funcții radiale - RBF - este un triplet de forma:

$$((\text{strat}, \Gamma), F_{W,P}, \Psi)$$

unde:

- (strat, Γ) este un graf orientat cu trei nivele având noduri în mulțimea Γ , cu $\text{Card}(\Gamma) = d + M + K$, dispuse astfel:

- pe stratul de intrare, d noduri;
- pe stratul funcțional asums, M percepții având un percepție;
- pe stratul de ieșire, K noduri percepție având un percepție.

- $F_{W,P}: \mathbb{R}^d \rightarrow \mathbb{R}^K$ este funcția implementată de rețea definită de

vezi versu

prin metode de antrenare supervizate, ca

$$F_{W,P}(x) = (y_1(x), y_2(x), \dots, y_K(x))' \quad \text{cu} \quad F_{W,P}(x) = (y_1(x), y_2(x), \dots, y_K(x))'$$

$$y_k(x) = \sum_{j=1}^M w_{kj} \phi_j(x) \quad \text{funcțiile de transfer ale celor } K \text{ perceptronuri de pe stratul de ieșire}$$

cu $\{w_{kj}\}_{k=1, \overline{K}}^{j=1, \overline{M}}$ coeficienții de ponderare conexiunile celor M perceptronuri de pe stratul ascuns cu cei K perceptronuri de pe stratul de ieșire ($\{w_{kj}\}_{k=1, \overline{K}}^{j=1, \overline{M}}$ sunt considerate deplasările perceptronurilor de pe stratul de ieșire).

$$\text{iar} \quad \phi_j(u) = \frac{\exp(-u)}{\sum_{j=1}^M \exp(-u_j)} \quad j = \overline{1, M}$$

$$\phi_0(x) = 1 \quad (\forall) x \in \mathbb{R}^d$$

sunt funcțiile radiale, se funcționează ca funcții de transfer ale celor M perceptronuri de pe stratul ascuns.

Funcțiile de integrare ale perceptronurilor de pe stratul ascuns sunt de forma

$$u \equiv G_j(x) = \frac{1}{2} (x - \mu_j)' \Sigma_j^{-1} (x - \mu_j) \quad j = \overline{1, M}$$

și depind de parametri $P = \{\mu_j, \Sigma_j\}_{j=1, \overline{M}}$ unde $\{\mu_j\}_{j=1, \overline{M}} \in \mathbb{R}^d$ se numesc vectori prototip iar $\Sigma_j \in \mathcal{M}_{d \times d}(\mathbb{R})$ simetrică.

Funcțiile de integrare ale perceptronurilor de pe stratul de ieșire sunt produsul scalar în \mathbb{R}^M și depind de parametri $W = \{w_{kj}\}$

- Ψ este regula de "învățare" de către rețea și funcției $F_{W,P}$ (adică de determinare a parametrilor W și P) și care se aplică în doi pași:

- 1) printr-o metodă nesupervizată pentru parametri P ;
- 2) printr-o metodă supervizată pentru parametri W .

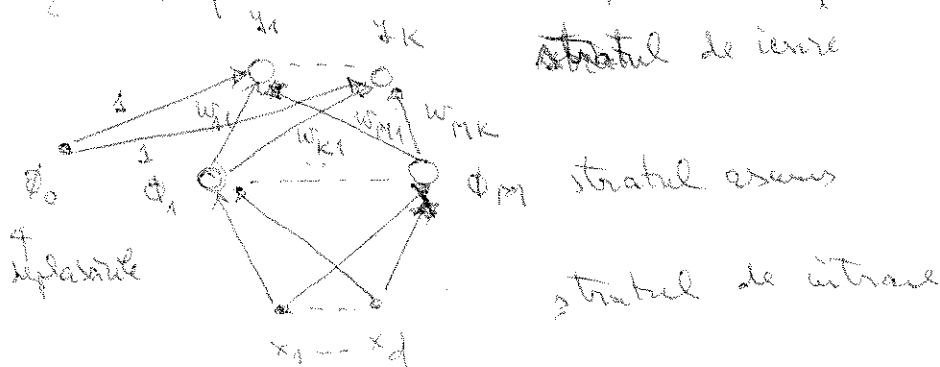


Fig. 14. Exemplu de RBF

Observații

Funcțiile radiale sunt cazuri particulare ale funcțiilor nucleare, adică funcții ^{nucleare} ~~ne~~ ^{non-negative} care ~~au pe 0 x ca~~ ^{au pe 0 x ca} asimptotă orizontală la $\pm \infty$, care își concentrează majoritatea "masei" în jurul unui punct dat; ~~Care este utilizat~~ acestea se mai numesc și funcții radiale locale.
O altă alegere de funcții radiale locale, decât cele exponențiale utilizată în RBF, este:

$$\phi(x) = (x^2 + \sigma^2)^{-\alpha}, \quad \alpha > 0.$$

În rețeaua funcțiilor radiale locale sunt utilizate, mai des, funcțiile

$$\phi(x) = x^2 \ln x \quad \text{sau} \quad \phi(x) = (x^2 + \sigma^2)^\beta \quad 0 < \beta < 1.$$

MATLAB este implementată funcția radială din def. 5

2) Spuneam, în introducere, că activarea unui perceptron este dată de distanța dintre vectorul de intrare și un vector prototip.

Funcția de integrare ~~este~~ implementată de percepții de pe stratul asums din def. 5 este inspirată din ^{densitatea} funcția de probabilitate a variabilei normale multidimensionale, motiv pentru care funcțiile

$\phi \circ G_j$ s.n. și funcții radiale gaussiene, și măsurarea distanței "reducere" (adimensionalizată datorită matricii Σ^{-1}) de la vectorul de intrare la vectorul prototip. Dacă vectorul de intrare provine dintr-o populație repartizată normal multivariat atunci vectorul prototip este ^{vectorul} ~~media~~ ^{al} variabilei aleatoare ~~care este matrice~~ ^{care este matrice} de varianță-covarianță.

Pentru reducerea numărului de parametri ce trebuie "măsurati" (d pt. μ și $\frac{d(d+1)}{2}$ pt. Σ , adică, $\frac{d(d+3)}{2}$ parametri pt. fiecare subpopulație perceptron $j=1, M$ de pe stratul asums) în practică se utilizează funcții de integrare de tipul

$$G_j(x) = \frac{1}{2\sigma_j^2} \|x - \mu_j\|^2 \quad (\text{matricea de varianță-covarianță}$$

se reduce la $\sigma_j^2 \mathbb{I}_{d \times d}$, în acest caz fiind necesară doar doi parametri de măsurat pt. fiecare din cei M perceptroni de pe stratul asums).

Cum percepții prezentate în acest curs sunt percepții fără conexiuni coeficienți $\{ \mu_j, \frac{1}{2\sigma_j^2} \}_{j=1}^M$ care trebuie să fie dați (vezi ~~observație~~ ^{observație} 3) și estimati sunt "memorați" de rețea sub forma ~~parametrelor~~ ^{coeficienților}

vezi versu

În RBF-ul cu ponderare este un plus în faptul că este ușor de înțeles

ce ponderează raexilele intrărilor cu percepționii de pe stratul respectiv deplasările acestor percepționii, astfel:

$$w_j^0(i) = \frac{1}{2\sigma_j^2} \quad \text{deplasările - bias-urile - percepțiilor } j=1, \overline{M}$$

$$w_j^k(i) = \mu_{jk}^i \quad \text{cu } k=1, d \quad \text{iar } \mu_j = (\mu_{j1}^1, \mu_{j2}^1, \dots, \mu_{jd}^1)^T$$

3) Dacă $M=N$ (adică numărul de percepționii de pe stratul ascuns este egal cu numărul de vectori din selecția de învățare) atunci conform lui Towell 1987 citat de Bishop în 1995, RBF implementează funcția & interpolatează exact punctele din mulțimea de învățare S , adică

$$F_{W,S}(x) : \mathbb{R}^d \rightarrow \mathbb{R} \quad \text{cu} \quad F_{W,S}(X_i) = Y_i \quad , i=1, \overline{M} \text{ în } (X_i, Y_i) \in S$$

Int-adevăr, dacă M scum pe $F_{W,S}$ (care simplă notat cu F când nu este nici un risc de confuzie) $F(x) = \sum_{i=1}^M w_i \phi(\|x - X_i\|)$ cu ϕ funcție radială iar $\{w_i\}$ parametrii atunci

punând condiția de interpolare exactă (adică $F(X_i) = Y_i$)

$$\Phi W = t \quad \text{unde} \quad \Phi \in \mu_{M \times M}(\mathbb{R}) \text{ cu } \phi_{ij} = \phi(\|X_i - X_j\|)$$

$$W \in \mu_{M \times 1}(\mathbb{R}) \text{ iar } t \in \mu_{M \times 1}(\mathbb{R}) \text{ cu } t_i = Y_i$$

se obține sistemul cu W vectorul necunoscutele, iar dacă Φ^{-1} există (cum arătat de Micchelli 1986 pentru o clasă largă de funcții, în ipoteza că punctele din S sunt distincte) atunci $w = \Phi^{-1} t$ ceea ce demonstrează afirmația referitoare la interpolarea exactă a unei mulțimi de puncte.

În mare număr de lucrări (Broomhead & Lowe 1988, Moody & Darken 1988 etc.) au arătat că dacă vectorii prototip nu sunt dați a priori atunci M , numărul de percepționii de pe stratul ascuns, poate fi ales mult mai mic decât N , numărul de vectori din selecția de învățare fără ca acest lucru să împiedice asupra proprietății de aproximație optimă pe care RBF-ul îl are (conform lui Girosi & Poggio, citat de Bishop 1995) (O schemă de aproximație este optimă dacă în mulțimea punctelor implementate - mulțimea punctelor generate de toate valorile parametrilor W - există o funcție care minimizează eroarea de aproximație).

• Antrenarea RBF

tratat în def 9 ca antrenarea ~~rețelei~~ ^{RBF} presupune două etape:

o etapă supervizată pentru determinarea ponderilor de conectare dintre percepționii de pe stratul 1 și cei de pe stratul 2, și o etapă nesupervizată pentru determinarea vectorilor prototip și a factorilor de scară (σ_j^2). Vom discuta mai întâi etapa de învățare supervizată, care am discutat-o în detaliu în cazul RPN cu funcții de transfer liniare.

1) ~~Invățarea supervizată~~

Să considerăm cele K percepționii de pe stratul 2

$$y_k(x) = \sum_{j=0}^M w_{kj} \cdot \phi_j(x) \quad (\text{cu convenția din def 9 referitoare la } \phi_0(x))$$

ier matricial

$$y(x) = W \phi \quad \text{unde } W = \{w_{kj}\} \text{ și } \phi = (\phi_j)$$

Se consideră riscul funcțional empiric

$$E = \frac{1}{2} \sum_{i=1}^N \sum_{k=1}^K (y_k(x_i) - \gamma_{ik})^2 \quad (\text{ținta este aici vectori, nu scalar})$$

Derivând și impunând condiția de anulare a derivatelor (ca în § "Rețele de percepționii pt. problemele clasării liniare") rezultă

$$\phi^T \phi W^T = \phi^T T \quad \text{cu } T = (\gamma_{ik}), \quad \phi = (\phi_j(x_k))$$

$$\Rightarrow W^T = \phi^+ T \quad \text{cu } \phi^+ \text{ - pseudo-inversa lui } \phi.$$

În practică sistemul de mai sus se rezolvă folosind descompunerea în valori singulare pentru a evita problemele legate de proasta condiționare a matricii ϕ .

2) Învățarea nesupervizată

Borshap 1995 enumeră mai multe metode de învățare nesupervizată: metoda celor K -medii (K -means clustering algorithm); metoda celor mai mici pătrate ortogonale (~~best~~ orthogonal least squares); metoda amestecurilor gaussiene (gaussian mixture models).

În MATLAB este implementată metoda celor mai mici pătrate ortogonale - OLS - (funcția solve) motiv pt. care o vom prezenta în cele ce urmează. În plus, metoda OLS optimizează și arhitectura rețelei în sensul că, cu titlu informativ metoda K -mediilor este implementată în SPAD pentru metodele de clasificare automată.

Stabilitate numerică în M de percepționii de pe stratul 1 asigură necesare în riscul empiric pentru funcția implementată de RBF în def 9.

Procedura de estimare a parametrilor P și a lui M este iterativă; în plus are ca "efect lateral" - side effect - și determinarea parametrilor W .

Pașul 0 (Initializare)

$j=0$ - indice ce numără percepștroni pe stratul asums
 $\bar{E}_j = \inf$ - valoare ce memorează pasul empiric al RBF găsite de algoritmul OLS
 $S_j = \emptyset$ - mulțimea vectorilor prototip (submulțime a vectorilor S din setul de antrenare)
 $\bar{S}_j = S$ sau def $\bar{S}_j = \begin{pmatrix} S_j \\ S \end{pmatrix}$
 $I_j = \{1, 2, \dots, M\}$ - mulțimea indicilor vectorilor din \bar{S}_j

Pașul 1 (Antrenarea rețelei)

Pentru $i \in I_j$

Pașul 1.1 (antrenarea nesupervizată)

~~Se construiește rețeaua cu~~

$$S_c = S_j \cup \{x_i\}; \quad \bar{S}_c = \bar{S}_j \cup \{x_i\}; \quad P = \{S_c, 1\}$$

Se construiește RBF cu mulțimea $P = P$

Pașul 1.2 (antrenarea supervizată)

Pe rețeaua de la pașul 1.1 se calculează, prin antrenare supervizată, mulțimea ponderilor W_i și riscul empiric funcțional E_i pe baza selecției de învățare S_c .

Pașul 2 (construirea rețelei)

Se calculează $i^* = \arg \min E_i$.

~~Rețeaua găsită este de~~

$$S_{j+1} = S_j \cup \{x_{i^*}\}; \quad \bar{S}_{j+1} = \begin{pmatrix} S_{j+1} \\ S \end{pmatrix}; \quad I_{j+1} = \text{indicii vectorilor din } \bar{S}_{j+1}$$

$$\bar{E}_j = E_{i^*}; \quad \text{Rețeaua } M = j+1;$$

$$P = \{S_{j+1}, 1\}; \quad W = \{W_{i^*}\}$$

Rețeaua găsită are parametri (P, W) are M percepștroni pe stratul asums, și realizează un risc empiric funcțional $= E$.

Pașul 3 (S-a terminat algoritmul?)

Dacă $[(E > \varepsilon - \text{dat}) \vee j < \text{MAX}]$ atunci $j \leftarrow j+1$ și salt la Pașul 1;
 altfel STOP

mențin numărul de percepștroni pe stratul asums, adăugând

→ vefi versu

Observație 1) Dacă $NMP = N$ atunci $E = 0$ dar pentru a obține o bună generalizare a algoritmului $NMP < N$ chiar dacă este posibil ca $E > 0$ pentru ~~soluția~~ de autotestare S .

2) Procedura OLS implică resurse de calcul deosebite - este o metodă "computer intensivă" - datorită pașilor 1.2 care se execută de card(S) $= N - j$ ori, pt. un j fixat, unde trebuie rezolvat un sistem algebric liniar de jk necunoscute prin metoda descompunerii în valori singulare.

3) Datorită faptului că $NMP \neq N$ metoda OLS nu este în general optimă, căci la pasul 1.1 parametri β sunt stabiliți independent de valoare T .

Curs 2

Satorită observatilor de mai sus (ptul 3) în special) s-a încercat implementarea unei metode de antrenare a RBF total supervizată în sensul că și parametrii P sunt estimați astfel încât să minimizeze rezidul funcției funcțional evaluat pe selecția de intrări. Urmand acest demers la sistemul de ecuații liniare

$$W' = \Phi^T T \text{ se mai adaugă sistemele de ecuații}$$

neliniare

$$\frac{\partial E}{\partial w_j} = \sum_{i=1}^N \sum_{k=1}^K \{y_k(X_i) - y_i^k\} w_{kj} \exp\left(-\frac{\|X_i - \mu_j\|^2}{2\sigma_j^2}\right) \cdot \frac{\|X_i - \mu_j\|^2}{\sigma_j^3} = 0 \quad j = \overline{1, M}$$

și

$$\frac{\partial E}{\partial \mu_j^l} = \sum_{i=1}^N \sum_{k=1}^K \{y_k(X_i) - y_i^k\} w_{kj} \exp\left(-\frac{\|X_i - \mu_j\|^2}{2\sigma_j^2}\right) \cdot \frac{(X_i^l - \mu_j^l)}{\sigma_j^2} = 0 \quad \begin{matrix} j = \overline{1, M} \\ l = \overline{1, d} \end{matrix}$$

Observații 1) Spre deosebire de OSL metoda de antrenare total supervizată necesită cunoașterea prealabilă a numărului de percepțiuni de pe stratul ascuns.

2) Satorită sistemelor neliniare și metoda total supervizată este "computationally intensive" și duce la apariția unor minime locale ale rezidului funcției. Totuși există speranța că dacă funcțiile radiale sunt bine localizate atunci pentru un vector din selecția de antrenare doar o mică parte din ~~totalul~~ funcțiilor radiale vor fi semnificativ activate și deci doar acele funcții vor fi actualizate în răspuns la respectivel vector de intrare. În această linie procedura de antrenare poate fi accelerată simțitor dacă, pentru a evita mulțimea calculilor nesemnificative, din start sunt alese acele funcții radiale ce sunt relevante pentru selecția de antrenare. S. Omohundro (1987) a publicat tehnici eficiente de stabilirea funcțiilor radiale relevante.

3) Rezolvarea sistemelor liniare de mai sus implică folosirea unor metode iterative de tipul celor descrise în cursul 1 pt. percepțiunile formale. Aceste metode necesită inițializarea lor cu valori de start. În cazul RBF aceste valori erau generate aleator; în cazul RBF folosirea uneia dintre metodele de antrenare nesupervizate ~~pentru generarea valorilor inițiale este~~ este des recomandată în literatura de specialitate.

→ vezi versu

4) Antrenarea total supervizată a RBF nu garantează că funcțiile radiale rămân localizate (funcțiile sunt localizate dacă înținem $\{x \in \mathbb{R}^d / \| \phi_j(x) \| \geq \varepsilon\}$ este relativ mică; de exemplu ~~în~~ volumul unor hipersfere de rază ≈ 1 . ~~Pe de altă parte~~ ~~Se observă că, acest~~ ~~volum~~ ~~este~~ ~~conținut~~ ~~de~~ $\{ \phi_j \}$. ~~Pe de altă parte~~ ~~data~~ $v_j \rightarrow \pm \infty$ atunci sistemele neliniare sunt satisfăcătoare dar, din observații de mai sus, funcțiile nu mai sunt localizate). Simulările numerice au confirmat și ele această "dubioasă presunție" arătând că o submulțime a percepțiilor de pe stratul esous evoluează către funcții de transfer (adică funcții radiale) cu o plajă foarte largă de răspuns (Moody & Darken 1989). Adică, unul din avantajele de bază ale RBF, anume antrenarea rapidă ^{în doi pași} și interpretabilitatea stratului esous, se pierde dacă se adoptă antrenarea total supervizată.

5) MATLAB-ul raportează următorii timpi ^{de antrenare} ~~pe stratul esous~~

| Funcție | Metoda | Timpi cu xkls | Nr. ^{de iterații} max | Operații |
|----------|---|--------------------------|---|-----------|
| TRAINBP | Propagare înapoi | 259.1 | 4123 | 21603435 |
| TRAINBPX | Propagare înapoi rapidă | 42.4 | 570 | 2.971.696 |
| TRAINLM | — — — pe bază algoritmului Levenberg-Marquardt | 3-3 | 5 | 315.769 |
| SOLVERB | Funcții radiale | 1.9 | 5 | 38.305 |

Cm 2

- Comparabile între RBF și RPM

RBF și RPM ~~sunt~~ au valori similare în sensul că pentru dispozitive tehnice de aproximare a punctelor neliniare multidimensionale, în ambele cazuri punctele sunt exprimate sub formă de compuneri de aplicații scalare parametrizate. Structura particulară a celor două tipuri de rețele este însă foarte diferită. Unele din descrierile importante ale ~~rețelor~~ RBF comparativ cu RPM sunt:

- (1) Funcția de ieșire a unui perceptron asumat într-o RPM depinde de suma ponderată a intrărilor transformată de funcții de transfer monotone. Această funcție este ^{paralelă,} constantă pe suprafețe formate din ~~plan~~ hiperplane $(d-1)$ -dimensionale în spațiul d -dimensional. În contrast, perceptronul asumat într-o RBF folosește ca funcție de integrare distanța de la un vector prototip, iar ca funcție de transfer, o transformare ce este, de obicei, o funcție localizată. În ambele cazuri, această funcție este constantă pe hiperesfere $(d-1)$ -dimensionale concentrice (sau mai general pe hiperelipsoidi $(d-1)$ -dimensionali).
- (2) O RPM formează o reprezentare distribuită în spațiul valorilor funcțiilor de transfer ~~pentru~~ ^{ale} perceptronilor asumați deorece, pentru un vector de intrare dat, mai multe perceptronuri asumați contribuie la evaluarea ~~la~~ ^{la} ieșirii. În timpul procesului de învățare, valorile funcțiilor de transfer ale perceptronilor asumați sunt astfel ajustate încât, contribuția liniară ~~la~~ ^{la} ponderile ultimului strat să genereze ieșiri corecte pentru un interval de posibile valori de intrare. Punctele de minim local cât și erorile de trunchiere inerente unor metode numerice puternice neliniare fac ca procesul de convergență să fie lent dintr-un caz în altul, în cazul unor metode de optimizare avansate. În contrast, o RBF cu funcții radiale localizate formează o reprezentare locală în raport cu un vector de intrare dat, în mod obișnuit doar foarte puține funcții de transfer sunt semnificativ nenule.
- (3) RPM are în cele mai multe cazuri mai multe straturi, o topologie complexă și o ^{utilitate} varietate de funcții de activare. RBF are o topologie simplă formată din două straturi cu doar două tipuri de funcții de activare specifice.
- (reverso)

(2) ORPM furnizează o reprezentare distribuită în spațiul valorilor punctilor de transfer ~~pentru~~^{pentru} percepții ascunse deosebite, pentru un vector de intrare dat, mai multe percepții ascunse contribuie la evaluarea ~~pentru~~^{pentru} ieșiri. În timpul procesului de învățare, valorile punctilor de transfer ale percepțiilor ascunse sunt astfel ajustate încât, combinația liniară a ponderilor ultimului strat să genereze ieșiri corecte pentru un interval de posibile valori de intrare. Punctele de minim local cât și erorile de training sunt utilizate unor metode numerice puternice, neliniare, iar procesul de convergență să fie lent din cauza neliniarității procesului de optimizare avansată. În contrast, ORB cu puncte radiale localizate furnizează o reprezentare locală în raport cu un vector de intrare dat, în mod obisnuit doar foarte puține puncte de training sunt semnificativ neune.

(3) RPL are în ele mai multe căutări, mai multe structuri, o topologie complicată și o ^{utilitate} varietate de funcții de activare. RBF are o topologie simplă formată din două structuri cu doar două tipuri de funcții de activare specifice.

(4) Toți parametrii unei RPM sunt, în mod obișnuit, determinați în același timp într-o procedură de antrenare globală bazată pe o metodă de învățare supervizată. O RBT, în contrast, este în mod obișnuit antrenată în două etape, cu parametrii funcțiilor de activare determinați printr-o metodă de învățare nesupervizată ce folosește doar vectorii de intrare iar ~~parametrii funcțiilor de activare~~ ^{ponderile unității de activare} fiind ulterior determinați printr-o metodă de învățare supervizată liniară rapidă și exactă. ~~ce folosește doar~~