

Lab 7: Data Mining.
Serii temporale
Metode de tip ansamblu

1. Serii temporale

Analiza seriilor temporale are ca scop să modeleze și să explice dependența unor date de momente de timp successive. Exemple tipice de serii temporale sunt: temperatura înregistrată zilnic, curs de schimb valutar, prețul unor acțiuni etc.

Principalele prelucrări care pot fi efectuate asupra unei serii de timp sunt:

- *Pre-procesare* (de exemplu, transformarea seriei prin normalizare sau standardizare, completarea valorilor absente prin interpolare, eliminarea zgomotului prin netezire, eliminarea tendinței prin calcularea diferențelor dintre elemente succesive etc)
- *Predicție*: estimarea valorilor ulterioare din serie pe baza valorii curente și a celor anterioare (folosind un model care descrie dependența valorii curente din serie de valorile anterioare).

Un proces de predicție este caracterizat prin:

- *Intrare*: datele de intrare sunt valori anterioare din serie
- *Iesire*: rezultatul reprezintă valoarea/valorile următoare din serie
- *Model*: un model de regresie care descrie legătura dintre valoarea curentă a seriei și valorile anterioare (numărul de valori anterioare despre care se consideră că influențează valoarea curentă este denumit întârzierea seriei (*time-lag*))

Considerăm seria X_1, X_2, \dots, X_n și întârzierea T . Deci valoarea curentă X_i depinde de valorile $X_{i-1}, X_{i-2}, \dots, X_{i-T}$. Prin urmare secvența de valori din serie poate fi transformată într-un alt set de date în care sunt T atribute predictor și un atribut prezis:

<i>Atribute predictor</i>	<i>Atribut prezis</i>
$X_1 \ X_2 \ \dots \ X_i \ \dots \ X_T$	X_{T+1}
$X_2 \ X_3 \ \dots \ X_{i+1} \ \dots \ X_{T+1}$	X_{T+2}
\dots	
$X_{n-T} \ X_{n-T+1} \ \dots \ X_{n-i} \ \dots \ X_{n-1}$	X_n

Folosind acest set de date se poate construi un model de regresie (în aceeași manieră ca pentru date care nu sunt temporale). Una dintre principalele dificultăți este alegerea adecvată a valorii T .

Exercițiul 1.

- a) Deschideți fișierul [airlines.arff](#) (conținând nr de pasageri ai unei companii aeriene înregistrat lunar în perioada 1949 – 1960)
- b) Construiți un nou set de date folosind o întârziere $T=12$. *Indicație*: utilizați Weka pt eliminarea atributului corespunzător datei și Excel (sau un limbaj de programare) pt construirea noului set de date
- c) Aplicați un model de regresie pentru noul set de date și analizați rezultatele obținute

Exercitiul 2.

(doar pt versiune Weka >=3.7.3)

- Instalati pachetul **Time Series Forecasting** utilizand **Weka GUI Chooser ->Tools->Package manager** si selectand pentru instalare **timeSeriesForecasting**
- Deschideti fisierul **airlines.arff**
- Preziceti urmatoarele 6 valori utilizand unul dintre urmatoarele modele: (i) linear regression; (ii) multilayer perceptron; (iii) random forests. *Indicatie:* selectia modelului se realizeaza utilizand panelul **Advanced Configuration->Based Learner**

Obs: detalii privind pachetul **TimeSeriesForecasting** pot fi gasite la <http://wiki.pentaho.com/display/DATAMINING/Time+Series+Analysis+and+Forecasting+with+Weka>.

(pentru Weka 3.8.1) se utilizează **Forecast (Basic configuration)** și se specifică doar **Number of time units to forecast** (6 dacă se dorește estimarea următoarelor 6 valori).

2. Metode de tip ansamblu (ensemble models)

Sunt meta-modele care se obțin din câteva modele de bază antrenate pe același set sau pe seturi diferite de antrenare. Există mai multe variante de a construi modele de tip ansamblu:

- Utilizând modele bazate pe algoritmi diferiți antrenati pe același set de date (e.g. *bucket of models*)
- Utilizând modele bazate pe același algoritm dar antrenate pe seturi diferite de date (e.g. *bagging and boosting*)
- Utilizând diferite modele și împărțind setul de date (e.g. *stacking*)

Exercitiul 3. Utilizand Weka Experimenter comparați performanța următoarelor metamodele: **Vote, Bagging, Random Forest, AdaBoost** si **Stacking** pt seturile de date: **iris.arff, glass.arff**

- Utilizați valorile implicite ale parametrilor
- Imbunătățiți comportamentul pt **Vote, Bagging** si **AdaBoost** înlocuind clasificatorul de bază cu alt clasificator.