

Data Mining

Lab 5:

Reguli de asociere
Modele de regresie

1. Reguli de asociere/

Exemplu (problema cosului de cumparaturi). Se consideră un set de tranzacții (T_1, T_2, \dots, T_n), fiecare tranzacție conținând un set de produse achiziționate. De exemplu:

T_1 : {pâine, lapte, apă}

T_2 : {pâine, carne, apă}

T_3 : {pâine, unt, carne, apă}

T_4 : {fructe, apă}

Se pune problema identificării produselor care sunt frecvent cumpărate împreună și a unor reguli de asociere de forma **IF pâine AND apă THEN carne**.

Regulile de asociere sunt de forma **IF A THEN B**. Termenul A joacă rol de antecedent, iar B rol de consecință (totuși nu exprimă relații de cauzalitate ci doar de corelare).

Dintr-un set de tranzacții pot fi extrase numeroase reguli – e necesar să poată fi evaluată relevanța lor pentru a fi ierarhizate.

Pentru evaluarea relevanței unei reguli se folosesc cel puțin două mărimi:

- **Support (support):** $\text{supp}(A \rightarrow B) = \text{numărul de tranzacții ce conțin pe A și B} / \text{numărul total de tranzacții}$
- **Incredere (confidence):** $\text{conf}(A \rightarrow B) = \text{numărul de tranzacții ce conțin pe A și B} / \text{numărul de tranzacții ce conțin pe A}$

Exemplu: IF pâine și apă THEN carne

$A = \{\text{pâine, apă}\}$, $B = \{\text{carne}\}$

$\text{Supp}(A \rightarrow B) = 2/4 = 0.5$

$\text{Conf}(A \rightarrow B) = 2/3 = 0.6$

Obs: pe lângă aceste măsuri se folosesc și indicatori ai noutății regulii (cât este de interesantă sau neobișnuită regula); un exemplu de astfel de indicator este cel denumit **lift**:

$\text{Lift}(A \rightarrow B) = \text{prob}(A, B) / (\text{prob}(A) \text{prob}(B))$

(probabilitatea se estimează prin frecvența relativă)

Regula e considerată interesantă dacă valoarea măsurii lift este mare. Dacă valoarea e apropiată de 1 aceasta sugerează că A și B nu sunt corelate deci nu poate fi extrasă o regulă utilă de forma $A \rightarrow B$

Exemplu: R=IF pâine AND carne THEN apă
Conf(R)=2/2=1
Lift(R)=0.5/(0.5*1)=1

Algoritm de extragere a regulilor de asociere din date (algoritmul APRIORI)

Date de intrare: set de tranzacții (fiecare tranzacție conține o listă de entități)

Parametri de control:

- Prag pentru suport minim (ex: 0.2)
- Nivel minim de încredere (ex: 0.9)

Structura general a algoritmului Apriori:

Pas 1: identificare subseturi cu suport semnificativ (mai mare decât pragul minim) – “frequent itemsets”; identificarea acestor subseturi se bazează pe:

- Identificarea subseturilor frecvente cu un singur element (lista L_1)
- FOR $k=1, K$ DO construirea listei L_k cu subseturi frecvente având k elemente pornind de la lista L_{k-1} (subseturi frecvente cu $k-1$ elemente)

Pas 2: construirea regulilor prin partiționarea subseturilor identificate la pasul 1 în două părți (o parte pentru antecedentul regulii și o parte pentru consecință); se rețin doar regulile care au nivelul de încredere mai mare decât pragul

Exercițiul 1.

- Deschideți în Weka setul de date [supermarket.arff](#)
- Determinați reguli de asociere folosind [Associate->Apriori](#) și valorile implicite ale parameterilor
- Aplicați același algoritm pentru alte valori ale pragului pentru suport ([lowerBoundMinSupport=0.2](#)) și pentru încredere ([minMetric=0.75](#)).

2. Modele de regresie

2.1. Regresie liniara

În modelele liniare dependența dintre variabila (variabilele) prezise și cele predictor este descrisă printr-o funcție liniară de forma $Y=WX$. Parametrii modelului (elementele vectorului/matricii W) se determină pornind de la date folosind o tehnică de minimizare a sumei pătratelor erorilor

Exercițiul 2.

- Deschideți în Weka fisierul [autoPrice.arff](#)
- Utilizați [Class->Functions->SimpleLinearRegression](#) pentru a determina o dependență liniară simplă între atributul de ieșire (preț) și cel mai relevant dintre atributele de intrare. Analizați valorile corespunzătoare coeficientului de corelație și erorii [Correlation Coefficient](#) și [Mean Absolute Error](#).
- Utilizați [Class->Functions->LinearRegression](#) pentru a determina o dependență liniară între atributul de ieșire (preț) și cele mai relevante dintre atributele de intrare. Analizați valorile

corespunzătoare coeficientului de corelație și erorii **Correlation Coefficient** și **Mean Absolute Error**. Analizați influența metodei de selecție a atributelor (**attributeSelectionMethod**)

2.2. Regresie neliniară

Exercițiul 3. (tot pentru fisierul **autoPrice.arff**)

- d) Utilizați **Class->Functions->MultilayerPerceptron** cu valorile implicite ale parametrilor. Analizați valorile corespunzătoare coeficientului de corelație și erorii **Correlation Coefficient** și **Mean Absolute Error**.
- e) Identificați în categoria **Class->Trees** varianta care permite construirea unui arbore de regresie

Obs. Versiunile de Weka mai mici de 3.8 include și implementare a rețelelor RBF (**RBF Network**).

Exercițiul 4. Aplicați prelucrările de la Exercițiile 2 și 3 în cazul setului de date **autoMPG.arff** și analizați diferențele (în particular în ceea ce privește arborii de regresie).