

Botond BENEDEK,
Mihai-Constantin AVORNICULUI

Faculty of Economics and Business Administration, Babeş-Bolyai University of Cluj-Napoca

DESIGNING AN EVENT EXTRACTION SYSTEM

Case
Study

Keywords

*Event Extraction System,
Designing Methods,
UML*

JEL Classification

M15

Abstract

In the Internet world, the amount of information available reaches very high quotas. In order to find specific information, some tools were created that automatically scroll through the existing web pages and update their databases with the latest information on the Internet. In order to systematize the search and achieve a result in a concrete form, another step is needed for processing the information returned by the search engine and generating the response in a more organized form. Centralizing events of a certain type is useful first of all for creating a news service. Through this system we are pursuing a knowledge - events from the Internet documents - extraction system. The system will recognize events of a certain type (weather, sports, politics, text data mining, etc.) depending on how it will be trained (the concept it has in the dictionary). These events can be provided to the user, or it can also extract the context in which the event occurred, to indicate the initial form in which the event was embedded.

INTRODUCERE

Dimensiunea Internetului și conținutul său nestructurat și dinamic, precum și natura sa multilingvă, fac din extragerea cunoștințelor utile o problemă de cercetare provocatoare. Web-ul generează o cantitate mare de date în alte formate, care conțin informații valoroase. Astfel ar fi informația log-urilor serverelor Web despre șabloanele accesate de utilizatori care poate fi folosită pentru personalizarea informației sau pentru îmbunătățirea design-ului paginii Web.

Pentru a sistematiza căutarea pe Internet și a obține un rezultat într-o formă concretă este utilă o etapă în care se prelucerează informațiile întoarse de către motorul de căutare și se generează răspunsul într-o formă mai organizată (Avornicului, 2010).

Centralizarea evenimentelor de un anumit tip este utilă în primul rând pentru realizarea unor servicii de știri. Aceste servicii trebuie să ofere informații actualizate, dacă se poate chiar în timp real, despre un anumit tip de evenimente (Han & Kamber, 2006), (Lin, et al., 2008).

Să luăm de exemplu evenimentele sportive. Un utilizator al acestui serviciu de știri poate să dorească să afle ce jocuri sportive se desfășoară într-o anumită regiune de pe glob (Oraș, Țară, etc.) dintr-o perioadă anume: în acel moment, cu o zi în urmă, sau săptămâna viitoare, etc. Toate acestea informații trebuie obținute din informațiile deja centralizate. Se pot obține astfel date despre un anumit eveniment din mai multe surse. Deasemenea sursele (siturile web din care se extrag informațiile) se pot completa unele pe altele în conținutul de informații despre anumite evenimente.

Acest sistem urmărește extragerea de evenimente din documentele HTML. Sistemul va recunoaște evenimentele de un anumit tip în funcție de modul în care va fi antrenat.

Acest sistem urmărește să fie de un real ajutor în cazul căutării informației din mai multe pagini de web, când se dorește o centralizare a evenimentelor de un anumit tip.

De exemplu, dacă se dorește aflarea informațiilor despre vreme din mai multe pagini de web (yahoo, google, situl CNN, etc) pentru o anumită zonă și o anumită perioadă atunci sistemul se poate dovedi util (Avornicului, 2014).

Pentru fiecare tip de eveniment (vreme, sport, etc.) există câte un dicționar de concepte. Construirea unui dicționar de concepte este o sarcină care necesită multă muncă. Încercăm să realizăm un sistem care să fie extensibil, astfel încât să permită, la adăugarea unui nou dicționar, extragerea de evenimente specifice acestuia. La ora actuală prezintă un domeniu de interes construirea dicționarelor de concepte în mod automat, utilizând algoritmi de învățare.

MODELUL CONCEPTUAL

Pentru a putea realiza ceea ce ne propunem, sistemul (componentele) trebuie să se ocupe, în general, de trei aspecte (Avornicului, 2009):

- aducerea documentelor HTML de pe Internet și salvarea lor într-o bază de date pentru a putea fi prelucrate ulterior (a)
- prelucrarea documentelor și obținerea informațiilor necesare (b)
- oferirea utilizatorilor un mod de acces asupra informațiilor colectate (c)

a) Pentru a se putea realiza acest aspect sistemul trebuie să primească o serie de adrese web pe care trebuie să le parcurgă. Această listă de adrese poate să fie dependentă de tipul de evenimente pe care dorim să le identificăm. Aplicația va parcurge recursiv această listă de adrese și va salva în baza de date locală toate documentele întâlnite.

b) Prelucrarea documentelor și extragerea informațiilor necesare are la bază un dicționar de concepte care descriu tipurile de evenimente. Acest dicționar de concepte trebuie să fie flexibil pentru a putea identifica diferite variații de tiparuri de evenimente ce pot apărea în documente. Toate evenimentele identificate vor fi apoi stocate într-o bază de date pentru o referire ulterioară.

c) În final trebuie oferit utilizatorilor accesul la informațiile extrase. De exemplu, pentru evenimente din sport, se poate prezenta utilizatorilor o listă cu evenimentele sportive ordonate după data la care vor avea sau au avut loc. Se va permite deasemenea și o căutare a evenimentelor după diferite criterii. Astfel se va putea face un acces rapid la informația dorită, fără a fi necesară parcurgerea evenimentelor zi cu zi. Figura 1 prezintă relevant aspectele sistemului.

Utilizatorii vor trebui să acceseze aplicația prin intermediul Internetului. Astfel, dacă utilizatorii sunt umani, se vor oferi documente HTML cu informația dorită, iar dacă utilizatorii sunt programe atunci informația va putea fi transmisă sub o formă generică, document XML de exemplu, pentru a putea fi prelucrată cu ușurință (Markov & Larose, 2007).

CERINȚE FUNCȚIONALE

Fiind vorba de un serviciu, accesul la baza de date cu evenimentele colectate va putea fi disponibil oricărui utilizator prin intermediul Internetului. Bineînțeles că se poate restrânge această arie la utilizatorii înregistrați să acceseze serviciul respectiv, de exemplu. Deci virtual numărul utilizatorilor care vor folosi acest serviciu este nelimitat (Kovács, 2012). Trebuie să se permită

deasemenea accesul simultan al mai multor utilizatori la baza de date cu evenimentele extrase.

Având în vedere că toate procesele consumatoare de timp și putere de calcul se vor desfășura offline (transparent pentru utilizatorii finali), timpul de răspuns al sistemului la cererile utilizatorilor trebuie să fie redus (de ordinul milisecundelor/secundelor).

Sistemul trebuie să fie capabil să proceseze un volum de date suficient, astfel încât, de exemplu în cazul evenimentelor sportive, să se poată face o actualizare a acestora din oră în oră (Lin, 2007).

Scopul sistemului este să ofere un răspuns cât mai complet la interogările utilizatorilor (Zaki, 2000). Teoretic, acest răspuns este corect și complet, însă în practică acest lucru este imposibil de realizat de un sistem, fără a se cunoaște date despre formatul cunoștințelor din text (adică modul de prezentare a informației în pagina HTML). Din acest motiv indicatorii de eficiență ai sistemului vor varia în funcție de sursele din care se extrage informația. Ceea ce ne propunem este să păstrăm acești indicatori între anumite limite astfel încât aplicația noastră să aibă o utilitate reală.

Indicatorii pe care sunt luate în considerare sunt:

- **acoperirea** - reprezintă raportul: numărul de rezultate corecte găsite / numărul de rezultate corecte existente;
- **precizia** – reprezintă raportul: numărul de rezultate corecte găsite / numărul total de rezultate găsite.

În ceea ce privește specificațiile interfeței cu utilizatorul, ceea ce se poate spune la momentul actual este că se permite utilizatorului introducerea de interogări în diverse formate și, în același timp, controlarea modului de afișare a rezultatelor (Sommerville, 2010). Aspectul exact al interfeței cu utilizatorul va fi stabilit în faza de testare cu scopul de a determina un mod de lucru cât mai ergonomic (Han & Kamber, 2006).

ASPECTE LEGATE DE PROIECTARE

Analiza și proiectarea unui sistem informatic reprezintă una dintre cele mai importante etape ale realizării unui sistem. Dacă un sistem nu este bine gândit de la bun început, este foarte posibil ca, la implementarea sa sau la exploatarea sa, să constatăm că datele pe care ne așteptăm să le furnizeze respectivul sistem nu pot fi obținute, prelucrările durează mai mult decât ne-am fi așteptat, situațiile conțin date eronate care nu ne permit să luăm o decizie, sau că, sistemul, în loc să simplifice activitatea mai mult o îngreunează. Puțini sunt aceia care realizează că de fapt, toate aceste neajunsuri nu sunt datorate nici programatorilor și nici implementatorilor sistemului informatic ci chiar celor care au „gândit” respectivul sistem. De multe ori, etapa de analiză și

modelare este „sărită” considerându-se ca fiind neproductivă și consumatoare de timp. Într-adevăr, dacă privim activitatea de analiză și proiectare având drept finalitate o serie de documente cu scheme și descrieri cu privire la cum cum ar arăta sistemul în realitate, putem considera că atât timpul alocat cât și banii investiți sunt niște pierderi inutile. Dacă însă luăm în calcul faptul că după etapa de analiză și proiectare echipele de dezvoltatori au un plan de abordare a problemei, multe dintre module sunt demarate în baza unei analize de optimizare a termenului de finalizare iar posibilitatea de a relua anumite aspecte este mult diminuată vom începe să acordăm atenția cuvenită acestei prime etape (Jing & McKeown, 1999). Ceea ce trebuie remarcat este faptul că, adevărata returnare a investiției în analiză și proiectare nu are loc în momentul dezvoltării ci în ultimele etape: mentenanța și exploatarea. Un proces de analiză și proiectare derulat după toate rigorile metodologice conferă o temelie robustă a ceea ce urmează să fie făcut și asigură continuitate sistemului în momentul în care acesta va trebui extins sau îmbunătățit. Fără o documentație atent întocmită, respectivul sistem va avea serios periclitată acele etape care au un impact decisiv asupra clientului. Dacă va apărea o problemă, aceasta nu va putea fi soluționată în mod optim fără o cunoaștere în detaliu a arhitecturii sistemului. Se va încerca rezolvarea unui aspect semnalat de către utilizatori dar modificările introduse pentru rezolvarea aceluși aspect se vor propaga în lanț și vor pune echipa în fața unor situații uneori imposibile (Avornicului, 2010).

Trebuie menționat că, proiectarea unui sistem este influențată de o serie de factori. De multe ori, proiectarea unui sistem depinde de experiența proiectantului sau a echipei de proiectanți. Există chiar un moto care afirmă că, cel de-al doilea sistem conceput de un proiectant pentru un anumit domeniu este întotdeauna mai bun decât primul. Afirmatia are un suport real dat fiind faptul că, echipa de proiectanți deja cunoaște anumite aspecte și evită greșelile făcute anterior. Un alt aspect ce trebuie luat în considerare este acela că, anumite metodologii de proiectare se pretează mai bine unui anumit tip de problemă, echipa de proiectanți reușind să obțină rezultate bune utilizând mult mai puține resurse. Deși ridicate la un nivel de formalizare suficient de mare, anumite metodologii sunt preferate altora din motive legate strict de derularea unui anumit proiect. De exemplu, este dorit ca o echipă de proiectanți să recurgă la utilizarea UML (Avornicului, 2014). Totuși, dată fiind experiența pe care componenții respectivei echipe o au în utilizarea metodelor sistemice și intervalul de timp extrem de redus în care trebuie parcursă etapa de proiectare, este posibil ca aceștia să recurgă la o metodă cunoscută lor. Echipa respectivă ar fi trebuit să facă un important salt calitativ dublat de un effort susținut de

documentare pentru a putea prelua din plin avantajele noilor concepte care stau la baza UML cum ar fi: polimorfismul, evenimentele, interfața, mesajele, componentă comportamentală etc. Pentru aceeași problemă sunt oferite soluții diferite. Diferența este făcută pe considerente care țin de aspecte precum (Avornicului, 2009):

- cât de flexibilă este soluția adoptată
- cât de rapidă este în prelucrarea datelor
- cât de mulți utilizatori concurenți suportă respectiva soluție
- cât de costisitoare sunt dispozitivele hard necesare rulării unei anumite soluții
- cât de extensibilă este respectiva soluție
- care este efortul pe care-l necesită pentru a o adapta unor cerințe specifice
- cât de mare este redundanța datelor sau care este gradul de normalizare etc..

De multe ori, o anumită soluție reflectă personalitatea celui care proiectează, felul în care respectiva persoană sau grup de persoane percepe realitatea. Acest aspect i-a determinat pe mulți să spună că modelarea ca și programarea este o artă (Chen & Wei, 2002). De exemplu unii proiectanți, atunci când au modelat conceptul de partener au văzut necesar să menționeze dacă acesta este furnizor, client sau ambele (Cutting, et al., 1992).

DIAGRAMA USE CASE PENTRU SISTEMUL DE EXTRAGERE DE EVENIMENTE

Acestea au rolul de a reprezenta într-o formă grafică funcționalitățile pe care trebuie să le îndeplinească sistemul informatic în faza finală. De aceea modelul realizat de ele împreună cu descrierea succintă a fiecărui caz de utilizare determinat se numește model al cerințelor.

Cazurile de utilizare arată comportamentul sistemului sau a unei părți din sistem și este o descriere a unui set de secvențe de acțiuni, incluzând variante pe care un sistem le execută pentru a produce un rezultat observabil pentru un actor (Raffai, 2005), (Sommerville, 2010).

Diagramele acestea sunt formate din două categorii de entități (actori și cazuri de utilizare) și relațiile dintre acestea (Figura 2).

Ele implică interacțiunea dintre actori și sistem, iar actorii pot fi oameni sau sisteme automate.

Comportamentul unui caz de utilizare se poate specifica prin descrierea fluxului de evenimente. Acestea trebuie să indice cum și când începe și se încheie cazul de utilizare, când interacționează cu actorii, ce obiecte sunt schimbate, fluxul de bază și cel alternativ al comportamentului.

În continuare prezentăm detaliat funcția fiecărui modul (Figura 3):

- **Baza de Date cu Adrese de Pagini:** Păstrează lista paginilor Web pentru un anumit domeniu, pe care se va face căutarea de evenimente

- **Baza de Date de Pagini:** Păstrează paginile Web, indexate și categorizate după domenii

- **Baza de Date cu Evenimente:** Păstrează evenimentele identificate, indexate după domenii, timp, locații și alte criterii

- **Baza de Date cu Dictionarul de Concepte:** Păstrează conceptele care vor fi căutate, indexate după domenii

- **Dezvărcare pagini:** Aduce periodic pagini Web de pe Internet de la adresele prezente în Baza de Date cu Adrese de Pagini, și le stochează în Baza de Date de Pagini. Acest modul poate aduce în mod recursiv și paginile indicate de link-urile din pagina curentă, de pe același server, până la o anumită adâncime de aducere.

- **Segmentare:** Realizează împărțirea textului din paginile Web din Baza de Date cu Pagini în segmente

- **POST:** Acest modul are funcția de a prelua segmentele și a le adnota cu părțile de vorbire corespunzătoare cuvintelor din segment.

- **Extragere evenimente:** Extrage evenimente din segmentele adnotate, primite de la modulul POST, pe baza conceptelor din Dictionarul de Concepte și a WordNet-ului, depunându-le în Baza de Date cu Evenimente

- **Creare Dictionar Concepte:** Construiește conceptele specifice unui anumit domeniu, prin algoritmi de învățare, apoi le adaugă în Baza de Date cu Concepte.

- **WordNet:** Furnizează relații între cuvinte de natură semantică

- **Interfața cu aplicația client:** Primește cereri de căutare evenimente, pe care le obține din Baza de Date cu Evenimente

- **Interfața client:** Asigură interfața cu utilizatorul, introducerea de query-uri de evenimente, setarea de opțiuni, afisarea rezultatelor după diverse criterii

DIAGRAMA DE CLASĂ

Ele fac parte din categoria diagramelor statice și descriu structura internă a sistemului informatic prin identificarea claselor, a atributelor și a operațiilor acestora și a relațiilor dintre clase. Construcția lor are loc în faza de elaborare a sistemului informatic, fiind cele mai importante în această fază. Selectarea claselor necesită anumite deprinderi. O abordare metodică a determinării claselor ce modelează soluția unei probleme informatice este aceea de a extrage substantivele esențiale din documentul de specificație. O metodă mult mai rapidă este aceea de a extrage toate substantivele care apar în diagrama de cazuri de utilizare (atât în denumirile actorilor cât și a cazurilor de utilizare).

Clasele nu sunt elemente izolate ale sistemului. Astfel trebuie să identificăm relațiile între clase (Figura 4).

CONCLUZII

Cercetările teoretice și practice privind proiectarea și realizarea sistemelor de extragere de evenimente au început să se focalizeze tot mai mult pe utilizarea limbajului UML în modelarea acestor sisteme.

Tehnologia pentru extragerea evenimentelor din pagini web și arhitectura sistemului informatic influențează semnificativ felul în care trebuie proiectat și implementat sistemul de extragere de evenimente.

În acest sens propunem utilizarea modelelor UML pentru definirea cerințelor și determinarea contextului informațional, astfel încât să poată fi creat un scenariu cât mai complet de ceea ce va trebui să acopere sistemul de extragere de evenimente.

De asemenea propunem utilizarea în faza de proiectare, a tehnicii prototipizării prin care să se permită comunicarea directă dintre utilizator, sistem și proiectantul acestuia. Considerăm, că o determinare bine documentată și cât mai apropiată de ceea ce se dorește de la viitorul sistem, încă din faza de concepere, asigură succesul elaborării, implementării și utilizării sistemului.

Recomandăm abordarea proiectării și implementării sistemului de extragere de evenimente prin fluxul iterativ al procesului unificat, ceea ce va determina ca versiunea executabilă de bază să fie cât mai stabilă și să acopere cât mai bine cerințele utilizatorilor.

BIBLIOGRAFIE

- [1] Avornicului, M., 2009. Data mining în Internet.
- [2] Avornicului, M., 2010. *Informatikai rendszerek tervezése és menedzsmentje*. Kolozsvár: ÁBEL Kiadó.
- [3] Avornicului, M., 2014. Considerations On Objective Methods For Developing Applied Event Extraction Systems. *SEA-Practical Application of Science*, pp. Volume II, Issue 2 (4), pp. 447-456.
- [4] Chen, G. & Wei, Q., 2002. Fuzzy association rules ant the extended mining algorithms. *Information Sciences*, p. pp. 201–228..
- [5] Cutting, D. R., Pedersen, J. O., Karger, D. & Turkey, J., 1992. *A cluster-based approach to browsing large document collections*. Copenhagen, s.n., p. pp. 318–329.
- [6] Han, J. & Kamber, M., 2006. *Data Mining: Second Edition Concepts and Techniques*. s.l.:Morgan Kaufman Publishers.
- [7] Jing, H. & McKeown, K. R., 1999. *The decomposition of human-written summary sentences*. Berkeley, s.n., p. pp. 129–136.
- [8] Kovács, G., 2012. Productivity improvement by lean manufacturing philosophy. *ADVANCED LOGISTIC SYSTEMS: THEORY AND PRACTICE*, p. pp. 9–16..
- [9] Lin, B., 2007. Web Data Mining – Exploring Hyperlinks, Contents an Usage Data. *Springer*.
- [10] Lin, T. Y., Xie, Y., Wasilewska, A. & Lian, C. J., 2008. *Data mining: Foundations and Practice*. *Springer*.
- [11] Markov, Z. & Larose, D., 2007. *Data Mining the Web*. s.l.:John Wiley & Sons.
- [12] Raffai, M., 2005. *UML 2 modellező nyelv*. s.l.:Palatia Nyomda és Kiadó.
- [13] Sommerville, I., 2010. *Software Engineering, 9th Edition*. s.l.:Addison-Wesley.
- [14] Zaki, M., 2000. *Parallel and distributed data mining: An introduction..* s.l.:Springer-Verlag.

ANEXE

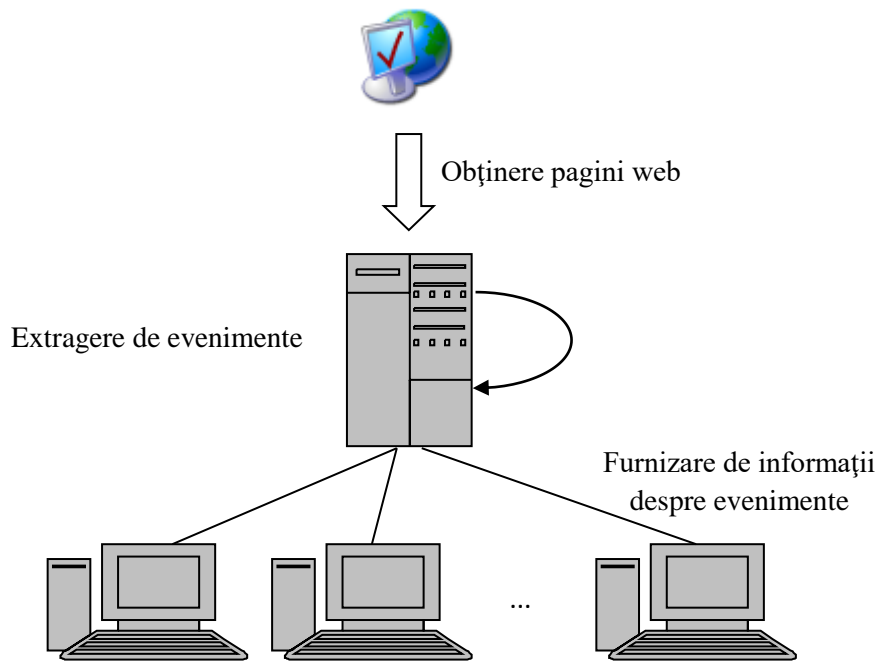


Figura 1 Aspectul sistemului

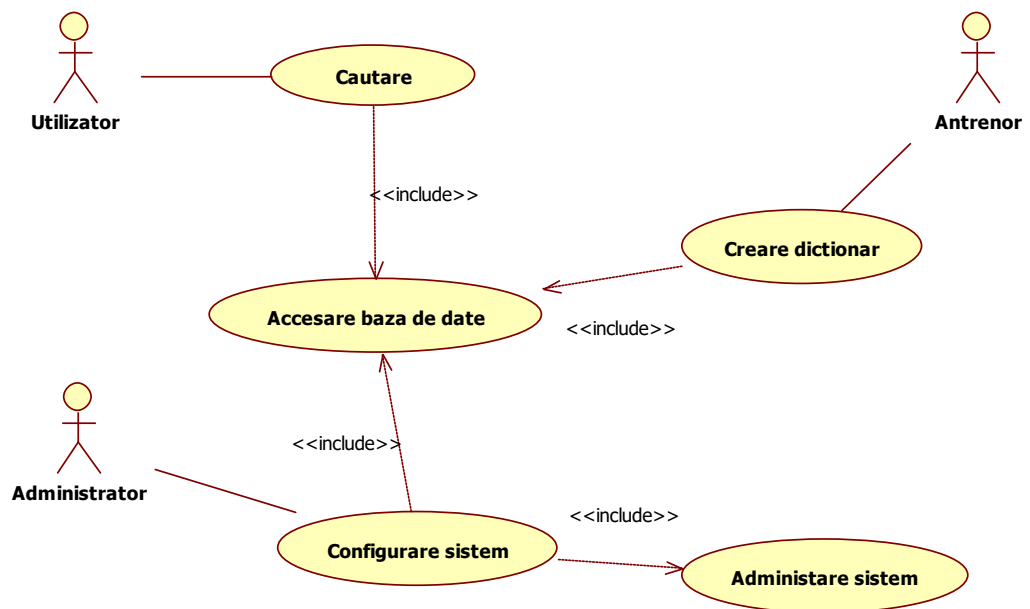


Figura 2. Diagrama "use case"



Figura 3. Dependente între module

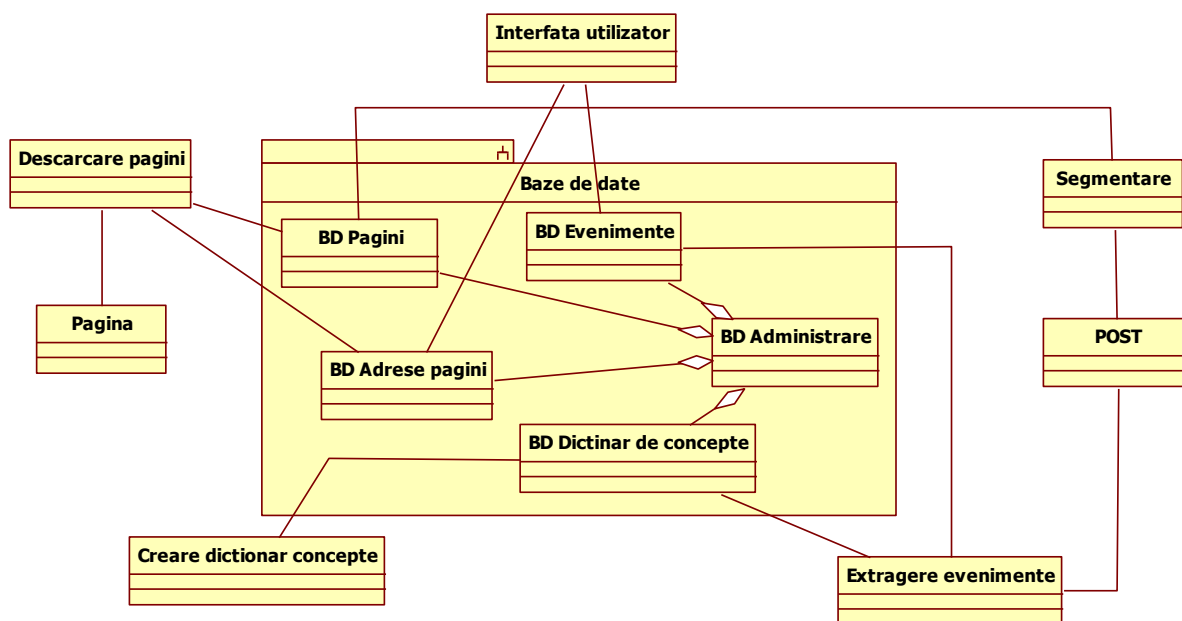


Figura 4. Relațiile între clase