



(12) CERERE DE BREVET DE INVENȚIE

(21) Nr. cerere: a 2019 00190

(22) Data de depozit: 26/03/2019

(41) Data publicării cererii:  
27/11/2020 BOPI nr. 11/2020

(71) Solicitant:  
• ZA CLOUD S.R.L., STR. GOVORA,  
NR. 16A, AP. 16, CLUJ-NAPOCA, CJ, RO

(72) Inventatori:  
• ȘTEȚCO EMIL IOAN, STR. GOVORA,  
NR. 16A, AP. 16, CLUJ-NAPOCA, CJ, RO;

• BARA GEORGE ANTONIU, STR. HAȚEG,  
NR. 28, BL. K2, AP. 24, CLUJ-NAPOCA, CJ,  
RO;

• SUCIU MIHAI ALEXANDRU,  
ALEEA PADIN, NR. 20, AP. 61,  
CLUJ-NAPOCA, CJ, RO

(54) METODĂ DE PUBLICITATE ONLINE DE TIP SEMANTIC

(57) Rezumat:

Invenția se referă la o metodă de publicitate online. Metoda conform invenției furnizează reclame pe baza contextului utilizatorului care este dedus prin analiza semnificativă a textului relevant din pagina web vizualizată de utilizator folosind metode de procesare a limbajului natural. Metadatele pe baza cărora se construiesc regulile de livrare sunt extrase din textul relevant folosind taxonomii standard, iar extragerea entităților și a sentimentului dedus, adică negativ, neutru sau pozitiv, se face cu algoritmi de învățare profundă.

Revendicări: 5  
Figuri: 2

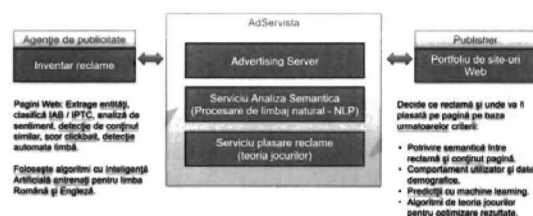


Fig. 1

Cu începere de la data publicării cererii de brevet, cererea asigură, în mod provizoriu, solicitantului, protecția conferită potrivit dispozițiilor art.32 din Legea nr.64/1991, cu excepția cazurilor în care cererea de brevet de invenție a fost respinsă, retrasă sau considerată ca fiind retrasă. Întinderea protecției conferite de cererea de brevet de invenție este determinată de revendicările conținute în cererea publicată în conformitate cu art.23 alin.(1) - (3).



8

BUCURIA DE STAT PENTRU INVENȚII ȘI MĂRCI	
Cerere de brevet de invenție	
Nr. a	2019 00190
Data depozit	26-03-2019

# Metodă de publicitate online de tip semantică

Invenția se referă la o metodă de publicitate online de tip semantic care se bazează pe algoritmi de inteligență artificială și de teoria jocurilor (Game Theory). Reclamele sunt servite utilizatorului pe baza contextului acestuia. Contextul este dedus pe baza datelor extrase din pagina online vizualizată de utilizator.

Internetul oferă un mediu versatil pentru publicitate. Un avantaj pentru publicitatea online față de mediile tradiționale (ex. afișe pe marginea drumului, ziar, etc.) este personalizarea conținutului livrat. Conținutul adaptat utilizatorului este mult mai eficient pentru furnizorul de publicitate deoarece reclamele furnizate sunt mai de interes pentru utilizator. Probabilitatea ca un utilizator care beneficiază de conținut personalizat să ignore reclame afișate în strânsă legătură cu pagina online vizualizată este mai mică față de un utilizator care primește conținut standard, fără corelarea conținutului vizualizat cu reclamele afișate. Astfel, în primul caz reclamele furnizate vor genera venituri mai mari pentru furnizorul de publicitate, iar în al doilea caz reclamele furnizate vor genera mai multe costuri pentru advertiser, costuri generate de vizualizări (impressions) fără acțiune directă asupra reclamelor. De asemenea reclamele personalizate pot reduce costurile pentru furnizor deoarece acestea sunt distribuite la un număr mai mic de utilizatori cu probabilitate foarte mare de a fi interesați de reclamele afișate care sunt legate exact de conținutul paginii online vizualizată în acel moment.

Abordarea semantică se bazează pe analiza automată a conținutului pus la dispoziție de către publisher (entitatea care integrează anunțurile publicitare în conținutul online publicat) unde va fi afișat anunțul digital și selectarea unui anunț cât mai relevant pentru pagina/aplicația respectivă. Acest algoritm este unul bazat pe analiză predictivă și implică folosirea de Procesare a Limbajului Natural (NLP), Învățare Automată (Machine Learning) și a Teoriei Jocurilor (Game Theory). Teoria Jocurilor este o ramură a matematicii aplicate care abordează problema comportamentului optim în jocurile cu două sau mai multe persoane, într-un cadru descris de un ansamblu de reguli precise care stabilesc posibilitățile de acțiune ale fiecărui jucător, precum și modul cum li se acordă acestora, în final, anumite valori.

Scopul abordării semantice este de a corela intenția utilizatorului, care se poate determina din corelarea conținutului pe care acesta îl vizualizează într-un anumit moment și / sau istoricul preferințelor sale cu un material de marketing relevant pentru intenție, livrat într-un mod personalizat și cu efect emoțional. Rezultatul abordării semantice poate fi cuantificat prin numărul de click-uri pe reclamă, numărul de vizualizări sau alte acțiuni ale utilizatorului pe pagină.

Prezenta invenție are ca scop furnizarea de reclame publicitare online utilizatorilor pe baza conținutului digital explorat de aceștia și protejarea mărcilor afișate pe baza sentimentului conținutului paginilor vizualizate.

În scopul furnizării reclamelor într-un mod personalizat se cunoaște de asemenea patentul US 0059713 A1/2012 *Matching advertisers and users based on their respective intents*, ce prezintă o metodă de selecție a scopului unui advertiser pe baza informațiilor primite de la utilizator. Dezavantajul metodei



constă în atributele statice pe baza cărora se deduce intenția utilizatorului și nu pe baza conținutului paginii vizualizate de utilizator la un moment dat.

Figura 1 prezintă platforma propusă. Textul din pagina accesată de utilizator este procesat. Pe baza meta-datelor NLP extrase se va alege reclama potrivită din inventarul de reclame corespunzătoare unei campanii de publicitate și va fi livrată în pagina utilizatorului. Teoria jocurilor intervine în alegerea celei mai potrivite reclame, astfel încât toți actorii implicați (utilizator, publisher și advertiser) să își optimizeze câștigurile.

Antrenarea algoritmilor de procesare de limbaj natural (NLP) s-a făcut în trei cicluri succesive, pentru a obține o acuratețe de cel puțin 80%. Pentru fiecare ciclu de antrenare / re-antrenare s-a folosit un set de texte relevante statistic extrase dintr-un corpus de date folosind metoda de eșantionare statistica simplă pentru a minimiza cantitatea de date necesară unui ciclu de antrenare / re-antrenare.

Pentru a putea antrena un algoritm de extragere entități dintr-un document din limba engleză, s-a utilizat un corpus de date ce constă din 50.000 de articole extrase din feed-uri RSS ale principalelor agenții de presă din SUA și Marea Britanie (Fox News, CNN, Reuters, New York Times, Washington Post, BBC, The Sun etc.). Articolele vin într-un format HTML, din care - pentru corpus - este necesar doar textul principal al articolului, lăsând la o parte conținutul irrelevant, cum ar fi reclamele sau textele laterale. Pentru a crea un corpus adecvat pentru antrenare algoritmi NLP, entitățile trebuie identificate și adnotate în textul articolelor. Entitățile din articol sunt identificate și marcate corespunzător. 70% din articole sunt folosite pentru antrenarea modelului iar restul sunt folosite pentru testarea acestuia. Textul este adnotat ca în următorul exemplu:

According to the former <START:NATIONALITY> Egyptian <END> <START:TITLE> Antiquities Minister <END> <START:PERSON> Zahi Hawassa <END> , who participated in the excavations <START:TEMPORAL\_TIME> 30 years <END> ago.

Pentru antrenarea modelului în limba română s-au urmat aceiași pași ca în varianta pentru limba engleză, corpusul pentru limba română este obținut din site-uri și agenții de știri din România (adevarul.ro, hotnews.ro, stirileprotv.ro, zf.ro, mediafax.ro, agerpres.ro, etc.). S-a folosit un corpus de 20.000 de articole.

Pașii urmăți în procesul de antrenare / re-antrenare modelului sunt:

1. Pregătirea corpusului folosit pentru antrenare și testare.
2. Extragerea textului principal al articolului din sursa HTML a articolelor.
3. Identificarea entităților pentru articolele din corpus.
4. Adnotarea entităților.
5. Antrenarea modelului.
6. Verificarea modelului.

Modelele obținute pentru limbile engleză și română au fost încărcate de un micro-serviciu Java Spring care rulează sub Java Virtual Machine (JVM). Serviciul numit "nlp-service" poate fi accesat de serviciul API Gateway, care este - la rândul lui - un alt micro-serviciu, numit "api-gateway-service". API Gateway are un endpoint /rest/v1/entities. End-pointul accepta un JSON în formatul:

DocumentRequest {

content (string, optional): conținutul în format text

contentUri (string, optional): URI către conținut (content și contentUri se exclud reciproc), în acest caz textul articolului este obținut prin descărcarea articolului, după care se extrage textul din codul HTML

language (string, optional): limbajul textului în format ISO 639-3

}

Rezultatul procesării este întors în format JSON și conține entitățile extrase, se întoarce o listă cu entități, fiecare element din listă având un câmp ce descrie tipul, numărul aparițiilor și entitate extrasă. Un posibil rezultat:

{

"entities": [{

"type": "PERSON",

"count": 13,

"mention": "Trump"

}, {

"type": "TITLE",

"count": 10,

"mention": "President"

}, {

"type": "LOCATION",

"count": 10,

"mention": "US"

}, {

"type": "ORGANIZATION",

"count": 3,

"mention": "Democratic"

}, {

"type": "NATIONALITY",

"count": 2,



```
        "mention": "Americans"
      }, {
        "type": "TEMPORAL:DATE",
        "count": 1,
        "mention": "January"
      }, {
        "type": "TEMPORAL_TIME",
        "count": 1,
        "mention": "two minutes"
      }, {
        "type": "IDENTIFIER_EMAIL",
        "count": 1,
        "mention": "donald.j.trump@gov.usa.com"
      }
    ]
  }
}
```

Setul de meta-date ce trebuie extrase pentru personalizarea semantică a reclamelor livrate conține, pe lângă entitățile extrase, și clasificarea conținutului după taxonomii standard (IAB – The Interactive Advertising Bureau și, respectiv, IPTC - News Categories Taxonomy for the Media), împreună cu sentimentul articolului (negativ, neutru, pozitiv).

Pentru antrenarea algoritmilor de clasificare IAB, IPTC și sentiment s-au folosit algoritmi cu învățarea profundă (Deep Learning – CNN). În învățarea profundă (Deep Learning), o rețea neuronală convoluțională (CNN sau ConvNet) este o clasă de rețele neuronale profunde, aplicată cel mai frecvent la analiza imaginilor vizuale. Cercetările noastre au arătat că rețelele neuronale convoluționale se pretează foarte bine și domeniului NLP. Corpusul de date folosit pentru antrenarea rețelelor CNN a fost clasificat automat, cu intervenție umană pentru calibrare, iar eșantioanele folosite pentru ciclurile de antrenare / re-antrenare au fost extrase folosind aceeași metodă de eșantionare simplă aleatorie.

După extragerea meta-datelor NLP din pagina vizualizată de utilizator se pot livra reclamele personalizate. Soluția permite selecția reclamelor pe baza unor reguli predefinite. Entitățile extrase din pagina vizitată sunt potrivite cu regulile de livrare pe baza entităților (cu tip), clasificare a conținutului (taxonomiile IAB și IPTC) și analiza sentimentului (neutru, negativ sau pozitiv).

Regulile de livrare create pot asigura siguranța unui client de publicitate (a mărcii). De exemplu dacă sentimentul extras din text este pozitiv sau neutru față de marca Adidas, se vor livra reclame despre marca Adidas. Prin impunerea unui sentiment neutru sau pozitiv se asigură ca marca pentru care se face

publicitate nu este asociată cu un articol posibil negativ. Dacă sentimentul extras din pagină este negativ față de marca Adidas, platforma livrează un anunț publicitar pentru altă marcă de pantofi sport.

Regulile de livrare ce folosesc clasificarea conținutului furnizează reclame pe baza categoriei asociate cu conținutul paginii vizitate. Dacă pentru o pagină se obține clasificarea "Travel" (IAB) or "holiday or vacation" (IPTC), entitatea extrasă este „Zurich” de tipul LOCATION, iar articolul citit este neutru sau pozitiv față de o călătorie în Zurich, se vor livra reclame legate de călătorie sau vacanțe în Zurich. Se poate ține cont și de locație în alcătuirea regulilor de livrare.

Toate regulile de livrare semantice (entități, sentiment, clasificari) se pot combina între ele pentru a obține reguli mai complexe pentru definirea exactă a contextului în care reclamele asociate unei campanii de publicitate să apară pentru maximizarea ROI (Return On Investment).

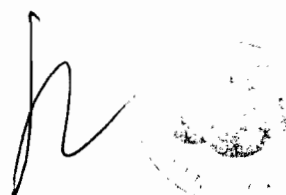
Figura 2 prezintă cum funcționează serviciul de publicitate online de tip semantic bazat pe inteligență artificială.

Pe baza metadatelor extrase algoritmul de alegere a reclamelor încearcă o potrivire între cuvintele cheie asociate reclamelor și aceste metadata. Printr-un mecanism de licitație în timp real se determină prețul pentru reclamă și se asigură echitabilitatea între actorii implicați.



## Revendicări

1. Metodă de furnizare a reclamelor pe baza contextului utilizatorului, **caracterizată prin aceea că** contextul utilizatorului este dedus prin analiza semantica a textului relevant din pagina web vizualizată de utilizator folosind metode de procesare a limbajului natural cu Inteligență Artificială. Algoritmi inspirați din Teoria Jocurilor intervin în alegerea celei mai potrivite reclame, astfel încât toți actorii implicați (utilizator, publisher și advertiser) să își optimizeze câștigurile.
2. Metodă ce permite extragerea textului relevant din sursa brută, neprelucrată HTML a articolelor online, **caracterizată prin aceea că** textul relevant este extras din sursa brută folosind algoritmi euristici, independenți de structura HTML a articolelor procesate și a limbii folosite în conținutul articolelor. Articolele vizualizate de utilizator sunt preluate în format brut HTML de pe site-ul public al Publisherului, din care, pentru analiza semantică pentru furnizarea reclamelor, este necesar doar textul relevant articolului lăsând la o parte părțile irelevante din pagina brută HTML, cum ar fi reclamele, textele laterale sau alte referințe.
3. Metodă de antrenare / re-antrenare algoritmi de Inteligență Artificială cu învățarea profundă (Deep Learning – CNN), **caracterizată prin aceea că** antrenarea se face în cicluri succesive, pentru a obține o acuratețe de cel puțin 80%. Pentru fiecare ciclu de antrenare / re-antrenare se folosesc un set de texte relevante statistic extrase dintr-un corpus de date folosind metoda de eșantionare statistica simplă pentru a minimiza cantitatea de date necesară unui ciclu de antrenare / re-antrenare. Folosind un set inițial de date, relevant statistic, este antrenat un model. Setul inițial este obținut folosind o metodă de eșantionare aleatoare, proporția acestuia din corpusul datelor fiind 11%. Mai apoi, 30% din date sunt folosite pentru testare, fază în care se decide care dintre algoritmi (extragere entități, clasificare, sentiment) vor fi reantrenați în ciclurile următoare. Se repetă ciclul de antrenare/testare de un număr prestabilit ( $\geq 3$ ) până se obține o acuratețe satisfăcătoare ( $> 80\%$ ).
4. Metodă de extragere a meta-datelor NLP multi-limbă din pagina online, **caracterizată prin aceea că** textul preluat din pagina online vizualizată de utilizator este clasificat automat folosind taxonomii standard (IAB și IPTC), iar extragerea entităților și a sentimentului (pozitiv, negativ și neutru) din text se face cu algoritmi cu învățarea profundă (Deep Learning – CNN).
5. Metodă de alegere a reclamei relevante în funcție de meta-datele NLP extrase din pagina online vizualizată de utilizator, **caracterizată prin aceea că**, pentru fiecare reclama din inventarul de reclame asociate unei campanii, se definesc reguli de stricte de livrare a reclamelor. În plus, regulile de livrare pot fi combinate cu alte atribute statice legate de comportamentul utilizatorului, cum ar fi: țara de origine a utilizatorului, limba acestuia, setări temporare. Algoritmi inspirați din Teoria Jocurilor intervin în alegerea celei mai potrivite reclame, astfel încât toți actorii implicați (utilizator, publisher și advertiser) să își optimizeze câștigurile. Printr-un model de licitație în timp real se determină prețul și asigură echitabilitatea între jucători. Pentru a găsi oferta optimă, problema este formulată ca o problemă de optimizare.



# Desene



Figura 1

The screenshot shows a web page with a navigation bar (Sun TV & SHOWBIZ, NEWS, FABULOUS, MONEY, MOTORS, TECH, DEAR DEBTS, PUZZLES, TOPICS & /) and a main article titled "MAKING A SPLASH TripAdvisor reveals Europe's best beaches... and Bournemouth is at number FIVE". The article text mentions Bournemouth as the fifth best beach in Europe. Below the article are social media sharing icons and a "COMMENTS" section.

Overlaid on the right side of the page is an advertisement for "GENERIC HOLIDAYS inc." with the text: "Haven't booked a summer vacation yet? Experience the best beach in Spain with GENERIC HOLIDAYS inc. >BOOK NOW<".

Below the advertisement, there is a list of criteria for how AdServista functions:

**Cum funcționează AdServista?**

1. Utilizatorul accesează pagina Web
2. Pagina este analizată semantic:
  - Conținutul relevant este extras din pagina Web
  - Limba conținutului este identificată automat (ENG)
  - Analiza de sentiment (POZITIV)
  - Este generat scorul clickbait (22%)
  - Categorisire: shopping or vacation
  - Extragere entități:
3. Mixare cu informații de comportament.
4. Potrivire semantică între pagină și reclame disponibile
5. Selecția celei mai potrivite reclame (teoria jocurilor)..

Figura 2

