

CURS 1: EXTRAȚIA ȘI TRANSFORMAREA DATELOR

1: EXTRACTIA ȘI TRANSFORMAREA DATELOR

Extracția și transformarea datelor

Curs 1 – Web mining și data analysis Durată: 1 oră și 20 minute

2: IMPORTANȚA DATELOR ÎN AI

- Modelele AI se bazează pe date de calitate pentru acuratețe.
- Garbage in -> Garbage out: datele slabe duc la modele slabe.
- Datele sunt necesare pentru:
 - Antrenarea modelelor de învățare automată
 - Realizarea de predicții
 - Înțelegerea tendințelor și comportamentelor
- Exemplu: Un model de detectare a spamului eșuează atunci când este antrenat pe etichete de e-mail incorecte

3: TEHNICI DE COLECTARE A DATELOR

API-uri (Interfețe de Programare a Aplicațiilor) – Oferă acces direct la date structurate din diverse surse. – Exemplu:

Utilizarea API-ului Coindesk pentru a colecta titluri despre bitcoin pentru analiza sentimentului

Web Scraping – Extrage informații de pe site-uri web atunci când API-urile nu sunt disponibile. – Necesită parsarea HTML. – Unele: BeautifulSoup, Scrapy. – Exemplu: Scraping al prețurilor produselor de pe un site de e-commerce

Baze de date – Bazele de date SQL și NoSQL stochează date structurate. – Exemplu: Interogarea unei baze de date Sqlite pentru a obține date despre tranzacțiile utilizatorilor

4: PREZENTARE GENERALĂ A PREPROCESĂRII DATELOR

- Datele nu sunt rareori curate; preprocesarea este esențială.
- Pași importanți:
 - Gestionarea valorilor lipsă
 - Normalizare și scalare
 - Codificarea variabilelor categorice
 - Detectarea și eliminarea valorilor aberante
- Exemplu: Curățarea unei baze de date de clienți prin eliminarea înregistrărilor duplicate

5: GESTIONAREA VALORILOR LIPSĂ

- Datele lipsă pot afecta performanța modelului.
- Strategii:
 - Eliminarea rândurilor/coloanelor cu prea multe valori lipsă.
 - Imputarea valorilor lipsă folosind media, mediana sau moda.
 - Prezicerea valorilor lipsă folosind modele de învățare automată.
- Exemplu: Gestionarea valorilor lipsă pentru vârstă într-un set de date despre supraviețuirea pe Titanic

6: NORMALIZARE ȘI SCALARE

- Normalizarea scalează datele la un interval fix (ex. 0-1).
- Standardizarea centrează datele în jurul valorii zero.
- Esențială pentru modelele care se bazează pe distanțe (ex. k-NN, SVM).
- Exemplu: Scalarea caracteristicilor numerice într-un set de date despre prețurile locuințelor

7: CODIFICAREA CARIABILELOR CATEGORICE

- Modelele de învățare automată necesită date numerice.
- Strategii:
 - One-hot encoding (pentru date nominale)
 - Label encoding (pentru date ordinale)
- Exemplu: Codificarea genului (masculin/feminin) în valori binare

8: INGINERIA ȘI SELECTIA CARACTERISTICILOR

- **Ingineria Caracteristicilor:** Crearea unor noi caracteristici semnificative din date existente.
- **Selectia Caracteristicilor:** Eliminarea caracteristicilor irelevante sau redundante.
- Tehnici:
 - Analiza corelației
 - Eliminarea recursivă a caracteristicilor (RFE)
 - Analiza Componentelor Principale (PCA)
- Exemplu: Extragerea cuvintelor cheie din recenziile produselor pentru îmbunătățirea analizei sentimentului

9: GESTIONAREA SETURILOR DE DATE DEZECHILIBRATE

- Problema: O clasă domină setul de date.
- Soluții:
 - Resampling (supra-eșantionarea minorității, sub-eșantionarea majorității).
 - Utilizarea funcțiilor de pierdere ponderate în antrenarea modelului.
 - Generarea de date sintetice (tehnica SMOTE).
- Exemplu: Gestionarea dezechilibrului de clasă în seturile de date pentru detectarea fraudelor

10: REZUMAT & CONCLUZII

- Datele de înaltă calitate sunt esențiale pentru succesul AI.
- Metode de colectare a datelor: API-uri, web scraping, baze de date.
- Pași de preprocesare: curățare, gestionarea valorilor lipsă, scalare, codificare.
- Ingineria și selecția caracteristicilor îmbunătățesc performanța modelului.
- Gestionarea dezechilibrului de clasă este crucială pentru predicții corecte.

11: CONȚINUTUL SEMINARIULUI

- Extracția datelor dintr-un API (ex. Twitter, OpenWeather).
- Exemplu de web scraping: extragerea detaliilor produselor de pe un site de e-commerce.
- Preprocesarea și curățarea unui set de date din lumea reală.
- Gestionarea valorilor lipsă și ingineria caracteristicilor.
- Aplicarea tehnicilor de scalare și codificare în practică.

12: ÎNTREBĂRI & DISCUȚII

- Sesiune deschisă pentru întrebări și discuții.