

INTELIGENȚĂ ARTIFICIALĂ

**Rezolvarea problemei de optimizare a unei mașini cu vectori
suport cu ajutorul unui algoritm evolutiv**

Clasificarea mașinilor

Autor: APETRIA George-Laurențiu

Cuprins

SECȚIUNEA 1: Descrierea problemei considerate	1
SECȚIUNEA 2: Aspecte teoretice privind algoritmul.....	2
SECȚIUNEA 3: Modalitatea de rezolvare.....	4
SECȚIUNEA 4: Listarea părților semnificative din codul sursă însoțite de explicații și comentarii	5
SECȚIUNEA 5: Rezultatele obținute prin rularea programului în diverse situații, capturi de ecran și comentarii asupra rezultatelor obținute.....	13
SECȚIUNEA 6: Concluzii	14
SECȚIUNEA 7: Bibliografie	15

SECȚIUNEA 1: Descrierea problemei considerate

Problematica de bază a proiectului este reprezentată de clasificarea unui set de mașini care se împarte în 4 categorii: unacceptable, acceptable, good și very good, fiecare categorie având un anumit număr de instanțe în setul de antrenare, respectiv testare. Mai mult decât atât, fiecare instanță conține un număr de 6 atribute prin intermediul cărora se poate face clasificarea (buying, maint, doors, persons, lug_boot, safety). Soluția propusă pentru a rezolva problema menționată mai sus este utilizarea unui algoritm evolutiv cu ajutorul căruia să se determine parametrii SVM-ului, necesari ulterior în procesul de clasificare. După determinarea ecuației caracteristice SVM-ului (parametrii α , bias, W), ecuație care separă planul în două regiuni, poate începe procesul de clasificare.

SECȚIUNEA 2: Aspecte teoretice privind algoritmul

Din punct de vedere tehnic, trebuie rezolvată problema duală a SVM, care presupune o maximizare: $W(\alpha)$. Maximizarea se realizează cu ajutorul algoritmului evolutiv. SVM este la bază un perceptron cu un singur strat, deci va trebui să determinăm dreapta care maximizează marginea ce desparte cele două clase de obiecte. Datele (obiectele care sunt intersectate de margine) se numesc vectori suport. Sunt vectori deoarece sunt puncte într-un spațiu n dimensional – în acest caz 2 dimensiuni și sunt suport pentru că sprijină marginea. Scopul modelului este: detectarea valorilor suport care maximizează marginea. Deci există o singură margine de lățime maximă. Pentru algoritmul evolutiv, se utilizează varinata standard prezentată în curs.

Problema duală a SVM:

$$\begin{aligned} \max_{\alpha} \quad & W(\alpha) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m y^{(i)} y^{(j)} \alpha_i \alpha_j \langle x^{(i)}, x^{(j)} \rangle. \\ \text{s.t.} \quad & \alpha_i \geq 0, \quad i = 1, \dots, m \\ & \sum_{i=1}^m \alpha_i y^{(i)} = 0, \end{aligned}$$

Genele din cromozomi vor fi alpha-urile, funcția de fitness va fi $W(\alpha)$.

Funcția de fitness:

$$W(\alpha) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m y^{(i)} y^{(j)} \alpha_i \alpha_j \langle x^{(i)}, x^{(j)} \rangle.$$

Soluția va fi cromozomul care maximizează W , apoi din alpha-uri se refac parametrii w și b și apoi se face clasificarea.

$$w = \sum_{i=1}^m \alpha_i y^{(i)} x^{(i)}.$$

Deoarece setul de date de antrenare are 1728 de instanțe – un număr foarte mare pentru a putea fi considerat ca număr de gene - fiecare individ α va fi compus din 100 de gene (numărul de gene este dat de numărul de instanțe, adică 100 (62 instanțe din clasa unacc, 20 instanțe din clasa acc, 9 instanțe din clasa good, 9 instanțe din clasa vgood)):

$$\alpha = [\alpha_1 \ \alpha_2 \ \alpha_3 \ \dots \ \alpha_{100}]$$

În plus, fiecare individ trebuie să respecte două constrângeri:

$$0 \leq \alpha_i \leq C \quad (1)$$

$$\sum_{i=1}^n \alpha_i y_i = 0 \quad (2)$$

După ce obținem multiplicatorii Lagrange prin intermediul algoritmului genetic, se calculează bias-ul cu ajutorul formulei:

$$b = \frac{1}{n} \sum_{i=1}^n \left(y_i - \sum_{j=1}^m (y_i \alpha_j \mathbf{x}_i^T \mathbf{x}_j) \right)$$

Parametrii specifici SVM utilizați în cadrul implementării sunt:

- Nucleul polinomial de grad II.
- Costul $C = 1e-5$.

SECȚIUNEA 3: Modalitatea de rezolvare

Mașinile cu vector suport pleacă de la o fundamentare matematică riguroasă, deci nu sunt modele empirice în care se fac încercări pentru testarea funcționalității sau pentru modificarea numărului de straturi sau de neuroni pentru a mai face o încercare.

Ideea de bază a SVM este: cum putem avea garanții că un model obține soluții bune pe mulțimea de test?

Cheia problemei este că datele de test sunt din aceeași distribuție de probabilitate cu datele de antrenare. Acestea nu sunt complet diferite, ci trebuie să respecte o anumită structură. Așadar, atât datele de antrenare, cât și datele de test trebuie să fie eșantionate din aceeași distribuție de probabilitate, adică din aceeași clasă mare de obiecte.

SVM este prin excelență un clasificator binar, însă problema de clasificare ce trebuie rezolvată implică existența a patru clase distincte. Așadar, se impune o problemă de clasificare cu clase multiple. În acest sens, am optat pentru abordarea „una versus toate”. Astfel, am obținut cele 4 modele:

- unacc VS others
- acc VS others
- good VS others
- v-good VS others

Pentru fiecare dintre cele 4 modele se utilizează separat algoritmul evolutiv pentru a determina coeficienții Lagrange, respectiv bias-ul.

În final, pentru a clasifica o nouă instanță, se apelează toate cele 4 modele și se alege clasa cu funcția de decizie maximă:

$$C = \operatorname{argmax}_{i=1..k} (f(\mathbf{x}))$$

SECȚIUNEA 4: Listarea părților semnificative din codul sursă însoțite de explicații și comentarii

// Cuantificarea unui set de date

```
String line;

// Class 1 set quantification - unacc
string quantified = "";
try
{
    // Pass the file path and file name to the StreamReader constructor
    StreamReader sr = new StreamReader("D:\\Support_Vector_Machine\\dataSet.data");
    // Read the first line of text
    line = sr.ReadLine();

    char[] delimiterChars = { ',' };
    int[] quantifiedValues = new int[7];
    string quantifiedValuesString = "";

    // Continue to read until you reach end of file
    while (line != null)
    {
        string[] words = line.Split(delimiterChars);

        switch (words[0])
        {
            case "low":
                quantifiedValues[0] = 1;
                break;
            case "med":
                quantifiedValues[0] = 2;
                break;
            case "high":
                quantifiedValues[0] = 3;
                break;
            case "vhigh":
                quantifiedValues[0] = 4;
                break;
        }

        switch (words[1])
        {
            case "low":
                quantifiedValues[1] = 1;
                break;
            case "med":
                quantifiedValues[1] = 2;
                break;
            case "high":
                quantifiedValues[1] = 3;
                break;
            case "vhigh":
                quantifiedValues[1] = 4;
                break;
        }
    }
}
```

```

switch (words[2])
{
    case "2":
        quantifiedValues[2] = 1;
        break;
    case "3":
        quantifiedValues[2] = 2;
        break;
    case "4":
        quantifiedValues[2] = 3;
        break;
    case "5more":
        quantifiedValues[2] = 4;
        break;
}

switch (words[3])
{
    case "2":
        quantifiedValues[3] = 1;
        break;
    case "4":
        quantifiedValues[3] = 2;
        break;
    case "more":
        quantifiedValues[3] = 3;
        break;
}

switch (words[4])
{
    case "small":
        quantifiedValues[4] = 1;
        break;
    case "med":
        quantifiedValues[4] = 2;
        break;
    case "big":
        quantifiedValues[4] = 3;
        break;
}

switch (words[5])
{
    case "low":
        quantifiedValues[5] = 1;
        break;
    case "med":
        quantifiedValues[5] = 2;
        break;
    case "big":
        quantifiedValues[5] = 3;
        break;
}

if (words[6] == "unacc")
{
    quantifiedValues[6] = 1;
}

```



```

        else
        {
            quantifiedValues[6] = -1;
        }

        for (int i = 0; i < 6; ++i)
        {
            quantifiedValuesString += quantifiedValues[i].ToString();
            quantifiedValuesString += ",";
        }
        quantifiedValuesString += quantifiedValues[6].ToString();

        quantified += quantifiedValuesString;
        quantified += '\n';

        // Read the next line
        line = sr.ReadLine();
        quantifiedValuesString = "";
    }

    using (StreamWriter sw = new
StreamWriter(@"D:\Support_Vector_Machine\dataSetQuantifiedUnacc.data"))
    {
        quantified = quantified.Trim();
        sw.Write(quantified);
    }

    // Close the file
    sr.Close();
}

catch (Exception e)
{
    Console.WriteLine("Exception: " + e.Message);
}
finally
{
    Console.WriteLine("Writing to the dataSetQuantifiedUnacc.data file was completed
successfully.");
}

```

// Funcția de fitness

```
public void ComputeFitness(Chromosome c, int[,] data)
{
    double sumAlfa = 0;

    int[] v1 = new int[6];
    int[] v2 = new int[6];

    for (int i = 0; i < c.NoGenes; ++i)
    {
        sumAlfa += c.Genes[i];
    }

    double sumAlfaiAlfaj = 0;
    for (int i = 0; i < c.NoGenes; ++i)
    {
        for (int j = 0; j < c.NoGenes; ++j)
        {
            for (int ind1 = 0; ind1 < 6; ++ind1)
            {
                v1[ind1] = data[i, ind1];
            }

            for (int ind2 = 0; ind2 < 6; ++ind2)
            {
                v2[ind2] = data[j, ind2];
            }
            sumAlfaiAlfaj += c.Genes[i] * c.Genes[j] * data[i, 6] * data[j, 6] *
kernel1.PolinomialKernel(v1, v2);
        }
    }
    c.Fitness = sumAlfa - 0.5 * sumAlfaiAlfaj;
}
```

// Determinarea soluției problemei – determinarea parametrilor alpha

```
Chromosome[] population = new Chromosome[populationSize];
for (int i = 0; i < population.Length; i++)
{
    population[i] = p.MakeChromosome();
    p.ComputeFitness(population[i], data);
}

for (int gen = 0; gen < maxGenerations; gen++)
{
    Chromosome[] newPopulation = new Chromosome[populationSize];
    newPopulation[0] = Selection.GetBest(population); // elitism

    for (int i = 1; i < populationSize; i++)
    {
        // Select 2 parents: Selection.Tournament
        Chromosome p1 = new Chromosome(Selection.Tournament(population));
        Chromosome p2 = new Chromosome(Selection.Tournament(population));

        // Generating a child by crossover application: Crossover.Arithmetic
        Chromosome c = new Chromosome(Crossover.Arithmetic(p1, p2, crossoverRate));

        // Applying mutation to the child: Mutation.Reset
        Mutation.Reset(c, mutationRate);

        // Calculating the value of the fitness function for the child:
        // ComputeFitness from the problem p
        p.ComputeFitness(c, data);

        // Introduction of child in newPopulation
        newPopulation[i] = c;
    }

    for (int i = 0; i < populationSize; i++)
        population[i] = newPopulation[i];

    Console.WriteLine("{0}% \r", (gen + 1));
}
Console.WriteLine();
Console.WriteLine("Complete processing.");
return Selection.GetBest(population);
```

// Funcția de calcul pentru nucleul polinomial de grad II

```
public int PolinomialKernel(int[] v1, int[] v2)
{
    int result = 0;
    for (int i = 0; i < v1.Length; ++i)
    {
        result += (v1[i] * v2[i]);
    }
    result += 1;
    result *= result;

    return result;
}
```

// Script de ajustare a valorilor alpha pentru a respecta cele 2 condiții impuse de SVM

```
int k;
double sPlus;
double sMinus;

// Repair algorithm
while (Math.Abs(s) > limitaSuma)
{
    sPlus = 0.0;
    sMinus = 0.0;

    for (int i = 0; i < trainInstances; ++i)
    {
        if (i >= 91)
        {
            sPlus += newAlfa[i] * data[i, 6];
        }
        else
        {
            sMinus += newAlfa[i] * data[i, 6];
        }
    }
    sMinus = Math.Abs(sMinus);

    if (Math.Abs(sPlus) > Math.Abs(sMinus))
    {
        k = rand.Next(91, trainInstances);
    }
    else
    {
        k = rand.Next(0, 91);
    }

    if (newAlfa[k] > Math.Abs(s))
    {
        newAlfa[k] = newAlfa[k] - Math.Abs(s);
    }
    else
    {
        newAlfa[k] = 0;
    }

    s = 0;
    for (int i = 0; i < trainInstances; ++i)
    {
        s += newAlfa[i] * data[i, 6];
    }
}
```

// Calcul Bias

```
double bias = 0;
double biasSum;

int[] v1 = new int[6];
int[] v2 = new int[6];

for (int i = 0; i < trainInstances; ++i)
{
    biasSum = 0;
    for (int j = 0; j < trainInstances; ++j)
    {
        for (int ind1 = 0; ind1 < 6; ++ind1)
        {
            v1[ind1] = data[i, ind1];
        }

        for (int ind2 = 0; ind2 < 6; ++ind2)
        {
            v2[ind2] = data[j, ind2];
        }

        biasSum += data[j, 6] * newAlfa[j] * kernel.PolinomialKernel(v1, v2);
    }
    bias += data[i, 6] - biasSum;
}
bias /= trainInstances;
```

// Calculare valoare de clasificare

```
Console.WriteLine("Calculation of classification value: ");
double[] classificationValue = new double[noInstances];

int[] trainV = new int[6];
int[] testV = new int[6];

double val;
for (int j = 0; j < noInstances; ++j)
{
    val = 0;
    for (int i = 0; i < trainInstances; ++i)
    {
        for (int ind1 = 0; ind1 < 6; ++ind1)
        {
            trainV[ind1] = data[i, ind1];
        }

        for (int ind2 = 0; ind2 < 6; ++ind2)
        {
            testV[ind2] = dataTest[j, ind2];
        }
        val += newAlfa[i] * data[i, 6] * kernel.PolinomialKernel(trainV, testV);
    }
    classificationValue[j] = val;
}
```

// Determinarea claselor

```
int countUnacc = 0;
int countAcc = 0;
int countGood = 0;
int countVgood = 0;

int[] classType = new int[noInstances];
for (int i = 0; i < noInstances; ++i)
{
    if (classifiedValues1[i] >= classifiedValues2[i] && classifiedValues1[i] >=
classifiedValues3[i] && classifiedValues1[i] >= classifiedValues4[i])
    {
        classType[i] = 1;
        countUnacc++;
    }

    else if (classifiedValues2[i] >= classifiedValues1[i] && classifiedValues2[i] >=
classifiedValues3[i] && classifiedValues2[i] >= classifiedValues4[i])
    {
        classType[i] = 2;
        countAcc++;
    }

    else if (classifiedValues3[i] >= classifiedValues1[i] && classifiedValues3[i] >=
classifiedValues2[i] && classifiedValues3[i] >= classifiedValues4[i])
    {
        classType[i] = 3;
        countGood++;
    }

    else if (classifiedValues4[i] >= classifiedValues1[i] && classifiedValues4[i] >=
classifiedValues2[i] && classifiedValues4[i] >= classifiedValues3[i])
    {
        classType[i] = 4;
        countVgood++;
    }
}
```

Mai jos, este prezentat rezultatul rulării algoritmului evolutiv pentru clasa acceptabile. Se poate observa că numărul de zerouri este unul rezonabil, astfel numărul vectorilor suport va fi mai mic decât jumătatea numărului total de soluții alpha obținute. Ideal, ar fi fost un număr chiar mai redus de vectori suport. Se observă că valorile alpha-urilor respectă ambele condiții impuse de problema duală a SVM.

Rezultatele obținute în urma rulării script-ului de clasificare se pot observa în imaginea de mai jos. Este posibil ca o serie de modificări asupra parametrilor SVM-ului să conducă la rezultate mai precise.

SECȚIUNEA 6: Concluzii

În concluzie, SVM reprezintă o metodă precisă de clasificare, însă dificultatea acestui model apare dintr-o serie de motive precum: alegerea potrivită a tipurilor parametrilor (nucleu), dar și a valorilor acestora (cost), precum și ajustarea seturilor de date utilizate pentru antrenare și testare, astfel încât acestea să respecte aceeași distribuție de probabilitate. Așadar, plecând de la teoria de învățare statistică în încercarea de a rezolva o problemă într-un mod riguros explicit, se obțin rezultate bune și din punct de vedere practic.

SECȚIUNEA 7: Bibliografie

1. <https://archive.ics.uci.edu/ml/index.php>
2. <https://archive.ics.uci.edu/ml/datasets/Car+Evaluation>
3. Evolutionary Support Vector Machines: A Dual Approach - Madson Luiz Dantas Dias, Ajalmar R. Rocha Neto - Department of Teleinformatics, Federal Institute of Ceará