

Tema practica

Taga Stefan Razvan - 3A2

January 10, 2024

Algorithm 1 Soluționarea Problemei de Clasificare a Fișierelor Text

```
1: Pregătirea Datelor:
2: for fiecare subdir in directorul "lemn" do
3:   for fiecare fișier cu extensia ".txt" în subdir do
4:     citește conținutul fișierului
5:     adaugă conținutul la lista de date (data)
6:     adaugă subdir la lista de etichete (labels)
7:   end for
8: end for
9: Extragerea de Caracteristici:
10: folosește un vectorizator pentru a extrage caracteristici din date
11: Implementarea Algoritmului de Învățare Automată:
12: împarte datele în seturi de antrenare și testare
13: inițializează un clasificator Naive Bayes (MultinomialNB)
14: antrenează clasificatorul pe setul de antrenare
15: Antrenarea Modelului:
16: folosește setul de testare pentru a face predicții
17: Evaluarea Modelului:
18: calculează acuratețea modelului
19: Predicții:
20: folosește modelul antrenat pentru a face predicții pentru noi fișiere
21: Analiza și Interpretare:
22: afișează raportul de clasificare
```

Vom utiliza un clasificator Naive Bayes, care este potrivit pentru problemele de clasificare a textelor.

```
1 import os
2 from sklearn.feature_extraction.text import
  CountVectorizer
3 from sklearn.model_selection import train_test_split
4 from sklearn.naive_bayes import MultinomialNB
5 from sklearn.metrics import accuracy_score,
  classification_report
```

```

6
7 # Preg tiera Datelor
8 data = []
9 labels = []
10
11 directory = "lemn"
12
13 for subdir in os.listdir(directory):
14     if os.path.isdir(os.path.join(directory, subdir)):
15         for filename in
16             os.listdir(os.path.join(directory,
17                                     subdir)):
18                 if filename.endswith(".txt"):
19                     with open(os.path.join(directory,
20                                             subdir, filename), "r") as file:
21                         content = file.read()
22                         data.append(content)
23                         labels.append(subdir)
24
25 # Extragerea de Caracteristici
26 vectorizer = CountVectorizer()
27 X = vectorizer.fit_transform(data)
28
29 # Implementarea Algoritmului de nvare Automat
30 X_train, X_test, y_train, y_test =
31     train_test_split(X, labels, test_size=0.2,
32                     random_state=42)
33
34 classifier = MultinomialNB()
35 classifier.fit(X_train, y_train)
36
37 # Antrenarea Modelului
38 y_pred = classifier.predict(X_test)
39
40 # Evaluarea Modelului
41 accuracy = accuracy_score(y_test, y_pred)
42 print(f"Accuracy: {accuracy:.2f}")
43
44 # Analiza i Interpretare
45 print("\nClassification Report:")
46 print(classification_report(y_test, y_pred))

```

Listing 1: Implementarea algoritmului în Python

Justificarea Alegerii Clasificatorului Naive Bayes

1. **Text Classification:** Problema pare a fi una de clasificare a textului, deoarece dorim să clasificăm fișierele text în funcție de numele lor.
2. **Simplicitate:** Algoritmii Naive Bayes sunt cunoscuți pentru simplitatea lor și funcționează bine în multe probleme de clasificare a textului. Aceștia sunt eficienți și nu necesită multe eforturi de ajustare a parametrilor.
3. **Număr de Caracteristici Limitat:** În acest exemplu, am ales să extragem doar numărul de apariții ale cuvintelor (termenilor) din fișierele text, ceea ce este potrivit pentru un clasificator Naive Bayes.
4. **Scalabilitate:** Naive Bayes este scalabil și poate funcționa bine și cu seturi de date mai mari, ceea ce îl face potrivit pentru probleme de dimensiuni moderate.
5. **Bun Început:** Naive Bayes este adesea considerat un "bun început" pentru problemele de clasificare a textului, iar rezultatele pot fi satisfăcătoare, mai ales dacă datele îndeplinesc asumțiunile naive ale algoritmului.