

CZU: 004.434:004.82

PREZICEREA PERFORMANȚELOR STUDENȚILOR FOLOSIND ÎNVĂȚAREA AUTOMATĂ (Machine Learning)

Maria CRISTEI, Ghenadie MARIN, Victor STELEA

Universitatea de Stat din Moldova

În prezent, învățarea automată (*machine learning*) ocupă un loc important în inteligența artificială, preocupându-se de dezvoltarea algoritmilor ce permit unui sistem informatic să învețe date, reguli și algoritmi. Învățarea automată presupune în primul rând identificarea și implementarea unei modalități cât mai eficiente de reprezentare a informațiilor, în sensul facilitării căutării, reorganizării și modificării acestora. În acest sens, în prezentul articol se descrie utilitatea și aplicabilitatea tehnicilor de învățare automată supervizată la problemele de predicție și implementarea acestora în dezvoltarea aplicațiilor informatice. Aplicația elaborată este unică prin felul ei de executare a modelului *machine learning* de predicție. Metodologia folosită în aplicația elaborată este mixtă, cuprinzând tehnologii complexe de ultimă oră: mediul de dezvoltare Jupyter Notebook, limbajul de programare Python împreună cu cele mai populare librării ale acestuia utilizate în *machine learning*, instrumente de dezvoltare a aplicației web Flask.

Cuvinte-cheie: *învățare automată, inteligență artificială, sisteme inteligente.*

PREDICTING STUDENT PERFORMANCE USING MACHINE LEARNING

At present, machine learning occupies an important place in artificial intelligence, and is concerned with the development of algorithms that allow an information system to learn data, rules, and algorithms. Automatic learning involves first and foremost the identification and implementation of a more efficient way of representing information in order to facilitate search, reorganization and change. In this respect, this article describes the utility and applicability of supervised automated learning techniques to prediction problems and their implementation in the development of computer applications. The elaborate application is unique in its way of executing the Machine learning prediction model. The methodology used in the developed application is mixed, including state-of-the-art complex technologies: the Jupyter Notebook development environment, the Python programming language along with its most popular bookkeeping libraries used in machine learning, Flask web application development tools.

Keywords: *machine learning, artificial intelligence, intelligent systems.*

Considerații generale privind învățarea automată (*machine learning*)

De-a lungul timpului, savanții au definit diferit noțiunea de *inteligență artificială*. Este de remarcat că toate definițiile pot fi împărțite în patru clase, în funcție de cum este înțeleasă „inteligența”:

- a) sisteme care gândesc la fel ca oamenii;
- b) sisteme care se comportă la fel ca oamenii;
- c) sisteme care gândesc rațional;
- d) sisteme care se comportă rațional [1].

Inteligența artificială este un concept care se referă la capacitatea unui sistem de a acționa într-un mod „inteligent”, cum ar fi: să învețe singur, să se adapteze și să reacționeze în situații total noi. Practic, aceasta ar însemna că inteligența artificială ar ieși din limitele în care a fost programată de cineva și ar acționa independent, indiferent de consecințe (pozitive sau negative). Astăzi un astfel de sistem nu există în realitate. Scopurile de bază ale cercetărilor în inteligența artificială sunt evidențiate de subdomeniile acesteia: prelucrările simbolice, procesarea simbolică și limbajele de procesare simbolică, procesarea limbajului natural, raționamentul automat și demonstrarea teoremelor etc.

Un sistem înzestrat cu „inteligență” este capabil să învețe pentru a-și îmbunătăți interacțiunea cu mediul. Mediul oferă stimuli sau informație elementului de învățare, care folosește această informație pentru a îmbunătăți cunoștințele (explicite) din baza de cunoștințe. Aceste cunoștințe sunt utilizate de componenta de prelucrare (rezolvare) în rezolvarea problemei. Sistemul învață dacă își îmbunătățește performanțele la îndeplinirea sarcinii pe baza experienței. Învățarea denotă schimbările dintr-un sistem, astfel încât acesta să poată realiza: (a) aceeași sarcină, mai eficient; (b) sarcini noi, posibil similare.

În funcție de diferența dintre nivelul informației oferite de mediu și cel al informației din baza de cunoștințe, pot fi identificate patru tipuri de învățare: (1) învățarea prin memorare; (2) învățarea prin instruire; (3) învățarea

prin inducție (din exemple); (4) învățarea prin analogie. Astfel, un sistem inteligent poate învăța fie din greșeli, prin autoperfecțiune, ghidat de profesor, generalizând, prin analogie, sau din exemple pozitive sau negative.

Metodele și tehnicile de achiziționare a cunoștințelor sunt *manuale* (analistul de cunoștințe) și *automate* (machine learning). Domeniul învățării automate formalizează aceste metode și caută aplicarea lor la sistemele automate.

Sistemele de învățare automată reprezintă varianta de inteligență artificială aplicabilă în prezent: în baza unor algoritmi matematici ce sunt capabili să opereze cu un volum foarte mare de date, să învețe singure date, reguli, algoritmi și să-și perfecționeze acțiunile. Învățarea automată presupune în primul rând identificarea și implementarea unei modalități cât mai eficiente de a reprezenta informații, în sensul facilitării căutării, reorganizării și modificării acestora. Diferența esențială față de inteligența artificială „veritabilă”, care ar trebui să acționeze complet independent, rezidă în faptul că *sistemele bazate pe învățare automată își desfășoară activitatea sub controlul uman și în cadrul parametrilor stabiliți de cercetători*. Aplicarea sistemelor inteligente bazate pe învățare automată este destul de vastă: *recunoaștere de imagini și semnal vocal* (recunoașterea scrisului de mână, detecția fețelor, înțelegerea limbajului vorbit); *vedere artificială* (detecția obstacolelor, recunoașterea amprentelor); *supraveghere bio*; *controlul roboților*; *predicția vremii*; *diagnosticare medicală*; *detecția fraudelor*.

Studiul învățării automate a dus la descrierea a numeroase metode, variind după scop, date de antrenament, a strategiei de învățare și a modalității de reprezentare a datelor. Astfel, tipologia sistemelor inteligente (SI) bazate pe învățare automată este următoarea:

- în funcție de scopul urmărit:
 - ✓ SI pentru predicții, care au ca scop predicția ieșirii pentru o intrare nouă, folosind un model învățat anterior;
 - ✓ SI pentru regresii, care au ca scop să estimeze forma unei funcții uni- sau multivariată, folosind un model învățat anterior;
 - ✓ SI pentru clasificare, care au ca scop clasificarea unui obiect în una sau mai multe categorii (clase) – cunoscute sau nu anterior, pe baza caracteristicilor (atributelor, proprietăților) lui;
 - ✓ SI pentru planificare, care au ca scop generarea unei succesiuni optime de acțiuni pentru efectuarea unei sarcini;
- în funcție de experiența acumulată în timpul învățării:
 - ✓ SI cu învățare supervizată;
 - ✓ SI cu învățare nesupervizată;
 - ✓ SI cu învățare activă;
 - ✓ SI cu învățare cu întărire;
- în funcție de modelul învățat (algoritmul de învățare):
 - ✓ Arbori de decizie;
 - ✓ Rețele neuronale artificiale;
 - ✓ Algoritmi evolutivi;
 - ✓ Mașini cu suport vectorial;
 - ✓ Modele Markov ascunse.

Una dintre principalele direcții – *învățarea supervizată* – este o clasă de metode și tehnici de învățare automată prin care se induce/determină o funcție de evaluare (șablon) dintr-un set de date (exemple, instanțe, cazuri) cu valori deja atribuite. Cu alte cuvinte, este formalizarea învățării din exemple, adică se folosește un set de instanțe rezolvate ale problemei pentru a antrena sistemul în vederea rezolvării unor instanțe noi. Aceste instanțe rezolvate se numesc *instanțe de antrenament*. Formal, setul de instanțe de antrenament este o mulțime de perechi atribut-valoare $(x, F(x))$, unde x este instanța, iar $F(x)$ clasa căreia îi aparține instanța respectivă.

Ideea constă în faptul că procesul de învățare decurge în două etape: *antrenare* și *testare*. Programul învață dintr-un set de exemple cu etichete cunoscute din setul de antrenament pentru a identifica exemple neetichetate din setul de teste cu o acuratețe cât mai mare. Un algoritm de învățare supervizată analizează datele de antrenament și produce o funcție-șablon, numită *funcție de clasificare* – în cazul în care datele de ieșire sunt dintr-o mulțime discretă (predefinită), sau *funcție de regresie* – dacă datele de ieșire aparțin unei mulțimi continue. Funcția produsă clasifică corect instanțele-exemplu, iar pentru un x pentru care nu se cunoaște $F(x)$ ea trebuie să propună/prezică o aproximare cât mai corectă a valorii $F(x)$. Acest lucru cere de la algoritmul de învățare să generalizeze de la datele de antrenament la orice date de intrare noi din mulțimea indicată în problemă.

Mașinile cu suport vectorial (*Support vector machines, SVM*) sunt bazate pe teoria învățării statistice. Ideea de bază este de a mapa datele originale într-un spațiu inițial printr-o funcție neliniară și de a construi un hiperplan optim într-un spațiu nou cu mai multe dimensiuni. Acești algoritmi pot fi utilizați atât pentru clasificare, cât și pentru regresie. Mașinile cu suport vectorial sunt o metodă de învățare supervizată *state-of-the-art*. Algoritmii SVM sunt folosiți pe larg în bioinformatică datorită acurateții înalte și capacității de a se descurca cu date multidimensionale, cum ar fi descrierea genelor, flexibilitatea modelării diverselor surse de date [2].

În prezent, aplicațiile practice ale SVM sunt numeroase. Cele mai cunoscute aplicații reale în care au fost implementate SVM sunt: recunoașterea scrisului de mână, clasificarea textului și a hipertextului, identificarea sexului din imagini și algoritmi de ordonare a paginilor după relevanță [3].

Analiza predictivă și construirea unui model de predicție al performanțelor studenților

Analiza predictivă este un domeniu al analizei statistice, care se ocupă cu extragerea informațiilor din date și folosirea acestora pentru a prezice comportamentul unor fenomene din viitor. Analiza predictivă cuprinde tehnici statistice din învățarea automată, data mining, teoria jocurilor, prin care se caută descoperirea unor relații și șabloane în volum mare de date, care pot fi folosite pentru a prezice comportamente și evenimente. Spre deosebire de alte tehnici, analiza predictivă este anticipativă, folosind evenimentele precedente (curente și din trecut) pentru a anticipa evenimentele din viitor.

Performanța studenților este unul dintre cele mai importante scopuri ale instituțiilor de învățământ, însă acest lucru nu este deloc simplu. Studenții, având note joase sau chiar negative, sunt nevoiți să repete anul academic sau, în cel mai rău caz, sunt exmatriculați. Aceste rezultate duc la consecințe neplăcute atât pentru instituția de învățământ, cât și pentru studenții acesteia, implicând cheltuieli financiare semnificative. Rezultatele nesatisfăcătoare ale studenților pot fi consecința mai multor factori, printre care: scăderea motivației acestora, probleme sociale (în familie, în relație), consumul relativ excesiv de alcool etc. Din această perspectivă, *ne propunem, prin aplicarea tehnicilor de învățare automată, elaborarea unui sistem inteligent care, în baza modelului de predicție, prin găsirea unor asocieri între performanțele și obiceiurile studenților, să poată fi utilizat la prevenirea eșecurilor posibile în condiții de incertitudine a stării universitare a studenților.*

Tehnicile de achiziționare automată a cunoștințelor necesită un set mare de date. Cu cât sunt mai puține date, cu atât scade posibilitatea oricărui model de învățare supervizată de a face vreo predicție. Învățarea supervizată face predicții în baza unor fapte precedente; astfel, dacă date sunt puține, rezultă că modelul nu are din ce învăța. Învățarea trebuie să ducă la formularea de suficiente „reguli”, astfel încât acestea să permită rezolvarea unor probleme dintr-un spațiu mai larg decât cel pe baza căruia s-a făcut învățarea. Adică, învățarea trebuie să îmbunătățească performanța unui sistem nu doar în rezolvarea repetată a aceluiași set de probleme, ci și în rezolvarea unor probleme noi. Acest lucru presupune o generalizare a unei metode de rezolvare pentru a acoperi un număr cât mai mare de instanțe posibile, dar și păstrarea unei specializări suficiente pentru a fi identificate corect instanțele acceptate. Aceasta se poate face fie inductiv, generalizând o problemă plecând de la un set de exemple, fie deductiv, plecând de la o bază de cunoștințe suficiente asupra universului problemei și extrăgând date și reguli esențiale. Pentru a putea face acest lucru, un algoritm de învățare trebuie să fie capabil să selecteze acele elemente semnificative pentru rezolvarea unei instanțe viitoare a problemei.

Procesul de creare a modelului predictiv include următorii pași:

1. **Definirea proiectului** – include formularea obiectivelor și rezultatelor așteptate pentru proiect și traducerea acestora în obiective și sarcini din analiza predictivă.
2. **Explorare** – analizarea sursei de date și determinarea celui mai potrivit set de date și a celei mai bune abordări pentru a construi modelul.
3. **Pregătirea datelor** – selectarea, extragerea și transformarea datelor ce vor fi folosite pentru crearea modelului.
4. **Construirea modelului** – crearea, testarea și validarea modelului. Evaluarea acestuia în baza unui set de date de antrenament și satisfacerea cerințelor proiectului.
5. **Exploatare/utilizare** – aplicarea regulilor (rezultatele/deciziile efectuate în model) formate pentru etichetarea unor noi date.
6. **Gestiunea modelului** – gestionarea și întreținerea modelelor construite pentru a îmbunătăți performanța, a promova reutilizabilitatea și a minimiza activitățile redundante.

Exemplele (instanțele) folosite pentru acest model: un total de 1044 de studenți, dintre care 395 au frecventat cursul „Inteligența artificială” și 649 – cursul „Limba engleză”, și s-a considerat că toți sunt consumatori de alcool (cel puțin o dată pe săptămână). Setul de date a fost preluat din baza de date de pe site-ul UCI Machine Learning (<https://archive.ics.uci.edu/ml/datasets/student+performance>). Instanțele sunt reprezentate printr-un număr fix (31) de attribute/caracteristici, care evaluează fiecare student și fiecare atribut putând avea un număr limitat de valori. Funcția obiectiv ia valori de tip discret, iar arborele de decizie reprezintă o disjuncție de mai multe conjuncții.

Datele din baza de cunoștințe mai sunt caracterizate și de acuratețea și calitatea învățării (metode de evaluare, măsurări de performanță). O componentă esențială a unui algoritm de învățare este metoda de verificare, o metodă capabilă să confirme dacă generalizările făcute sau regulile deduse se apropie mai mult de soluția ideală decât starea anterioară a sistemului.

Fie F funcția pe care dorim s-o aproximăm, S un set de instanțe de antrenament $(x, F(x))$, iar H funcția indusă prin învățare. Este sau nu $H(x)$ aproape de $F(x)$, pentru orice x din spațiul instanțelor posibile? Argumentul standard este acela că H nu poate fi prea departe de F , deoarece H clasifică corect instanțele de antrenament, deci probabil și pe celelalte. Deci, H este *probabil aproximativ corect* (PAC). Un concept este învățabil PAC dacă există un algoritm eficient care are o probabilitate mare de a găsi o aproximare a conceptului *probabil aproximativ corectă*. Termenul „învățabil PAC” a fost introdus de Valiant (1984).

Presupunerea aflată la baza acestei justificări este aceea că instanțele de antrenament și instanțele de testare sunt uniform distribuite în spațiul instanțelor posibile ale problemei. Acest lucru este fundamental în justificarea oricărui rezultat al unei învățări supervizate, prezentată grafic în Figura 1.

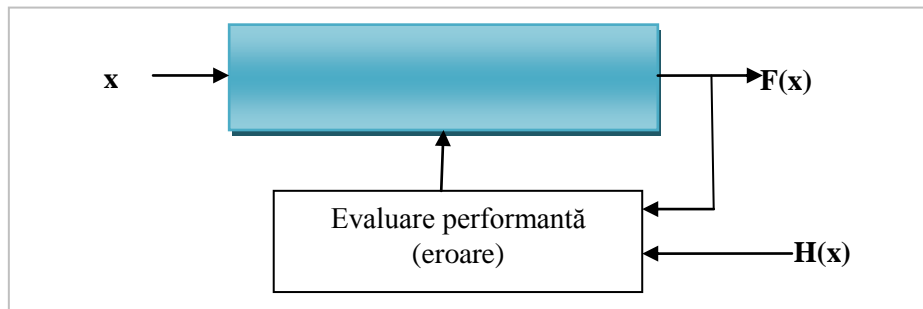


Fig.1. Învățarea supervizată.

După construirea mai multor modele de regresie diferite, există o multitudine de criterii prin care pot fi evaluate și comparate. În continuare, vom menționa cele mai importante măsurări de evaluare folosite a modelului predictiv de regresie elaborat. Vom nota cu p vectorul valorilor prezise, cu a – vectorul valorilor așteptate, cu n – numărul de instanțe.

Pentru măsurarea erorii medii pătrate (EMP), care poate fi comparată doar între modelele de regresie, a căror eroare măsurată are aceeași unitate de măsură, folosim [4]:

$$EMP = \frac{\sum_{i=1}^n (p_i - a_i)^2}{n} \quad [1]$$

Pentru a compara două modele pentru care erorile sunt măsurate în unități diferite folosim eroarea relativă pătrată (ERP):

$$ERP = \frac{\sum_{i=1}^n (p_i - a_i)^2}{\sum_{i=1}^n (a_i - a_i)^2} \quad [2]$$

Eroarea medie absolută (EMA) are aceeași unitate ca și datele originale și poate fi comparată numai între modelele ale căror erori sunt măsurate în aceeași unități:

$$EMA = \frac{\sum_{i=1}^n |p_i - a_i|}{n} \quad [3]$$

La fel ca ERP, eroarea relativă absolută (ERA) poate fi comparată între modele ale căror erori sunt măsurate în unități diferite.

$$ERA = \frac{\sum_{i=1}^n |p_i - a_i|}{\sum_{i=1}^n |a_i - a_i|} \quad [4]$$

Selectarea tehnologiilor informatice pentru realizarea aplicației

Metodologia folosită în realizarea informatică a aplicației elaborate a fost mixtă, bazându-se pe tehnologii software de ultimă oră. Decizia noastră a avut ca prim argument utilizarea instrumentelor dedicate implementării algoritmilor din *machine learning* și analiza datelor, fără necesitate de mari resurse. În acest sens, astăzi una dintre cele mai populare tehnologii este *Jupyter Notebook* (inițial fiind numit *IPython*, mai apoi redenumit în *Jupyter*). Ca limbaj de programare a fost ales *Python* și distribuția *Anaconda*, ca tehnologii inovatoare ce s-au remarcat prin simplitate, portabilitate și robustețe și au început să fie utilizate pentru producerea și dezvoltarea unui număr mare de aplicații software, sub diverse forme, implementate în diferite domenii. Limbajul *Python* poate fi folosit pentru programarea *obiect-orientată*, *funcțională* și *procedurală*. Ca web framework pentru dezvoltarea aplicației s-a folosit *Flask*.

Rezultate finale obținute

În urma pregătirii și analizei setului de date, cu ajutorul *IPython API* citim, executăm și transformăm *notebook*-urile *Jupyter* în *HTML*. Din pagina de start a aplicației (Fig.2) suntem invitați să selectăm și să încărcăm fișierul cu date pentru care se va efectua predicția.



Fig.2. Pagina de start a aplicației.

În urma studiului comparativ, se poate afirma cu siguranță că așa tip de aplicație nu a fost până în prezent elaborată, deoarece cel puțin un argument este faptul că tehnologia API este relativ nouă (motiv ce explică lipsa documentației de exploatare). Aplicația elaborată, combinând o interfață web cu un model de *machine learning*, facilitează enorm utilizarea modelului de predicție de către diverși utilizatori.

În urma încărcării și exploatării setului de date, au fost obținute rezultate/preziceri destul de pozitive (Fig.3, 4), folosind ca caracteristici pentru predicție performanțele (notele) anului 1 (A1) și ale anului 2 (A2).

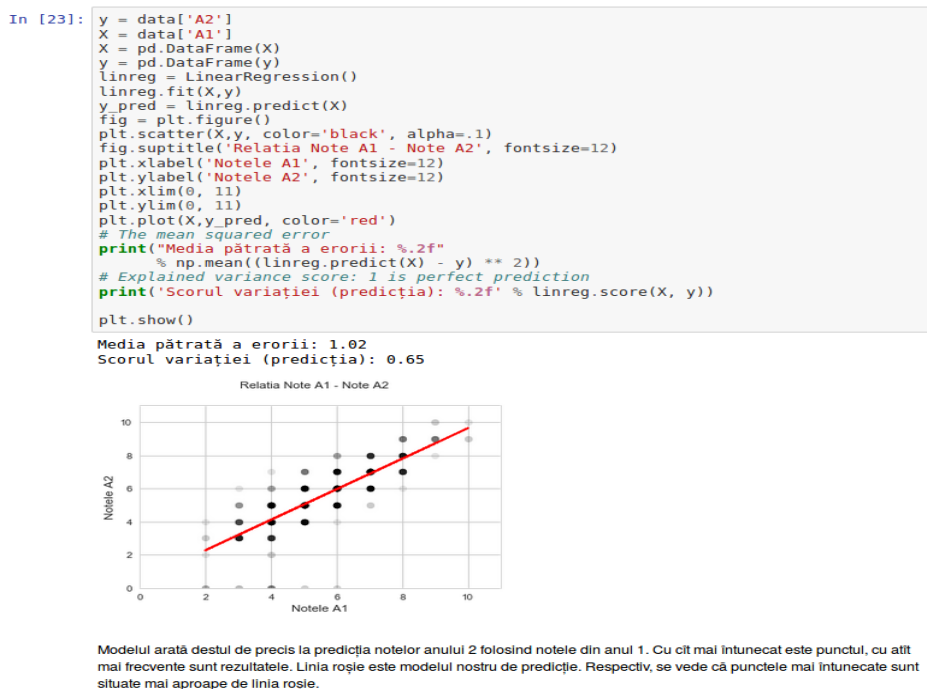


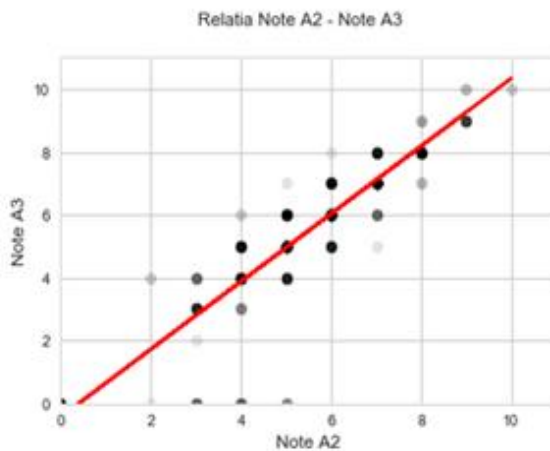
Fig.3. Rezultatele predicției prin **regresia liniară**, folosind caracteristicile **A1** și **A2**.

Analizând dacă această regresie funcționează și pentru predicția performanțelor (notelor) anului 3 (A3), obținem următoarele rezultate:

```
In [24]: y = data['A3']
X = data['A2']
X = pd.DataFrame(X)
y = pd.DataFrame(y)
linreg = LinearRegression()
linreg.fit(X,y)
y_pred = linreg.predict(X)
fig = plt.figure()
plt.scatter(X,y, color='black', alpha=.1)
fig.suptitle('Relatia Note A2 - Note A3', fontsize=12)
plt.xlabel('Note A2', fontsize=12)
plt.ylabel('Note A3', fontsize=12)
plt.xlim(0, 11)
plt.ylim(0, 11)
plt.plot(X,y_pred, color='red')
# The mean squared error
print("Media pătrată a erorii: %.2f"
      % np.mean((linreg.predict(X) - y) ** 2))
# Explained variance score: 1 is perfect prediction
print('Scorul variației (predicția): %.2f' % linreg.score(X, y))

plt.show()
```

Media pătrată a erorii: 1.03
Scorul variației (predicția): 0.76



Lucrează chiar mai bine pentru notele anului 3. Aceasta este din cauza corelației dintre A2 și A3 (~ 87%) mai mari decât celelalte dintre A1 și A2 (~ 80%).

Concluzii

Nu am avut posibilitatea de a găsi careva corelații semnificative între obiceiurile studenților și notele acestora la universitate. Aceasta poate fi datorită dimensiunii prea mici a seturilor de date, sau poate datorită vreunei variabile care noi nu o considerăm în aceste date. Dar se pare că există un model similar pentru fiecare dintre perioadele de note (A1, A2 și A3), arătând o oarecare continuitate în rezultatele universitare. Însă mai trebuie careva cercetare în acest domeniu pentru a putea avea rezultate mai clare.

Fig.4. Rezultatele predicției prin regresia liniară, folosind caracteristicile A2 și A3.

Concluzii

Există numeroase probleme în care nu s-au exploatat îndeajuns metodele și tehnicile de învățare automată existente. Una dintre aceste probleme este previziunea performanțelor studenților. Deși s-au întreprins încercări de a utiliza învățarea automată pentru astfel de predicție, această problemă rămâne în continuare un subiect important de cercetare. Studiul de caz din acest articol a demonstrat cu succes că pot fi aplicate tehnici de învățare automată pentru astfel de probleme. Algoritmii de învățare folosesc date precedente ca să construiască modele matematice care să fie capabile să facă predicții de performanță universitară.

Procesul de predicție este unul complex, care necesită un volum mare de informații, pe care specialistul nu poate să le dețină în totalitate. Elaborarea aplicației informatice prin implementarea tehnicilor de învățare automată permite realizarea de experimente pentru studierea problemei de predicție a performanțelor studenților.

În final, conchidem că utilizarea tehnologiilor de *machine learning* și, în special, a modelelor de predicție poate facilita enorm rezolvarea anumitor probleme (cum ar fi intervenția la timp a unor factori de decizie pentru a preîntâmpina dificultățile cu care se confruntă unii studenți în ce privește continuarea studiilor), ceea ce va duce la rezolvarea lor la timp.

Referințe:

1. RUSSEL, S.J., NORVIG, P. *Artificial Intelligence: A Modern Approach*, Prentice Hall. New Jersey, 1995, p.4-8, 567-570.
2. BEN-HUR, A., WESTON, J. *A User's Guide to Support Vector Machines*. Colorado State University, Colorado, NEC Labs America, Princeton, p.1-3.
3. KAR, P. *Support Vector Machines and their Applications*. Indian Institute of Technology Kanpur, 2009, p.56-70.
4. SAYAD, S. *Model Evaluation – Regression*. University of Toronto, 2010, http://chem-eng.utoronto.ca/~datamining/dmc/model_evaluation_r.htm

Prezentat la 13.07.2017