

# CLASIFICAREA FORMELOR PLANE UTILIZÂND DISTANȚE SINTACTICE

dr. ing. Ovidiu Grigore  
dr. ing. Octavian Grigore

Universitatea Politehnica București

**Rezumat:** În lucrare sunt prezentate metode de clasificare sintactice, bazate pe compararea cu modele ("matching"). În general, formele plane (caractere, imagini de obiecte etc.) au o alcătuire complexă, care, într-o abordare globală de tip statistic, conduc la discriminări greoaie, cu rezultate destul de slabe. Din acest motiv, s-a utilizat o descriere sintactică a formelor, care pune în evidență modul lor de alcătuire atât din punct de vedere al componentelor, cât și al relațiilor dintre acestea. Pentru clasificarea structurilor clasice rezultate, s-au folosit distanțe corespunzătoare. S-a introdus o nouă distanță sintactică *normată*, care a condus la îmbunătățirea rezultatelor de clasificare obținute prin utilizarea distanței clasice Wagner-Fischer. În final, sunt prezentate rezultatele obținute într-o aplicație de clasificare a rezultatelor scrise de mână.

**Cuvinte cheie:** recunoaștere sintactică a formelor, distanțe sintactice, clasificarea formelor plane.

## 1. Calcularea distanțelor între șiruri de primitive

Printre cei mai simpli algoritmi de clasificare a formelor ce pot fi descriși prin intermediul șirurilor de primitive sunt cei bazați pe compararea (potrivirea) șirurilor atașate formelor. Problema este simplă în cazul în care structura formei necunoscute se potrivește perfect cu una din structurile formelor lotului de învățare, însă, dacă acest lucru nu se realizează, atunci clasificarea trebuie făcută pe baza unei distanțe minime dintre șirul formei necunoscute și șirurile de primitive corespunzătoare formelor lotului de învățare.

În general vorbind, a potrivi un șir de caractere A cu altul B, înseamnă să transformăm simbolurile cuvintelor A în literele lui B cu un cost minim al operațiilor posibile de realizat. În general, se presupune că există 3 tipuri de operații ce se pot efectua pentru a transforma un cuvânt în altul:

1. *Înlocuirea caracterelor*, ceea ce presupune schimbarea unui simbol  $a$  cu un altul  $b$ , operație ce se poate nota  $a \rightarrow b$ ;
2. *Inserarea caracterelor*, introducerea unui simbol  $a$  în șir, echivalentă cu operația înlocuire:  $\Lambda \rightarrow a$  în care  $\Lambda$  este simbolul nul;
3. *Eliminarea caracterelor*, adică ștergerea unei litere "b" din șir, operație echivalentă cu înlocuirea simbolului "b" cu cel vid:  $b \rightarrow \Lambda$ .

Fiecare dintre operațiile de modificare a simbolurilor unui șir presupune că este cotate cu un anumit cost:

- $C(a, b)$  pentru o înlocuire  $a \rightarrow b$ ;
- $C_1(a) = C(\Lambda, a)$  pentru o inserare  $\Lambda \rightarrow a$ ;
- $C_2(b) = C(b, \Lambda)$  pentru o eliminare  $b \rightarrow \Lambda$ .

Vom admite că, pentru orice trei simboluri a,b,d din alfabetul V sunt adevărate:

- a)  $C(a, b) > 0$ ;
- b)  $C(a, b) = 0$ , dacă și numai dacă  $a=b$ ;
- c)  $C(a, b) = C(b, a)$ ;
- d)  $C(a, d) < C(a, b) + C(b, d)$ .

Dacă pentru transformarea unui șir de caractere  $A = a_1 a_2 a_3 \dots a_n$  în alt șir  $B = b_1 b_2 b_3 \dots b_m$ , este necesară succesiunea de operații  $s_1, s_2, \dots, s_r$  în care  $S = \{s_i \mid i = 1, r\}$  sunt transformări de tipul celor definite anterior, atunci costul modificării lui A în B prin secvența de transformări  $S = [s_1, s_2, \dots, s_r]$  este:

$$C_t(S) = \sum_{i=1}^r C(S_i)$$

în care  $C(S_i)$  este costul transformării  $S_i = a_i \rightarrow b_j$ , cu  $a_i$  sau  $b_j$  eventual egale cu cuvântul vid  $\Lambda$ .

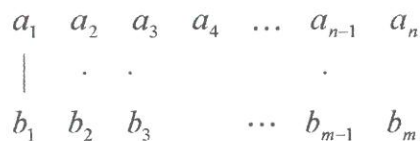
Se definește distanța  $d(A, B)$  între șirurile A și B ca minimul costurilor tuturor secvențelor de transformări care se pot face pentru a-l schimba pe A în B:

$D(A, B) = \min \{C_t(S) | S = s_1, \dots, s_r \text{ este o secven\c{t}e de transformare a lui A \u00een B}\}.$

Dac\u0103  $A = a_1 a_2 a_3 \dots a_n$  \u015fi  $B = b_1 b_2 b_3 \dots b_m$ , atunci se definesc transform\u0103rile speciale ale lui A \u00een B, secven\c{t}ele de modific\u0103ri de forma:

$$a_1 \rightarrow b_1, \quad a_2 \rightarrow \Lambda, \quad a_3 \rightarrow b_2, \dots, a_n \rightarrow b_{m-1}, \quad \Lambda \rightarrow b_m,$$

care sunt reprezentate \u00een mod sugestiv \u00een diagrama din figura 1, \u00een care barele indic\u0103 substitu\c{t}iile,



**Figura 1. Secven\c{t}a de opera\c{t}ii necesare transform\u0103rii cuv\u00e2ntului A \u00een B**

iar literele izolate sunt inserate sau eliminate. \u00c0ntr-o transformare special\u0103, folosind reprezentarea schematic\u0103 din figura 1, barele de substitu\c{t}ire nu trebuie s\u0103 se intersecteze \u015fi niciodat\u0103 dou\u0103 sau mai multe bare nu pornesc \u015fi nu ajung la acela\u015fi caracter.

Se poate demonstra [18] c\u0103 distan\c{t}a  $d(A, B)$  dintre dou\u0103 \u015firuri A \u015fi B este egal\u0103 cu minimumul costurilor transform\u0103rilor speciale \u00entre A \u015fi B. Aceast\u0103 afirma\c{t}ie este de un real folos, deoarece num\u0103rul transform\u0103rilor oarecare \u00entre componentele celor dou\u0103 \u015firuri A \u015fi B este mult mai mare fa\c{t}a de num\u0103rul total al transform\u0103rilor speciale \u00entre A \u015fi B, astfel c\u0103 se restr\u00e2ng destul de mult posibilit\u0103\c{t}ile de c\u00e2utare a costului minim \u00een scopul determin\u0103rii distan\c{t}ei  $d(A, B)$ .

Fie  $A = a_1 a_2 a_3 \dots a_n$  \u015fi  $B = b_1 b_2 b_3 \dots b_m$  dou\u0103 cuvinte din dic\c{t}ionarul  $V^*$ . Vom nota prefixele acestor cuvinte cu:

$$A(i) = a_1 a_2 a_3 \dots a_i$$

$$B(j) = b_1 b_2 b_3 \dots b_j$$

Principiul de baz\u0103 al algoritmului de calcul a distan\c{t}ei \u00entre dou\u0103 cuvinte, care va fi prezentat \u00een continuare, se bazeaz\u0103 pe evaluarea distan\c{t}ei \u00entre prefixele A (i) \u015fi B (j) pe baza distan\c{t}elor \u00entre prefixele cu lungime mai mic\u0103 cu o unitate, astfel:

$$d(A(i), B(j)) = D(i, j) = \min \begin{cases} D(i-1, j-1) + C(a_i, b_j), \\ D(i-1, j) + C_1(a_i), \\ D(i, j-1) + C_2(b_j) \end{cases}$$

### 1.1. Algoritm Wagner – Fischer [20]

P1.  $D(0, 0) = 0$

P2. Pentru i de la 1 la n se calculeaz\u0103:  $D(i, 0) = D(i-1, 0) + C_1(a_i)$

P3. Pentru j de la 1 la m se calculeaz\u0103:  $D(0, j) = D(0, j-1) + C_2(b_j)$

P4. Pentru i = 1, n \u015fi j = 1, m se determin\u0103:

$$d_1 = D(i-1, j-1) + C(a_i, b_j)$$

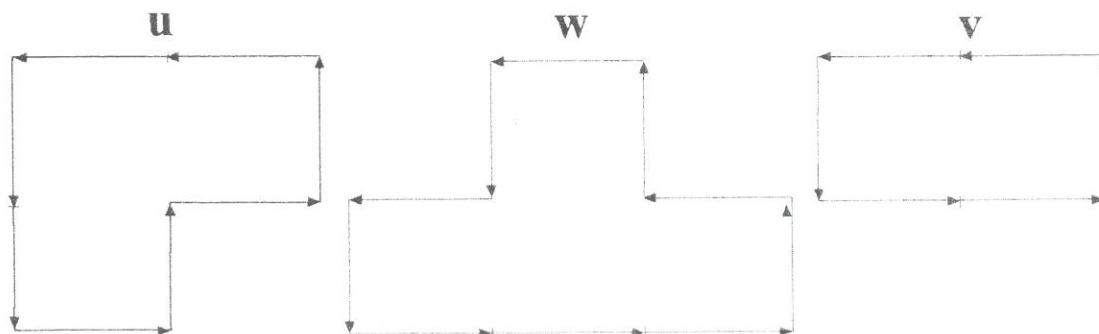
$$d_2 = D(i-1, j) + C_1(a_i)$$

$$d_3 = D(i, j-1) + C_2(b_j)$$

$$D(i, j) = \min \{d_1, d_2, d_3\}$$

P5. Distan\c{t}a \u00entre \u015firurile A \u015fi B este  $d(A, B) = D(n, m)$ .

Pentru a ar\u0103ta posibilitatea de aplicare a algoritmului Wagner – Fischer \u00een clasificarea formelor vom apela la urm\u0103toarea aplica\c{t}ie de recunoa\u015ftere a conturilor. Vom presupune



**Figura 2. Reprezentarea sintactică a conturilor unor forme plane**

ca formele  $u$  și  $w$  prezentate în figura 2 sunt prototipurile a două clase  $\omega_1$  și  $\omega_2$ , iar forma  $v$  din aceeași figură trebuie clasificată într-una din cele două clase. Pentru a lua decizia de clasificare pe baza distanței între șiruri de simboluri, se face descompunerea celor trei forme în raport cu primitivele prezentate în figura 3, împreună cu simbolurile atașate:  $a, b, c, d$ .



**Figura 3. Primitivele în raport cu care se descompun conturile din figura 2**

Astfel, se obțin următoarele șiruri de primitive atașate celor trei forme studiate:

$$u = a a b b c d c d$$

$$w = a b a b c c c d a d$$

$$v = a a b c c d$$

În raport de diferența dintre direcțiile celor patru primitive, vom atribui următoarele costuri pentru operațiile de schimbare a simbolurilor:

$$C(e \rightarrow \Lambda) = 1$$

$$C(\Lambda \rightarrow e) = 1$$

$$C(e \rightarrow e) = 0$$

pentru orice  $e \in \{a, b, c, d\}$  și:

$$C(a \rightarrow b) = 1$$

$$C(a \rightarrow c) = 2$$

$$C(a \rightarrow d) = 1$$

$$C(b \rightarrow c) = 1$$

$$C(b \rightarrow d) = 2$$

$$C(b \rightarrow a) = 1$$

$$C(c \rightarrow d) = 1$$

$$C(c \rightarrow a) = 2$$

$$C(c \rightarrow b) = 1$$

$$C(d \rightarrow a) = 1$$

$$C(d \rightarrow b) = 2$$

$$C(b \rightarrow c) = 1$$

Rezultatul final al aplicării algoritmului Wagner – Fischer pentru calcularea distanței dintre cele trei forme luate două câte două este prezentat în tabelul 1.

**Tabelul 1. Distanțele Wagner – Fischer între conturile din figura 1**

| $d_{W-F}$ | u | w | v |
|-----------|---|---|---|
| u         | 0 | 5 | 2 |
| w         | 5 | 0 | 4 |
| v         | 2 | 4 | 0 |



Se observă că distanța minimă se obține între  $u$  și  $v$  și, prin urmare, se ia decizia ca  $v$  să aparțină clasei din care face parte prototipul  $u$ .

De asemenea, este interesant de arătat care sunt transformările speciale, care conduc la un cost minim. Astfel:

- pentru transformarea  $u \rightarrow v$  avem:

$$\begin{array}{cccccccc}
 u: & a & a & b & b & c & d & c & d \\
 & & | & | & | & & | & | & | \\
 v: & a & a & b & & c & & c & d
 \end{array}$$

- pentru modificarea  $w \rightarrow v$ :

$$\begin{array}{cccccccccc}
 w: & a & b & a & b & c & c & c & d & a & d \\
 & & | & & | & | & | & | & & & | \\
 v: & a & & a & b & c & c & & d & &
 \end{array}$$

- între prototipurile celor 2 clase,  $u \rightarrow w$ :

$$\begin{array}{cccccccccc}
 u: & a & & a & b & b & c & & d & c & d \\
 & & | & & | & | & & | & & | & | \\
 w: & a & b & a & b & c & c & c & d & a & d
 \end{array}$$

În reprezentările de mai sus ale transformărilor speciale, caracterele izolate sunt inserate sau eliminate, liniile subțiri reprezintă potriviri primitive, iar liniile îngroșate reprezintă înlocuiri de caractere.

## 1.2. Algoritmul Wagner – Fischer normat

Distanța dintre 2 cuvinte este puternic influențată de lungimea respectivelor șiruri de caractere. Pentru a proba această afirmație, trebuie să observăm că nu este echivalentă o anumită distanță între cuvinte de lungime mică cu aceeași distanță între cuvinte de lungime mare. Pentru a explica mai bine, să presupunem o distanță egală cu 2, care, pentru cuvinte de 3 caractere, conduce la o eroare de 67% între cuvinte, iar pentru cuvinte de 10 litere procentul de nepotrivire a simbolurilor din cele două cuvinte scade la 20%.

Din acest motiv, se propune o metodă de normare a distanței Wagner – Fischer între cuvinte, astfel încât aceasta să devină independentă de lungimea cuvintelor.

Fie  $u$  și  $v$  două șiruri de caractere de lungimi  $n_u$  și respectiv  $n_v$ .

Notăm  $d_{W-F}(u, v)$  distanța dintre cuvintele  $u$  și  $v$ , calculată cu algoritmul Wagner – Fischer. Distanța normalată se calculează cu formula:

$$d_{normală}(u, v) = \frac{d_{W-F}(u, v)}{\max\{n_u, n_v\}}$$

Această distanță are valoarea maximă egală cu 1 când toate simbolurile celor două cuvinte sunt diferite, indiferent de lungimea lor.

Valorile normale ale distanțelor între formele  $u$ ,  $w$  și  $v$  din figura 1, folosite anterior, sunt prezentate în tabelul 3.

Tabelul 3. Distanțele normale între contururile din figura 1

| $d_{normală}$ | $u$  | $w$ | $v$  |
|---------------|------|-----|------|
| $u$           | 0    | 0.5 | 0.25 |
| $w$           | 0.5  | 0   | 0.4  |
| $v$           | 0.25 | 0.4 | 0    |

## 2. Aplicație de recunoaștere a caracterelor alfa – numerice, scrise de mână

Pentru o comparare mai bună a celor două metode, s-a realizat o aplicație de recunoaștere a caracterelor alfa-numerice, scrise de mână, descrisă în [21].

Astfel, s-au folosit caractere din 3 clase: A, C, E. Descrierea sintactică s-a realizat mai întâi sub formă de arbore. Pentru aceasta, s-a realizat, mai întâi, o prelucrare de subțiere după care aceasta a fost aproximată poligonal, folosind un anumit prag de eroare. Fiecare segment al aproximării poligonale a fost codat ca primitivă în funcție de domeniul său unghiular, după cum urmează:

- Primitiva  $a$  pentru  $\alpha \in [-7\pi/8, -5\pi/8)$
- Primitiva  $b$  pentru  $\alpha \in [-5\pi/8, -3\pi/8)$
- Primitiva  $c$  pentru  $\alpha \in [-3\pi/8, -\pi/8)$
- Primitiva  $d$  pentru  $\alpha \in [-\pi/8, \pi/8)$
- Primitiva  $e$  pentru  $\alpha \in [\pi/8, 3\pi/8)$
- Primitiva  $f$  pentru  $\alpha \in [3\pi/8, 5\pi/8)$
- Primitiva  $g$  pentru  $\alpha \in [5\pi/8, 7\pi/8)$
- Primitiva  $h$  pentru  $\alpha \in [7\pi/8, 9\pi/8)$

Reprezentarea de tip arbore este transformată în reprezentare sub forma de șir de caractere. Așa cum s-a arătat în [21], metoda de transformare arbore-șir, care conduce la rezultatele cele mai bune este cea a delimitatorilor.

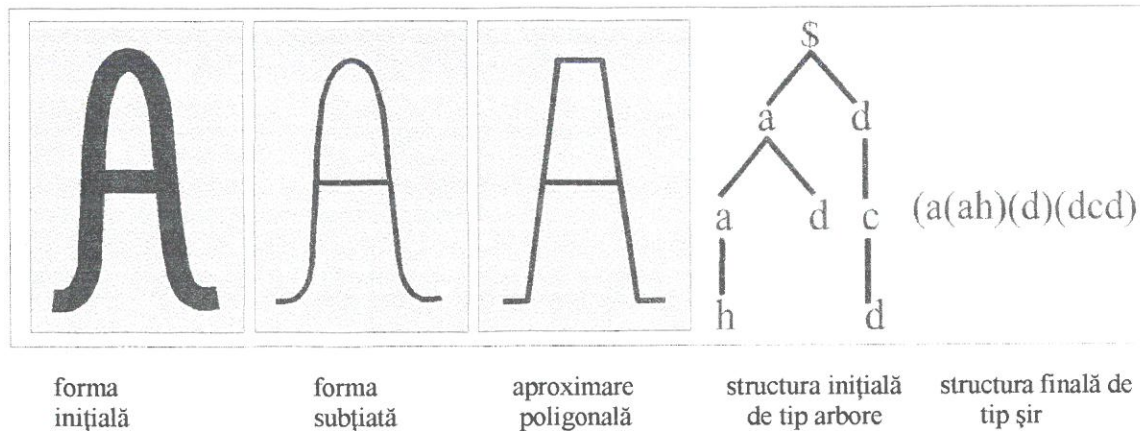


Figura 4. Lanțul de procesare pentru transformarea caracterelor de intrare în structura echivalentă de tip șir de primitive

Întreg lanțul de procesări necesare transformării caracterului de intrare în structura sa echivalentă de tip șir de primitive poate fi urmărit în exemplul dat în figura 4.

Rezultatele clasificării sunt prezentate în tabelul 4. Pentru comparare, s-au utilizat 4 parametri care pun în evidență capacitatea metodei de a mări sensibilitatea între clase, precum și calitatea discriminării formelor aflate în zona de graniță dintre clase.

Tabel 4. Rezultatul clasificării caracterelor

| Parametrii comparați  | Distanța Wagner - Fischer | Distanța normată |
|---|---------------------------|------------------|
| Distanța medie intraclasă                                       | 2.85                      | 0.21             |
| Distanța medie interclase                                       | 6.6                       | 0.5              |
| <u>distanța medie interclase</u><br>distanța medie intraclasă   | 2.31                      | 2.38             |
| <u>distanța maximă intraclasă</u><br>distanța minimă interclase | 0.66                      | 0.61             |



### 3. Concluzii

Se observă că, pentru metoda normată, se obține o valoare mai mare a raportului (distanța medie interclase/distanța medie interclase), ceea ce presupune o mărire a sensibilității între clase, respectiv o concentrare a formelor în jurul mediilor și o mărire a distanței între clase. De asemenea, metoda normată este mai bună și din punct de vedere al raportului (distanța maximă intraclase/distanța minimă interclase), ceea ce presupune ca algoritmul normat realizează și o clasificare mai bună a formelor aflate în vecinătatea granițelor claselor.

### Bibliografie

1. **ANDREWS, H.C.:** Introduction to Mathematical Techniques in Pattern Recognition, Wiley, New York, 1972.
2. **ATANASIU, I., A. MATEESCU:** Limbaje formale, București, 1990.
3. **BUNKE, H.:** Attributed Programmed Graph Grammars and their Application to Schematic Diagram Interpretation to Texture Analysis. În: IEEE Trans. Pattern Analysis Mach. Intell., Vol. PAMI – 4, No. 6, November, 1982, pp. 1103-1114.
4. **ESHERA, M. A., K.S. FU:** A Graph Distance Measure for Image Analysis. În: IEEE Trans. Syst. Man, Cybernetics, Vol. SMC – 13, No. 3, March, 1994, pp.398-408.
5. **EVANS, T.G.:** Grammatical Inference Techniques in Pattern Analysis. În: Software Engineering, Vol. 2, J.T.Tou (Ed.), Academic Press, New York, 1971.
6. **Filipski, A. J.:** a LEAST Mean – Squared Error Approach to Syntactic Classification. În: IEEE Trans. Pattern Mach. Intell., Vol. PAMI – 2No. 3, March 1980, pp. 252-255.
7. **GRIGORE, O.:** Syntactical Clustering Using the Guided Random Search Method, special issue of Buletinul Științific al Univerității “Politehnica” Timișoara and Tran. Automatic Control and Computer Science and Engineering, Nov. 1996, pp. 174-182.
8. **GRIGORE, O.:** Syntactical Clustering Using Genetic Algorithms. În: Real World Applications on Intelligence Technologies, Editura Academiei, 1996.
9. **GRIGORE, O.:** Syntactical Self – Organization Map. În: Computational Intelligence, Theory and Applications, B. Reusch (ED.), Springer Verlag, Dortmund, 1997, pp. 101-109.
10. **GRIGORE, O.:** The Cumulative Inference Algorithm for Regular Grammars. În: Proc. of the 10<sup>th</sup> Scandinavian Conference on Image Analysis, June 1997, Finland, pp. 191-199.
11. **GONZALES, R. C., P. WINTZ:** Digital Image Processing, Addison – Wesley, MA, 1996.
12. **LAM, L., S. W. LEE, S. Y. SUEN:** Thinning Methodologies – A Comprehensive Survez. În: IEEE Trans. Pattern Anal. Mach. Intell, Vol. PAMI – 14, No. 9, September 1992, pp. 869-885.
13. **LEMONE, K. A.:** Similarity Measures Between Extended to Sets of Strings. În: IEEE Trans. Pattern Anal. Mach. Intell, Vol. PAMI – 4, No. 3, May 1982, pp. 345-347.
14. **LU, S. Y., K. S. FU:** Error – Correcting Tree Automata for Syntactic Pattern Recognition and Image Processing, RPI; Troy, NY 1997, pp. 115-127.
15. **LU, S. Y., K. S. FU:** A Sentence – to –Sentence Clustering Procedure for Pattern Analysis. În: IEEE Trans. Syst. Man Cybern., Vol. SMC –8, No. 5, May 1994, pp. 219-227.
16. **Y. LU:** A Tree – to – Tree Distance and Its Applications to Cluster Analysis. În: IEEE Trans. Pattern Anal. Mach. Intell., Vol. PAMI – 1, No. 2, April 1979.
17. **NEAGOE, V. E., O. STĂNĂȘILĂ:** Teoria recunoașterii formelor, Editura Academiei Române, Bucharest 1992.
18. **PAVLIDIS, T.:** Structural Pattern Recognition, Springer Verla, Berlin, 1977.
19. **WAGNER, R. A., M. J. FISCHER:** The String to String Correction Problem. În: Journal Ass. Comput. Mach., Vol. 21, No. 1, January 1974, pp. 42-51.
20. **GRIGORE, OV., O. GRIGORE:** Metode de clasificare sintactică a formelor plane descrise prin structuri de tip arbore. În: Revista Română de Informatică și Automatică, vol. 10, nr. 2, 2000, pp. 23-32.