

---

*Metode statistice aplicate crossdisciplinar în  
serii enumerative sau atributive nenumerice.  
Analize frecvențiale statistice în (mini)corpusuri  
filologice*

**Prof. univ. dr. habil. Gheorghe Săvoiu**

*Universitatea din Pitești*

**Conf. univ. dr. habil. Constantin Manea**

*Universitatea din Pitești*

**Rezumat**

*Una dintre cele mai simple serii de date statistice, apărută ca rezultat al prelucrării prin grupare a realității lingvistice sau filologice, aprofundată științific fie sub denumirea de serie atributivă nenumerică fie sub aceea de serie enumerativă este supusă în articol investigații specifice a gândirii statistice de început, respectiv analizei frecvențelor de apariție a unor cuvinte sau expresii specifice. Metoda analizei frecvențiale în (mini) corpusuri filologice poate să devină astfel, prin aplicabilitate crossdisciplinară, o metodă utilă de validare în lingvistica modernă, cu șansa de a oferi soluții de alegere criterială riguroasă sau argumente pertinente, extrase din uzul limbii în situații de ambiguitate și chiar de incertitudine selectivă. Tipul de serie investigată în articol, tip dominant în universul filologic sau lingvistic, dobândește o utilitate specială ce poate fi asigurată prin prelucrarea rapidă cu ajutorul metodei sau analizei frecvențiale. Metoda frecvențială permite luarea de decizii prompte în context de incertitudine filologică sau lingvistică, oferind indicatori statistici diverși capabili să confere veridicitate argumentației probatorii axate pe uzul lingvistic, prin dezvoltarea statică, dinamică, spațială și structurală a cuantificării în (mini)corpusuri filologice a frecvențelor de apariție, prin confruntare de opțiuni, prin evaluarea concentrărilor sau diversificărilor și în final chiar prin valorificarea profilului statistic în lingvistica modernă. Autorii oferă câteva exemplificări pertinente ale utilității metodei prin crossdisciplinaritate, care vizează în acest caz simplificarea deciziei, atât în filologia tradițională, cât și în lingvistica modernă, valorifică informații statistice trecute, pentru a desprinde tendințele de utilizare de o manieră prospectivă cât mai corectă a limbajului științific clasic sau modern. Pentru aplicabilitatea (mini)corpusul filologic modern de tip Internet accesat cu motoare de căutare devine populația statistică observată și prelucrată, iar prezentarea tabelară și reprezentarea grafică adecvată delimitează imaginea*

---

*cantitativă a probatoriului. În final, autorii anticipează evoluția către o lichiditate semnificativă sau chiar către o volatilitate excesivă a limbajului modern și asimilând analiza fluxurilor de date ale unor serii enumerative, tratează și identifică statistic soluții subliniind impactul major al metodelor aplicate crosdisciplinar în cadrul oricărei cercetări științifice filologice sau lingvistice.*

**Cuvinte cheie:** crosdisciplinaritate, metodă statistică, seria enumerative, seria homogradă sau de attribute, analiză frecvențială, (mini) corpus filologic sau lingvistic, indicatori de concentrare – diversificare, metoda profilului statistic, uzul lingvistic.

**Coduri JEL:** C46, C49.

### 1. Introducere

Metodele statistice teoretice își extind permanent aplicabilitatea lor practică și devin în acest fel tot mai utile în variate domenii ale cercetării științifice, conturând o crosdisciplinaritate evidentă. Așa cum este concepută și descrisă aplicativ în acest articol, crosdisciplinaritatea “*constituie o abordare care selectează, îmbină, asociază, combină, aplică metode unice în realități științifice diverse sau pune în practica bine delimitată disciplinar metode ale unei științe în corpul metodologic al altora, devenind un concept generic în creativitatea demersului investigativ sau metodologic, și constituie, în final, un prim pas către apariția și delimitarea de noi discipline sau științe*” (Săvoiu, et al., 2020, p. 8).

Crosdisciplinaritate modernă oferă soluții complexe prin originalitatea abordărilor specifice sau prin creativitatea transpunerii lor în alte științe, pornind de la simpla constatare că simplifică un demers clasic devenit mult prea uzual și uneori tot mai ineficient, în raport cu evoluțiile realității. În investigațiile și argumentațiile cercetătorilor aflați în fața unor fenomene tot mai variate și complicate, de la cel sociologic, la cel filologic, de la cel biologic, la cel demografic etc., singura opțiune este transformarea viziunii izolatoare și consacrate unidisciplinar într-o investigați ă echipă de tip crosdisciplinar prin apelul la metode diverse, între care cele statistice, matematice și fizice apar ca prioritare. Astfel, multe din problemele lingvistice destul de dificil de soluționat pot fi rezolvate apelând la metode statistice dintre cele mai simple. Analiza textelor bazată pe recurențe statistice devine uneori o chestiune de viață și de moarte, așa cum ar putea învăța, astăzi, orice student aflat la primele sale lecții de limbă engleză, de la un uimitor lingvist, cercetător și scriitor, în același timp, ca David Crystal, care îi exemplifică expresia *a matter of life and death*, apelând la celebra sa lucrare *The English Language: A Guided Tour of the Language*, și descriind modalitatea prin care

---

același student atras de învățarea corectă a limbii engleze ar putea scăpa de la o gravă acuzație de plagiat sau de pedepsire prin spânzurare, oferind drept probă salvatoare o analiză profundă – simultan lingvistică și statistică – a unei scrisori care aparține criminalului și care i-a fost atribuită studentului în mod fals, urmărind doar frecvența utilizării unor expresii, a unui anumit stil sau coloratura distinctivă a limbajului ... (Crystal, 2002).

Acest articol este rezultatul unor căutări comune ale autorilor, apărute din dorința de a simplifica deciziile legate de uzul corect filologic sau lingvistic al limbajului și de a valida/invalida diverse ipoteze de analiză frecvențială sau a contura ca indicatori utili datele referitoare la frecvența în uz pentru unii termeni sensibili, în ceea ce privește normarea/standardizarea lingvistică. Patima investigării crosdisciplinare a unor recurențe și convergențe lingvistice, spiritul de echipă și dorința de a aplica metode statistice în domenii diferite, ca procese de investigație influențate de spațiu, timp și structură (Manea, Săvoiu, 2014), au generat două perioade de investigații prin accesări în 2013-2014 și respectiv în 2019-2020.

## 2. Recenzia literaturii de specialitate

În cercetarea filologică românească, mai veche sau mai recentă, având la bază metoda statistică a analizelor frecvențiale, primele studii au pornit fie de la dicționare ale limbii române, fie de la anumite texte, alese în așa fel încât să aibă un anumit grad de reprezentativitate pentru limba română. Primul filolog preocupat de aplicarea cuantificărilor statistice cu rol de pionier în domeniul asigurării unei metode sau analize probatorii în studiul structurii etimologice a lexicului românesc a fost A. de Cihac. Astfel acest evaluator a realizat o falsă statistică prin consultarea incorectă metodologic a unor dicționare și glosare selectate special pentru a obține anumite rezultate, fără justificarea științifică referitoare la lista completă din care s-a realizat prelevarea nereprezentativă pentru perioada 1870-1879, precum și la metoda de selecție în sine, evident dirijată și implicit subiectivă, mărturisită în prefața la volumul al doilea al dicționarului său etimologic, intitulat *Dictionnaire d'étymologie daco-romane. Éléments slaves, magyars, turcs, grecs-modernes et albanais*, Francfort s/M, apărut în 1879, care constituie practic primul dicționar științific al limbii române. A. de Cihac a realizat în final o apreciere aproximativă a unor frecvențe relative, luând în calcul cuvintele „*nederivate*” din dicționarul alcătuit de el, fără a se opri și asupra unităților lexicale ale limbii române. Pornind de la cele „*aproximativ 500 de cuvinte latine, 1.000 de cuvinte slave, 300 de cuvinte turcești, 280 de cuvinte grecești moderne și 20 până la 25 de cuvinte maghiare sau albaneze*” (Dimitriu, 1973, p. XIII), A. de Cihac apreciat eronat că elementul latin „*care constituie fără îndoială substanța*

---

*limbii române [...] nu numai că a rămas aproape staționar, după primirea sa, în ceea ce privește fondul vocabularului, dar acesta din urmă trebuie chiar să fi pierdut multe cuvinte ca urmare a atâtor tulburări la care aceste nefericite locuri au fost teatrul timp de secole”* (Dimitriu, 1973, p. VIII). Concluziile evaluărilor statistice aproximative și dirijate ale lui A. de Cihac diminuează locul și rolul elementului latin moștenit în vocabularul limbii române, în mod subiectiv. Din păcate, realitatea faptelor de limbă a fost grav viciată de datele furnizate de astfel de cuantificări de A. de Cihac, cât și de aceea în care Sextil Pușcariu, o justifica pe prima, folosind procentaje apropiate și o revalida în 1920. Frecvențele relative calculate și prezentate păreau să indice faptul că limba română ar avea o structură etimologică *preponderent slavă* a lexicului (cu o valoare cu puțin peste două cincimi), iar elementul *latinesc*, moștenit, nu era de fapt decât *una dintre celelalte* cincimi constitutive, celelalte *două cincimi* rămase reunind elemente lexicale turcești și elemente cu origine eterogenă: maghiară, neogrecescă, albaneză etc).

În primul rând, cuantificările statistice amintite nu s-au sprijinit nici pe o metodologie coerentă și nici pe un aparat statistico - matematic adecvat, deoarece, pe de o parte, ele nu corespund realității cifrelor date de „*indicele*” de la sfârșitul volumului al doilea, iar pe de altă parte numărarea cuvintelor nu a luat în considerație raportul dintre *cuvânt* și *variantă*. În legătură cu primul aspect, care pune în cauză și cuantificarea statistică făcută de Sextil Pușcariu după „*indicele*” de cuvinte al lui Cihac, Mircea Seche arată că numărătoarea nu corespunde realității din două motive clare (Seche, 1966, p. 107): i) în primul rând, cuvintele cuprinse în „*indicele*” lui A. de Cihac sunt în număr de peste 8900 (nu 5765, așa cum fals apreciase Sextil Pușcariu); ii) suma totală a cuvintelor înregistrate de dicționar (așadar, nu numai de *indicele* de la finalul volumului al doilea) este de 17645, excluzându-se nu numai toponimele, așa cum era și firesc, dar și variantele (fonetice sau lexicale) – Ambele argumente sunt de o importanță capitală pentru a sublinia lipsa completă de acuratețe a cuantificării statistice și în final lipsa unei aprecieri motivate și dovedite științific. Mai mult, lipsa unui tratament metodologic unitar sau absența unei comparabilități statistice asigurate este la fel de evidentă, pentru cuvintele de origine latină, A. de Cihac a luat în considerație numai bazele, iar pentru celelalte, și derivatele. Mircea Seche arată că totalitatea cuvintelor din dicționarul lui A. de Cihac se repartizează pe straturi etimologice după cum urmează: „*elemente de origine latină (și derivate ale acestora) – 6141; elemente de origine slavă – 4691; elemente de origine turcă (și derivate ale acestora) – 1250; elemente de origine greacă (și derivate ale acestora) – 1100; elemente de origine maghiară (și derivate ale acestora) – 1026; elemente comune cu albaneza (și derivate ale acestora) - 90*” (Seche, 1966,

---

p.108). Procentajele corespunzătoare acestor noi date sunt: elementul latin (și derivatele sale) reprezintă 45,6% din total, elementul slav (și derivatele acestuia) reprezintă 34,8%, elementul turc (și derivatele acestuia) reprezintă 7,1%, elementul neogrec (și derivatele acestuia) reprezintă 6,2%, elementul maghiar (și derivatele acestuia) reprezintă 5,8% din total, iar elementul comun cu albaneza (și derivatele sale) reprezintă doar 0,5% din totalul cuvintelor cercetate. Se observă așadar că diferența este sensibilă, atât față de frecvențe relative indicate de A. de Cihac, cât și față de cele indicate de Sextil Pușcariu.

În al doilea rând, statistica lui A. De Cihac pune pe același plan cuvinte inegale ca forță *circulatorie* și ca volum sau conținut *semantic*. Numărătoarea lui Cihac (ca și aceea a lui Pușcariu de mai târziu) nu ține cont de faptul esențial că unitățile lexicale ale unei limbi nu pot sta în nici un caz pe același plan în ceea ce privește importanța lor *relativă*. *Ponderea* pe care o au cuvintele de diverse origini marchează identitatea unui vocabular, ca și gradul lor de reprezentare în vocabularul de bază al limbii studiate. Această *pondere statistică*, măsurată ca frecvență relativă trebuie să fie studiată din punctul de vedere al dinamicii vocabularului limbii respective. Deși Sextil Pușcariu arată în mai multe lucrări, de exemplu *Locul limbii române între limbile romanice, Limba română*, vol. I, că latinitatea limbii române, vizibilă din întreaga sa structură, se poate deduce și din „*materialul de construcție*” al vocabularului, însă nu prin simpla numărare: „*Orice dicționar etimologic e unilateral, căci el ține seamă numai de originea, nu și de circulația cuvintelor în limbă. Într-un dicționar etimologic cuvinte cunoscute și întrebuințate zilnic de fiecare român, din orice parte a țării, ocupă o unitate de loc, întocmai ca vorbele întrebuințate numai în câte o regiune – și acolo foarte rar – și necunoscute tuturor celorlalte ținuturi*” (Pușcariu, 1976, p. 181). În aceeași lucrare (Seche, 1966), este menționat și un posibil reper al aproximărilor sau o sursă probabilă a estimărilor false ale lui A de Cihac, respectiv o altă «*statistică*» din 1840 a lingvistului rus I. Hinculov, de unde se pare că s-ar fi inspirat controversatul A. de Cihac în ceea ce privește procentajul diverselor elemente etimologice ale lexicului.

Lingvistul care a sesizat prima oară falsitatea fundamentală a statisticii lui A. de Cihac a fost savantul român Bogdan Petriceicu Hașdeu, care a dezvoltat teoria circulației ca o adaptare a teoriei circulației în economie, aplicând-o la cuvintele din lexicul limbii române. Bogdan Petriceicu Hașdeu a arătat în mod lămurit că frecvența sau circulația cuvintelor are o importanță decisivă în stabilirea fizionomiei lexicale a unei limbi. Criticând clasificarea etimologică făcută de A. de Cihac, pentru care: „*L'élément latin de la langue roumaine ne représente guère aujourd'hui qu'un cinquième de son vocabulaire, tandis que l'élément slave y entre pour le double ou pour*

---

2/5 à peu près...” (Hașdeu, 1984, p. 73). Bogdan Petriceicu Hașdeu a relevat defectul unor false «statistici» care pun pe același plan cuvinte inegale ca volum semantic și ca forță circulatorie: „*Dicționarul nu ne dă și nu ne poate da circulațunea limbei; și tocmai acesta este punctul cel esențial*”. Dimitrie Macrea a realizat în 1942 o altă statistică bazată pe cuvintele cuprinse în CADE (*Dicționarul enciclopedic ilustrat „Cartea Românească”*. Partea I: *Dicționarul limbii române din trecut și de astăzi* de I.-A. Candrea. Partea II: *Dicționarul istoric și geografic universal* de G. Adamescu, București, 1926-1931), a cărui concluzie, considerată cel puțin uluitoare, după părerea regretatului academician Alexandru Graur) era aceea că elementele latine reprezintă 20,58% din total, cele slave 16,41%, cele franceze 29,69%, iar restul de 33,32% ar fi format din cuvinte avându-și obârșia în limbi care nu contează procentual prea mult ori din cuvinte cu origine necunoscută. Este vorba așadar despre semnificația crucială, recunoscută unanim în lingvistica generală, a lexicului fundamental al unei limbi sau partea ei cea mai rezistentă și mai importantă, însuși *nucleul* lexicului limbii; această noțiune se opune *masei vocabularului* sau *vocabularului secundar* și pentru ea se folosesc mai multe denumiri, cum ar fi: *vocabular de bază*, *vocabular fundamental*, *vocabular esențial*, *fond principal lexical*, *fond principal de cuvinte*, mai rar *fond lexical uzual*. (Hristea et al., 1984, p. 14). Rolul și valoarea pe care le au diferitele elemente aflate în componența lexicală a limbii române pot fi mai bine precizate dacă se au în vedere datele pe care le oferă unele statistici ulterioare, efectuate pe baza unor lucrări de o mai mare profunzime, în aceeași măsură filologică, dar și statistică. Într-o astfel de cuantificare statistică întocmită de Sever Pop, în anul 1948, asupra *Dicționarului limbii române din trecut și de azi* de I. A. Candrea (apărut în Editura *Cartea Românească*, București, 1931), se constată ca numărul cuvintelor de origine latină se ridică la 8.800, la care trebuie să se adauge 14.000 de neologisme primite din limbile romanice, ceea ce dă un total de 22.800 de termeni (se impune o comparare cu situația din limba franceză, unde, după *Dictionnaire de l'Académie*, 1878, din 32.000 de cuvinte, 20.000 erau de origine savantă sau străină și numai 12.000 constituiau cuvinte franțuzești de origine populară). Numărul cuvintelor de origine slavă (preluate din slava comună ori din limbi slave moderne ca bulgara, sârba, ruteana, rusa, polona) era de asemenea mare, de aproape 7.800, dar mulți dintre acești termeni sunt căzuți în desuetudine sau sunt termeni de uz regional, restrâns. O interesantă cercetare statistică a Mihaelei Bîrlădeanu, realizată în lucrarea *Structura etimologică a două vocabulare reprezentative: român și francez*, apărută în *Studii și cercetări lingvistice*, nr. 6/1983, aceasta compară vocabularele reprezentative ale limbilor română și franceză, găsind următoarea structură etimologică pentru vocabularul reprezentativ al limbii

noastre: latină: 1. moștenite: 30,45%, 2. savante 1,77%; formații interne 24,81%; substrat 0,96%; superstrat vechi slav 8,91%; neogrecă 1,11%; împrumuturi din limbile slave moderne 1,80%, împrumuturi romanice din: franceză 7,64%, italiană 0,54%; turcă 0,73%; maghiară 1,27%; germane 0,27%; engleză; onomatopeice 0,23%; etimologie multiplă 17,36%; origine incertă 2,08%. Dintre lucrările relativ recente de cuantificare și de analiză statistică în domeniul filologic, câteva aparțin chiar autorilor articolului (Manea, 2004; 2009; Manea, Săvoiu, 2014). Metoda analizei statistice frecvențiale a constituit metoda majoră aplicată într-unul dintre segmentele principale ale volumului *Structura etimologică a vocabularului neologic (cu specială referire la anglicismele din limba română)*, apărut în 2004, cât și în *Încercare statistică asupra ortografierii cu â și î*, în 2009.

### 3. Metodologie

Seria enumerativă este cea mai simplă formă de prezentare statistică a unei populații grupate după cele mai banale soluții criteriale în universul noncantitativ, respectiv în lumea cuvintelor, specifică filologului sau lingvistului (Săvoiu, 2012). O serie de acest tip poate fi constituită din cea mai simplă listă a prenumelor și numelor unor oameni, grupați după un anumit criteriu organizatoric, administrativ, structural, spațial, temporal etc. În urma valorificării metodei grupării statistice a acestei populații în raport cu prenumele sau numele rezultă o serie de repartiție sau de distribuție de frecvențe ale prenumelor sau numelor identificate ca distincte, serie cunoscută și ca *serie homogradă sau atributivă nenumerică* primul șir fiind enumerativ sau calitativ iar cel de-al doilea cantitativ sau numeric (Săvoiu, 2003, p.118). *Într-o manieră simplificată de redare orice serie statistică enumerativă prelucrată prin grupare devine o serie atributivă nenumerică - cu referire la primul șir care este esențial fiind cel discriminant) – iar cel de-al doilea șir devine numeric ca urmare a faptului că s-au atribuit frecvențe de apariție unor cuvinte, expresii, calificative, ierahizări etc (Săvoiu, et al, 2006). În figura 1 sunt descrise într-un model general câteva serii atributive nenumerică sau enumerative uzuale:*

#### Forme de existență a seriei statistice atributive nenumerică sau enumerative prelucrate prin grupare

Fig. nr. 1

Varianta ( $x_i$ )	cuvânt	expresie	calificativ	ierahie	corect/incorect	da/nu
Frecvența ( $n_i$ )	$n_1$	$n_2$	$n_3$	$n_4$	$n_5$	$n_6$

Sursa: Realizat de autori

---

Seria statistică de acest gen cuprinde două șiruri paralele, primul de variante ale variabilei atributive nenumerice sau enumerative discriminate după un criteriu lingvistic, filologic, ierarhic, de evaluare, de cunoaștere etc. iar cel de-al doilea de frecvențe de apariție a variantelor din primul șir. Varianta ( $x_i$ ) poate reprezenta un cuvânt (substantiv, prenume, verb, adjectiv, adverb, etc), un substantiv (nume propriu uzual sau nu, un obiect comun) o formă de plural, o formă verbală, o formă de redactare sau pronunțare în uz, o soluție de scriere de uz local sau general, o sursă de intrare dintr-o altă limbă a unui neologism, o formă corectă sau incorectă de pronunțare sau redactare, un calificativ, o nuanță sau un calificativ ierarhic, un nivel de acoperire, un nivel de acceptare sau refuz, un raspuns mono sau plurisilabic (da, nu sau da, nu, nu știu/nu răspund), un mod de acceptare sau inacceptare etc. Frecvența ( $n_i$ ) este o informație numerică rezultată din agregarea cantitativă variantelor identice, suma tuturor frecvențelor reunind totalitatea cazurilor din serie.

Un exemplu la îndemână de serie enumerativă dintre cele mai simple și relevante prin simpla ei constituire, exemplu redescoperit *în ultimele trei decenii*, considerat celebru pentru elitismul celor înregistrați, dar și pentru asocierea lor cu militantismul antireligios sau ateismul de tip comunist *în România*, este acela al seriei de personalități incinerate (arse) și nu înhumate tradițional ortodox (fără a respecta ritualurile ortodoxe ale privegherii decedatului, *înmormântării și pomenirilor devenite cu trecerea timpului pomeni* cu un conținut dintre cele mai inexplicabile cu puțință)... Nu numai legal, dar și biblic, incinerarea constituie alături de înhumare o modalitate similară ca finalitate, prin care toți aceia care mor se întorc în țărână, prima descriind un proces accelerat de descompunere a unui cadavru prin ardere și transformare a unui corp uman în cenușă, iar cea de-a doua un proces mult mai lung în ani, unul natural de metamorfozare a omului în țărână ca o firească “întoarcere în pământul din care ești luat; căci pământ ești și în pământ te vei întoarce”. (*Biblia sau Sfânta scriptură*, 2001, versiune redactată și adnotată de Bartolomeu Valeriu Anania Facerea 3:19, p.22).

Această serie simplă enumerativă cuprinde ca variabile atributive nenumerice atât numele și prenumele a 2153 de personalități incinerate, cât și profesiunile acestora și prin gestul implicit al incinerării. Seria denotă absența credinței religioase ortodoxe în cazul tuturor celor înscrși în cadrul acestei liste a indivizilor într-o societate românească declarată la recensăminte permanent dominant ortodoxă. Lista românilor considerați altfel decât majoritatea locuitorilor României ascunde aparent mai multe informații interesante, pe care le poate dezvălui prompt dacă este prelucrată statistic elementar sau supusă unei analize frecvențiale și transformată într-o serie atributivă nenumerică. În termenii statisticii clasice, o listă simplă de persoane incinerate se transformă



---

*într-o serie homogradă* a profesiunilor sau ocupațiilor, amalgamate din păcate *în corpul acesteia. Seria homogradă rezultată* oferă nu doar simple statistici, ci veritabile paradoxuri contrastând cu opinia comună despre procesul și scopul incinerării așa cum se desprinde din analiza celor care au apelat ceva mai frecvent la această soluție ...

Incinerarea a devenit istoric o practică recunoscută oficial prin legile României, având statut juridic legal ca și înhumarea, un statut tot mai actual *în condiții de pandemie Covid-19, trăite în prezent* cu mare intensitate și de populația țării noastre. Incinerarea în raport cu înhumarea oferă unele avantaje nu numai sanitare, ci mai extinse, asigurând o ieșire “*mai demnă din scena vieții*” într-un context viciat de impactul unor boli infecțioase tot mai agresive, de boli transmisibile sau contagioase apte să genereze epidemii și chiar pandemii, de accidente ce desfigurează cadavrele, de rămășițe de corp uman grav afectate de războaie, explozii, inundații etc.: i) presupune costuri mult mai mici decât cele ale înhumării pe durată mai mare de timp eliminând costurile tradițiilor clasice înhumării; ii) este o practică mult mai estetică și ecologică; iii) induce un sentiment al egalității între oameni în final; iv) extinde cultul morților, dilatând accentul pus pe suflet comparativ cu cel pus pe trup; v) urna cu cenușa umană poate fi păstrată acasă sau îngropată în mormântul oricărui cimitir; vi) incinerarea presupune respectarea voinței exprimate liber de persoana decedată de a nu constitui ulterior mai mult decât o amintire și nu o obligație pentru familie etc. (<http://www.incinerareamurg.ro/romani-celebri-care-au-fost-incinerati>).

Fără a intenționa *în niciun fel* o promovare a incinerării în detrimentul înhumării clasice sau a face reclamă unor societăți comerciale al căror obiect de activitate este incinerarea și nu înhumarea, o analiză statistică a impactului celi dintâi *în Europa relevă* că peste 2/3 ajungând chiar și la 70% din oameni aleg soluția incinerării, în Marea Britanie, Suedia, Danemarca, Cehia, Ungaria, atingând chiar proporții optim paretiene în Elveția (85%). Se identifică multe ierarhii frecvențiale surprinzătoare sau aparent neașteptate și în seria enumerativă a incinerărilor în România, conform datelor disponibile online la <http://www.incinerareamurg.ro/romani-celebri-care-au-fost-incinerati> și prelucrate de autori. *Seria homogradă rezultată* este prezentată parțial prin segmentarea în raport cu profesiunea sau ocupația, declarată de familiile celor incinerati în două subgrupe limitative ce oferă informații statistice simple axate pe frecvențe maxime și minime (Tabel 1).

**Frecvențele minime sau aparițiile cele mai rare (stânga) și cele maxime sau cele mai dese (dreapta) calculate din seria enumerativă a incinerărilor în România**

*Tabel 1*

Profesiune/ocupație	Frecvență minimă ( $n_i$ )	Profesiune/ocupație	Frecvență maximă ( $n_i$ )
Oameni de afaceri	2	Profesori	333
Politolog	2	Generali	164
Piloți de curse	2	Scriitori	158
Tipografi	2	Actori	125
Preot ortodox, deputat	1	Ingineri	111
Bancher	1	Cercetători	109
Statistician	1 (Mircea Biji)	Militanți comunisti	108

Sursa: Realizat de autori prin prelucrarea frecvențelor de apariție la: <http://www.incinerareamurg.ro/romani-celebri-care-au-fost-incinerati>

Sistemul de indicatori statistici pentru analiza frecvențială într-o serie de repartiție sau distribuție de tip homograd este alcătuit dintr-o diversitate de componente: “frecvențe absolute ( $n_i$ ), frecvențe relative ( $n_i^*$ ), frecvențe cumulate crescător ( $n_i \uparrow$  sau  $n_i^* \uparrow$ ) sau descrescător ( $n_i \downarrow$  sau  $n_i^* \downarrow$ ), densități de repartiție a frecvențelor ( $n_i/h_i$  sau  $n_i^*/h_i$ ). Frecvențele cumulate crescător sau descrescător permit identificarea cuantilelor ( $C_v$ ) din a căror numeroasă familie fac parte: mediana ( $Me$ ), quartilele ( $Q_1, Q_2, Q_3$ ), decilele ( $D_1, \dots, D_9$ ) și centilele ( $C_1, \dots, C_{99}$ ), care împart populația statistică în două, în patru, în zece și într-o sută de părți egale” (Săvoiu, et al., 2020, p. 130).

O varietate de concepte și instrumente *statistice* sunt aplicabile și, așa cum s-a dovedit adesea, foarte utile practic în toate domeniile de cercetare și în toate tipurile de abordări științifice. Acest adevăr este dificil sau chiar imposibil de contestat în filologie sau lingvistică, pornind de la valorificarea sistemului de indicatori frecvențiali, trecând prin coeficienții de concentrare și diversificare Herfindahl – Hirschman și Gini – Struck, iar în final valorificând inovativ profilul statistic în filologie sau lingvistică, în cazul concret al unui (mini)corpus lingvistic dezpovărându-l de ambiguități și incertitudini legate de noramre, uz, corectitudine etc. conceptele care sunt supuse în continuare observării, cuantificării și analizelor statistice frecvențiale sunt de dimensiuni mici cu rol mai degrabă metodologic și de amplificare a rolului crosdisciplinarității (dimensiunea lor fiind cuprinsă între 25 și 50 de itemuri, aproape toate fiind termeni aparținând vocabularului științifico-tehnic al limbii engleze contemporane (așa cum reiese și din tabelele detaliate

---

ale rezultatelor). Motoarele de căutare utilizate au fost *Google* și *Ask*, iar corpusurile de cuvinte (texte) accesate au fost constituite din material (texte, articole etc.) de tip academic, găsite pe Internet.

O problemă metodologică majoră cum este aceea a subiectivismului sau a falsificării rezultatelor unei investigații *directe* într-un (mini)corpus filologic, într-un anumit sens dorit de cel ce realizează construcția (mini) corpusului filologic și prelevarea de eșantioane, observate și prelucrate ulterior, nu poate fi eliminată în lipsa onestității cercetătorului, indiferent de gradul de corectitudine a analizei frecvențiale sau de prelevare a eșantionului cu respectarea extracțiilor aleatoare ale fiecărui cuvânt sau expresii pornind de la probabilități practice cunoscute matematic. La toate aceste atitudini legate de aplicarea incorectă a teoriei sondajului, acțiuni premeditate sau nu de coordonatorii unor cercetări directe ce devin astfel neștiințifice sau simple opinii filologice sau lingvistice, se adaugă și existența unor premise permanente de apariție a efectului Hawthorne. Într-o cercetare directă sau pe teren, efectul Hawthorne este un efect clasic cunoscut ca formă de reactivitate psihologică a respondenților prin care subiecții unei cercetări parțiale sau experimentale își modifică anumite aspecte ale comportamentului lor, în uzul limbajului cu impact filologic sau lingvistic, ca urmare a faptului conștientizat că sunt studiați fără ca această atitudine să constituie un răspuns la manipuări cu scopuri subiective (Sesardić, 2018, p. 21).

Categoriile cele mai importante de erori metodologice grave, legate de apariția în cercetări *indirecte* sau *documentare*, axate pe analize frecvențiale pot fi considerate: i) erori cauzate de tehnica de generare de (mini)corpusuri filologice, respectiv ca urmare a inadecvării formei sau tehnicii de eșantionare aleatoare, dirijată sau mixtă; ii) erori metodologice sau sistematice, în strictă conexiune cu dificultățile de prelevare și cuantificare din Internet, tehnici cu rezultate eronate atâta timp cât sunt lipsite de instrumente și softuri de control sau verificare a înregistrărilor specifice, de la forme de plural, la forme verbale, de la expresii, la exprimări integrale sau abrevieri specifice etc; iii) erori cauzate de neîndeplinirea constrângerilor de programare temporară spațială și structurală a căutărilor; iv) erori generate de limitările diferențiate ale motoarelor de căutare; v) erori provocate de algoritmi de optimizare a căutărilor; vi) erori de pre-procesare inadecvată utilizată pentru a reduce dimensiunea spațiului soluției și în final rezultatele cercetărilor la cele validate, credibile, verificate etc.

#### 4. Analize frecvențiale statistice în (mini)corpusuri filologice

Analiza statistică frecvențială statică (absolută) este practic cea mai simplă modalitate de aplicare a unei metode crosdisciplinar și investigativ cu scopul de a soluționa o problemă filologică (lingvistică), cum poate fi aceea a identificării pluralului gramatical corect conform uzului dominant într-un (mini)corpus filologic sau într-o bază de date distinctivă (Internet). În tabelul 2 este prezentată o astfel de analiză cu relevanța și irelevanța ei specifică, pornind de la ideea de confruntare suplimentară pe principii statistice a cel puțin două motoare de căutare în cadrul acelorași (mini)corpusuri.

#### Analiză statistică frecvențială statică a unui plural gramatical corect conform uzului majoritar într-un (mini)corpus filologic (Internet)

Tabel 2

Cuvinte cu forme multiple de plural în limba engleză	Rezultate obținute cu motoare de căutare		Observații
	Google (search)	Ask (search)	
<i>apsides</i>	96,300	8,330	
<i>apses</i>	427,000	76,300	Relevant
<i>apsises</i>	12,200	-	
<i>octopuses</i>	651,000	225,000	Relevant
<i>octopodes</i>	124,000	12,800	
<i>octopi</i>	523,000	133,000	
<i>addenda</i>	5,700,000	430,000	Relevant
<i>addendums</i>	468,000	105,000	
<i>addendas</i>	115,000	-	
<i>criteria</i>	445,000,000	42,800,000	Relevant
<i>criteria</i>	511,000	131,000	
<i>criteria</i>	679,000	244,000	
<i>antennae</i>	7,570,000	832,000	Irelevant*
<i>antennas</i>	2,840,000	4,770,000	
<i>apexes</i>	416,000	67,600	
<i>apices</i>	607,000	193,000	Relevant
<i>apparatus</i>	169,000,000	12,700,000	Relevant
<i>apparatuses</i>	7,670,000	515,000	
<i>appendixes</i>	2,580,000	264,000	
<i>appendices</i>	17,400,000	2,150,000	Relevant
<i>aquariums</i>	29,300,000	3,210,000	Relevant
<i>aquaria</i>	9,790,000	1,260,000	
<i>automatons</i>	528,000	1,260,000	
<i>automata</i>	29,200,000	1,450,000	Relevant
<i>bureaux</i>	73,700,000	1,150,000	Irelevant**
<i>bureaus</i>	30,600,000	3,210,000	
<i>cerebellums</i>	44,800	6,510	
<i>cerebella</i>	393,000	50,200	Relevant
<i>curricula</i>	17,700,000	2,440,000	Relevant
<i>curriculum</i>	6,860,000	694,000	
<i>formulas</i>	57,000,000	8,200,000	Relevant

<b>formulae</b>	13,100,000	1,940,000	
<b>genera</b>	96,300,000	5,450,000	Relevant
<b>genuses</b>	117,000	15,900	
<b>hiatuses</b>	279,000	42,500	
<b>hiatus</b>	34,900,000	4,140,000	Relevant
<b>maximums</b>	2,860,000	407,000	
<b>maxima</b>	131,000,000	6,640,000	Relevant
<b>minimums</b>	8,030,000	1,080,000	
<b>minima</b>	63,500,000	1,340,000	Relevant
<b>nuclei</b>	22,200,000	4,050,000	Relevant
<b>nucleuses</b>	73,200	8,270	
<b>phenomena</b>	68,900,000	10,800,000	Relevant
<b>phenomenons</b>	462,000	99,900	
<b>syllabuses</b>	532,000	144,000	
<b>syllabi</b>	4,930,000	799,000	Relevant
<b>strata</b>	65,600,000	4,060,000	Relevant
<b>stratums</b>	336,000	24,200	
<b>vortexes</b>	486,000	113,000	
<b>vortices</b>	2,610,000	448,000	Relevant

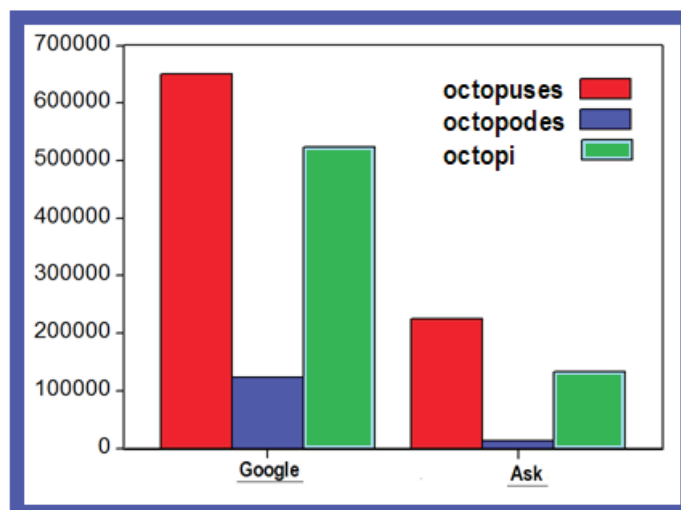
Sursa: Preluat de autori dintr-o investigație proprie mai veche (Manea, Săvoiu, 2014, pp.13-17).

Evaluând static sau absolut distribuțiile de frecvență ale pluralului gramatical adecvat în conformitate cu uzul sau utilizarea specifică a majorității subiecților într-un (mini)corpus filologic (Internet) pentru cuvinte ca *apse*, *addendum*, *antenna*, *apex*, *apparatus*, *appendix*, *automaton*, *bureau*, *criterion*, *cerebellum*, *curriculum*, *formula*, *genus*, *hiatus*, *maximum*, *nucleus*, *octopus*, *phenomenon*, *syllabus*, *stratum*, *vortex*), s-au identificat date relevante pentru majoritatea absolută (cu doar două excepții: *antenna* și *bureau*), s-au putut trage concluzii practice valoroase asupra utilității metodei analizei frecvențiale, abordată prin cel puțin două motoare de căutare. Reprezentările grafice ale unor valori frecvențiale comparate pot contribui printr-o mai bună vizibilitate la aprecierea pluralului relevant și probabil implicit corect, acolo unde informația cantitativă este relevantă, detaliind frecvențele absolute după mai multe motoare de cautare (de exemplu: Google și Ask).

---

**Analiză statistică grafică comparată a unui plural gramatical cu două motoare de căutare într-un (mini)corpus lingvistic (Internet)**

Fig. nr. 2



Sursa: Preluat de autori după o lucrare proprie mai veche (Manea, Săvoiu, 2014, pp.13-17)

Figura 2 subliniază că un grafic care compară distribuțiile formelor de plural studiate și cu ajutorul a două motoare de căutare poate fi mai util, prin vizibilitate superioară, în ceea ce privește evaluarea ipotezelor referitoare la adecvarea gramatical-fonetică a uneia dintre cele trei forme de plural, analizate statistic frecvențial. Informațiile cantitative relevante asigură că și frecvența ocurențelor va deveni relevantă, iar dacă apar practic stări de ambiguitate se poate apela firesc la mai multe motoare de căutare (Google și Ask).

În analize complexe se poate trece la o subîmpărțire a listei extinse de termeni specializați sau tehnici (dintre care unii au fost științifici doar inițial), după cum: a) ilustrează o chestiune legată de morfofonematică; b) ilustrează o simplă problemă de ortografie; c) ilustrează o categorie de termeni care nu pot fi numiți chiar termeni tipic tehnici sau științifici, deși sunt fără îndoială termeni livești.

“Iată câteva exemple de termeni din subcategoriile deduse din lista extinsă: a) *antenna, apex, apparatus, appendix, automaton, cactus, calyx, cerebellum, cerebrum, cicada / cicala, colloquium, cranium, criterion, curriculum, dilettante, discus, fauna, flora, formula, fungus, genus, hiatus, iambus, larynx, libretto, memorandum, novella, nucleus, palazzo, phenomenon, radius, radix, retina, rhombus, stratum, syllabus, tableau, tempo, trapezium,*

---

*vacuum, vertebra, vertex, vortex*; b) *bureau, flamingo, fresco, grotto, halo, manifesto, memento, motto*; c) *aquarium, candelabrum, cicerone, colossus, focus, grotto, gymnasium, hippopotamus, maximum, millennium, minimum, narcissus, persona grata, referendum, sanatorium, symposium, terminus, ultimum*” (Manea, Săvoiu, 2014).

Observațiile de detaliu pe care le-au putut face autorii în sprijinul ideii de concretizare faptică a unor ipoteze privind recurența cuvintelor în (mini) corpusuri de date încă din 2014, au relevat, că există de pildă foarte puțini termeni care au trei forme de plural dovedibile (viz. *octopus* „caracatiță”-pl. *octopuses, octopodes, octopi*, deși ultima formă este improprie / nerecomandată). La fel de instructive pot deveni și alte investigații axate pe o căutare asemănătoare privind distribuția substantivelor *apsis* „absidă” (pl. *apsides, apses, apsises*), *addendum* „adendă” (pl. *addenda, addendums, addendas*) și *criterion* „criteriu” (pl. *criteria, criterions, criterias*), posibil și *frustrum*.

În destule cazuri, marcarea uzului lingvistic, așa cum este înțeles și ilustrat de dicționarele anglo-americane, a fost într-atât de importantă încât lexicografuli respectivi au găsit de cuviință să gloseze (*MacMillan*) forme de plural neregulate drept *leme* legitime, de sine stătătoare (de exemplu, *bacteria, criteria, alga, data*). Când însă se ia în considerare și problema sensului, procedura în sine pare să își dovedească neputința totală: cum am putea verifica semnificația și utilizarea fiecărei ocurențe care apare în texte sau (mini)corpusuri pe care s-au făcut căutări? [ex: *genius* (pl. *geniuses* vs. pl. *genii*), *domino* – pl. *dominoes* („joc”) / *dominos* („mantie”), *index* – pl. *indexes/indices*, *stamen* – pl. *stamens/ stamina*, *milieu* – pl. *milieus/ ranc. milieux* [‘mi:ljø], *calculus* – pl. *calculuses*; med. *calculi* [‘kælkjulai], *polypus* – pl. *polypi*; *data* (pl., deși de obicei considerat nenumărabil) → sg. *datum*, *agenda* (pl. lui *agendum*; azi considerat sing.).

Analiza statistică frecvențială dinamică și comparată instrumental descrie o modalitate mai evoluată de investigație crosdisciplinară cu același rol de identificare a pluralului gramatical corect conform uzului dominant abordat în două momente diferite ale evoluției unui (mini)corpus filologic sau o bază de date distinctivă (Internet). În tabelul 3 este prezentată această analiză pe același fond tematic al pluralului corect (2014) dar de o manieră repetitivă (2020), cu un calcul suplimentar al *indicii evolutiv al relevanței aprecierii*:

**Analiză statistică frecvențială dinamică a pluralului gramatical corect  
conform uzului majoritar într-un (mini)corpus filologic (Internet)**

*Tabel 3*

Cuvinte cu forme multiple de plural în limba engleză	Rezultate obținute cu motoare de căutare în ani diferiți				Observații
	2014		2020		
	Google (search)	Ask (search)	Google (search)	Indice evolutiv (%)	
<i>apsides</i>	96300	8330	127000	131,88	
<i>apses</i>	427000	76300	<b>637000</b>	<b>149,18</b>	Relevant
<i>apsises</i>	12200	-	<i>5680</i>	<i>46,56</i>	
<i>octopuses</i>	651000	225000	<b>4450000</b>	<b>683,56</b>	Relevant
<i>octopodes</i>	124000	12800	162000	130,65	
<i>octopi</i>	523000	133000	2130000	407,27	
<i>addenda</i>	5700000	430000	<b>8560000</b>	<b>150,18</b>	Relevant
<i>addendums</i>	468000	105000	1170000	250,0	
<i>addendas</i>	115000	-	200000	173,91	
<i>criteria</i>	445000000	428000000	<b>698000000</b>	<b>154,83</b>	Relevant
<i>criteriaons</i>	511000	131000	1060000	207,43	
<i>criteriais</i>	679000	244000	2060000	303,39	
<i>antennae</i>	7570000	832000	8890000	113,21	
<i>antennas</i>	2840000	4770000	<b>74100000</b>	<b>26091,55</b>	**Relevant
<i>apexes</i>	416000	67600	597000	143,51	
<i>apices</i>	607000	193000	<b>2100000</b>	<b>349,42</b>	Relevant
<i>apparatus</i>	169000000	12700000	<b>190000000</b>	<b>112,43</b>	Relevant
<i>apparatuses</i>	7670000	515000	5320000	69,36	
<i>appendixes</i>	2580000	264000	3110000	120,54	
<i>appendices</i>	17400000	2150000	<b>22800000</b>	<b>131,03</b>	Relevant
<i>aquariums</i>	29300000	3210000	<b>59100000</b>	<b>201,71</b>	Relevant
<i>aquaria</i>	9790000	1260000	10600000	108,27	
<i>automatons</i>	528000	1260000	1370000	259,47	
<i>automata</i>	29200000	1450000	<b>52100000</b>	<b>178,42</b>	Relevant
<i>bureaux</i>	73700000	1150000	72800000	98,78	
<i>bureaus</i>	30600000	3210000	<b>118000000</b>	<b>385,62</b>	*Relevant
<i>cerebellums</i>	44800	6510	65500	146,21	
<i>cerebella</i>	393000	50200	<b>460000</b>	<b>117,05</b>	Relevant
<i>curricula</i>	17700000	2440000	<b>28000000</b>	<b>158,19</b>	Relevant
<i>curriculumus</i>	6860000	694000	8100000	118,08	
<i>formulas</i>	57000000	8200000	<b>155000000</b>	<b>271,93</b>	Relevant
<i>formulae</i>	13100000	1940000	52000000	396,94	
<i>genera</i>	96300000	5450000	<b>156000000</b>	<b>161,99</b>	Relevant
<i>genuses</i>	117000	15900	238000	203,42	
<i>hiatuses</i>	279000	42500	396000	141,94	
<i>hiatus</i>	34900000	4140000	<b>59400000</b>	<b>170,20</b>	Relevant
<i>maximums</i>	2860000	407000	3930000	137,41	
<i>maxima</i>	131000000	6640000	<b>548000000</b>	<b>418,32</b>	Relevant
<i>minimums</i>	8030000	1080000	9210000	114,70	
<i>minima</i>	63500000	1340000	<b>276000000</b>	<b>434,65</b>	Relevant
<i>nuclei</i>	22200000	4050000	<b>34800000</b>	<b>156,76</b>	Relevant
<i>nucleuses</i>	73200	8270	96600	131,97	
<i>phenomena</i>	68900000	10800000	<b>112000000</b>	<b>162,55</b>	Relevant
<i>phenomenons</i>	462000	99900	1210000	261,90	
<i>syllabuses</i>	532000	144000	1500000	281,95	
<i>syllabi</i>	4930000	799000	<b>6080000</b>	<b>123,33</b>	Relevant



<b>strata</b>	6560000	4060000	<b>6650000</b>	<b>101,37</b>	Relevant
<b>stratums</b>	336000	24200	928000	276,19	
<b>vortexes</b>	486000	113000	1020000	209,88	
<b>vortices</b>	2610000	448000	<b>4260000</b>	<b>163,22</b>	Relevant

Sursa: Realizat de autori pentru coloanele referitoare la 2020 și după (Manea, Săvoiu, 2014).

Notă\*: Agregarea în timp a datelor frecvențiale transformă situațiile irelevante în relevante într-o manieră dacă nu promptă cel puțin într-una neașteptată.

Notă \*\*Căutarea **antennae** a returnat 8890000 rezultate fiind devansată de căutarea **antennas** care a returnat 7410000 rezultate, conferind o extensie în uz multiplicată maximal în 6 ani de 260 de ori, respectiv mai mult cu **26091,55 – 100,00 = 25991,55 %**, în raport cu trecutul.

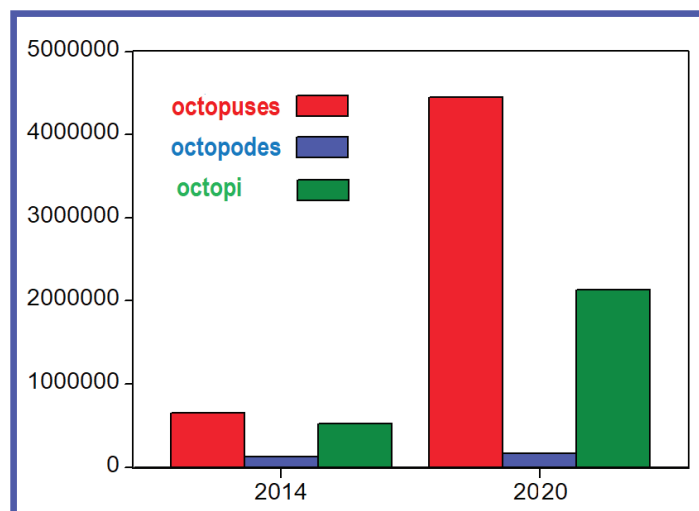
Timpul ucide irelevanța inițială și, așa cum se poate constata, se obțin soluții pentru toate exemplele, și chiar mai mult se pot identifica dinamici diferite conform analizei statistice frecvențiale. Astfel, unele cuvinte sau forme de plural exemplificate și în tabelul 2 și în 3 dețin o relevanță mult mai mare ca urmare a unei explozii în uzul acestora în (mini)corpusul investigat (Internet). Reluând analiza frecvențială în cazul triplului plural potențial: **octopuses, octopodes și octopi** se identifică, din datele finale, trei dinamici complet diferite, care reconfirmă relevanța formei **octopuses**, dar se mai constată și o accelerare a polarizării între formele **octopuses și octopi**, ceea ce oferă dovezi lingvistice ale continuității uzului lor lingvistic accentuat.

Reprezentările statistice grafice ale valorilor frecvențiale comparate în timp în figura 2 contribuie la aprecierea pluralului relevant (corect), printr-o mai bună vizibilitate, acolo unde informația cantitativă este relevantă, detaliind frecvențele absolute, după mai multe investigații realizate în momente diferite (ex. 2014 și 2020).

---

**Analiză statistică grafică de tip comparat în două momente diferite în timp a variantelor unui plural gramatical într-un (mini)corpus lingvistic (Internet)**

Fig. nr. 3



Sursa: Realizat de autori după datele din tabelul 3.

Calculul coeficienților de concentrare–diversificare, de tip HH (Herfindahl–Hirschman) sau  $C_{G-S}$  (Gini–Struck), descrie prin valorile ascendente (în cazul HH) sau care tind tot mai clar spre 1 (în cazul  $C_{G-S}$ ) o concentrare în corpusul lingvistic (Internet) analizat sau într-un (mini)corpus lingvistic specific (o bază de date distinctivă). Limitele unei concentrări în exces sunt conferite de valorile de 0,667 (HH) și 0,409 ( $C_{G-S}$ ) pentru situația în care oferta lingvistică se rezumă la 3 soluții ( $n = 3$ ) așa cum rezultă din calculele realizate în lucrarea publicată de Săvoiu, Crăciuneanu, Țaicu, în 2010.

În tabelul 4 se constată că *octopuses* este apt prin uzul lingvistic cuantificat cu ajutorul analizei statistice frecvențiale să reprezinte o soluție unică de plural, generând valori peste limitele unei concentrări excesive în 2020 (atât după valoarea efectivă a coeficientului HH cât și după aceea a  $C_{G-S}$ ).

**Analiză statistică frecvențială dinamică a pluralului gramatical corect cu ajutorul coeficienților statistici de concentrare - diversificare**

Tabel 4

Varianta	2014	(gi)	(gi) <sup>2</sup>	HH = 0,650  C <sub>G-S</sub> = 0,367	2020	(gi)	(gi) <sup>2</sup>	HH = 0,732  C <sub>G-S</sub> = 0,551
<i>octopuses</i>	651,000	0,501	0,251		4,450,000	0,660	0,436	
<i>octopodes</i>	124,000	0,096	0,009		162,000	0,024	0,001	
<i>octopi</i>	523,000	0,403	0,163		2,130,000	0,316	0,099	
Total	1,298,000	1,000	0,423	6,742,000	1,000	0,536		

Sursa: Realizat de autori

Metoda statistică anterioară motivează frecvențial riguros alegerea formei corecte *octopuses*, conform uzului lingvistic dominant (cu valori structurale  $g_i$  situate peste 0,5). Totodată, această metodă validând un proces de concentrare în exces, garantat de cuantificările coeficientului Gini-Struck ( $C_{G-S}$ ) care epășesc valori de 0,41 (Săvoiu, Crăciuneanu, Țaicu, 2010, pp. 15-27) conduce în final la ideea că analiza frecvențială completată cu analiza concentrării – diversificării asigură o credibilitate mult mai mare decât aplicarea exclusivă a analizei frecvențiale simple în filologie în investigarea variației uzului lingvistic în timp și spațiu.

O altă valorificare a analizelor statistice frecvențiale se poate realiza pentru identificarea sursei de unde a fost preluat un neologism, cu referire la literatura științifică ce l-a consacrat prin utilizare extinsă sau maximală. Pentru a exemplifica acest aspect, în tabelul 5 sunt sintetizate aparițiile unui neologism în limba română în paralel cu sinonimele (aproape omonime) din limbile franceză și din engleză:

**Analiză statistică frecvențială a sursei neologismelor, preluate în română din franceză sau engleză, în cadrul unui (mini)corpus lingvistic (Internet)**

Tabel 5

Limba română	Limba franceză	Limba engleză	Sursa (abrev.)
<i>fenomen</i> 30100000	<i>phénomène</i> 62300000	<i>phenomenon</i> 647000000	eng.
<i>teoremă</i> 12800000	<i>théorème</i> 3370000	<i>theorem</i> 60100000	eng.
<i>dilemă</i> 29600000	<i>dilemme</i> 3480000	<i>dilemma</i> 101000000	eng.

Sursa: Realizat de autori

Cu siguranță că etimologia sau originea greacă a celor trei cuvinte exemplificate nu poate fi pusă în discuție, dar analiza frecvențială subliniază că acești termeni științifici deși au intrat la noi prin limba franceză, conform uzului lingvistic aparțin cu siguranță de literatura științifică engleză. În mod evident, nu poate fi abandonată etimologia sau originea cuvintelor, deoarece se poate astfel ajunge astfel la a deveni nu informat ci dimpotrivă prost informat. În cazul descris în tabelul 5, analizei statistice frecvențiale nu i se pot aduce nici acuzații de rea informare și cu atât mai puțin de dezinformare, aceasta probând cantitativ valorificarea limbajului științific într-o limbă sau în alta.

Cât de mare este utilizarea unor expresii incorecte după normele aparent în vigoare, fără să uităm că limba și limbajul dețin o vitalitate extraordinară și când începe ca dinamica lor să depășească pe aceea a formelor considerate corecte? La o astfel de întrebare, tot o analiză cantitativă statistică de tip frecvențial găsește răspuns pornind de la echilibrul paretian 20/80 și credem că, de îndată ce forma incorectă ( $F_{INC}$ ) depășește 20 % din utilizările totale agregate și deține un *indice evolutiv al relevanței aprecierii* mai mare decât cea declarată corectă ( $F_{COR}$ ), o competiție evidentă se instalează în uzul celor două. În tabelul 6 sunt prezentate două exemple ale valorificării metodei analizei statistice frecvențiale în identificarea unor competiții lingvistice referitoare la uzul alternativ (corect-incorect) dar și unul de concurență neinstalată între o formă consacrată și alte două forme considerate toate împreună ca fiind corecte:

**Analiză statistică frecvențială a raportului paretian sau non paretian între forme corecte și incorecte ale aceluiași cuvânt într-un (mini)corpus filologic (Internet)**

*Tabel 6*

Frecvența $F_{COR}$	Frecvența $F_{INC}$	Uz agregat	Raport $F_{INC}/F_{COR}$
<i>Președinție</i> 819000	<i>Președenție</i> 78300	897300	8,7% vs 91,3%
<i>Să aibă</i> 25300000	<i>Să aivă + Să aibe</i> 2780000	28080000	9,9% vs 90,1%
<i>Găluști</i> 45700	<i>Găluște</i> 387000	*432700	10,6% vs 89,4%

Sursa: Realizat de autori.

\*Notă : Ambele forme sunt considerate corecte, dar chiar și așa nu s-a depășit limita de optim paretian care ar fi putut deschide o competiție lingvistică reală legată de uzul lor alternativ.

Argumentația care justifică folosirea raportului optim de tip paretian pornește de la conceptul de variație și implicit de la cel derivat de omogenitate în cazul variabilelor alternative sau binare, unde limita coeficientului de

omogenitate se atinge în evaluările de tip 20% cu 80% (DA vs NU). La valori mai mici de 80% comparative cu limita paretiană acestea probează posibilități certe de eterogenitate în cadrul oricărei populații, inclusive a unui (mini) corpus filologic.

Metoda profilului statistic, mai ales în varianta clasică de confruntare a unor profile statistice deține o aplicabilitate lingvistică sau filologică potențială destul de mare. Beneficiind de o mare putere de sinteză, metoda confruntării prin profile statistice poate sintetiza dialogul și argumentația filologică în mod esențial, așa cum se poate constata din tabelul nr. 7.

### **Confruntarea filologică a profilurilor statistice realizată în același (mini) corpus filologic**

*Tabel 7*

<b>Profilul statistic fals al structurii limbii române realizat de A. de Cihac și Sextil Pușcariu plagiat după lingvistul rus I. Hinculov</b>	<b>Profilul statistic al structurii limbii române realizat obiectiv și echilibrat al lui Mircea Seche</b>
Cuvinte de origine latină: 20,2%	Cuvinte de origine latină: 45,6%
Cuvinte de origine slavă: 41,0%	Cuvinte de origine slavă: 34,8%
Cuvinte de origine turcă: 16,7%	Cuvinte de origine turcă: 7,1%
Cuvinte de origine neogreacă: 11,0%	Cuvinte de origine neogreacă: 6,2%
Cuvinte de origine maghiară: 10,2%	Cuvinte de origine maghiară: 5,8%
Cuvinte de alte origini: 0,9%	Cuvinte de alte origini: 0,5%

Sursa: Realizat de autori cu rol sintetic și de confruntare finală.

Există, evident, mult mai multe metode statistice aplicabile în filologie sau lingvistică, ceea ce conferă un spațiu cu mult mai larg de mișcare crosdisciplinarității și chiar și multidisciplinarității, iar studiile și cercetările prezente și viitoare vor dovedi cu siguranță acest adevăr.

### **5. Concluzii**

Articolul este rezultatul aplicării crosdisciplinare a metodelor statistice, autorii fiind călăuziți de dorința de a identifica soluții utile, fie ele uneori și exclusiv parțiale, la câteva probleme controversate, cu rezultate obiective și fiabile, care să confirme sistematica logică și analogică a funcționării limbii naturale. Orice analiză frecvențială de tip statistic aplicată adecvat în corpusuri și (mini)corpusuri filologice conduce la final la sprijinirea utilizatorului comun al limbii, ale cărui capacități multiple de asociere anticipată, de comparare simplă cu scop de ierarhizare și de confruntare cu scop de eliminare, pot fi permanent optimizate oferind soluțiile lingvistice

---

probate și dovedind prin argumente statistice uzul practic al acestora. În opinia autorilor, această direcție de cercetare crosdisciplinară poate aduce cu sine clarificări suplimentare privind adecvarea metodelor statistice la investigația filologică, prin valorificarea directă a instrumentelor frecvențiale, a indicilor evolutivi ai relevanței aprecierii uzului lingvistic, a coeficienților de concentrare - diversificare pentru a identifica pluralul dominant al unor anglicisme, sursele de unde au fost preluate unele neologisme care dețin deja un caracter internaționalizat sau intensitatea competiției dintre formele corecte și încorecte ale aceluiași cuvânt în uzul lingvistic (Săvoiu, et al., 2020, p. 146)

Autorii își propun să continue investigațiile realizate cu altele similare prin crosdisciplinaritate, și chiar să inițieze într-o echipă mai largă ce tinde spre multidisciplinaritate, cercetări mai extinse prin multiplicarea metodelor, simultan cu extensia gradului de acoperire prin corpusuri de dimensiuni tot mai mari. Pe termen scurt, merită analizate, prin prisma uzului și a frecvenței de utilizare, variantele semantice hibride („*barbarisme*”) care apar frecvent în limba română, cel mai adesea împrumutate direct din limba engleză ori cel puțin calchiate după modele anglo-americe, cum ar fi: *oneros, intrepid, vocal, versatil*. Din păcate, nu se pot identifica nici chiar pe termen mediu soluții optime de cercetare a modalităților distincte de pronunțare a unor astfel de termeni căutați, întrucât pronunția/fonetica nu poate fi înregistrată în textele din corpusurile mari cu acces universal de tip Internet (Săvoiu, et al., 2020, p. 147)

Singura viabilă în viitorul apropiat rămâne așadar studiarea *metodică* a (mini)corpusurilor filologice.

Ce lucruri utile ar mai putea intra în vizorul analizelor statistice frecvențiale? Probabil că se pot identifica unele proiecte de lucru în echipă multidisciplinară în domeniul standardizării/ normării și al didacticii de tip comparativ-contrastiv devenite tot mai utile sau relevante. Analog simpla constituire a unor corpusuri de date, în principal, prin căutări efectuate pe multitudinea de texte găzduite de Internet, ar putea dovedi prompt și concret, care este situația reală a unor cuvinte care prezintă dificultăți de uz, cum se face acordul verbului/predicativului cu subiectul sau folosirea unor prepoziții sau conjuncții, precum și prevalența sau nu, în uzul real al limbii, a anumitor forme de plural declinate încorecte dar evolutive sau dinamice etc. Pornind de la modelul dicționarului compus de J. C. Wells, se pot realiza comparații lingvistice sau filologice internaționale a unor cuvinte românești cu a celor similar ca sens, folosire de exemplu în limba engleză, tot prin evaluare statistică frecvențială, dar și structurări sau ponderări ale unor cuvinte ce încep cu o anumită literă de dicționar, precum și cât la sută dintre substantive sau verbe au forme sau utilizări neregulate sau considerate „*aberante*”, pornind de la o riguroasă cuantificare statistică cu softuri tot mai performante.

---

O carte apărută destul de recent în cadrul unor cercetări crosdisciplinare și pe alocuri chiar multidisciplinare (Ross, Greenhill, Atkinson, 2013), unde s-au folosit simultan și genetica dar și istoria poveștilor populare a expus un fapt interesant cu consecințe legate de dinamica teritorială a uzului limbajului, constatând că poveștile sunt grupate geografic asemănător sau chiar similar genelor. Acest aspect ar putea permite determinarea ca areal în uz a limbii și a granițelor ei etnico-lingvistice, pornind pe urmele variantelor poveștilor, în paralel cu evoluția ADN-ului, impactul limbii și al poveștilor fiind mai mare decât acela al genelor sau cum exprima sintetic jurnalistă austriacă Christine Kenneally, în 2019: “*un cuplu de oameni își pot mai ușor amesteca genele fără a împărtăși aceeași limbă, cu mult mai ușor decât o poveste poate să treacă peste o barieră lingvistică*” (Kenneally, 2019, p. 385).

“*Cultura și limba sunt și rămân solide în esența lor, în timp ce genele sunt lichide (Khan, 2013), iar o metodă statistică poate studia aparent mai ușor un fenomen stabil sau momentul (...), dar nu trebuie să uite să o facă și de o manieră [metodologică] dinamică (la intervale periodice), deci cu intenții generatoare sau dătătoare de vitalitate și curgere filosofică (...)*” (Săvoiu, et al., 2020, p. 150)

#### Bibliografie

1. Graur, A., 1989. *Iarăși â și î*, în *Puțină gramatică*, vol. II, București: Editura Academiei R.S.R.
2. Hașdeu, B. P. 1984. *Cuvente den bătrâni*, București: Editura Didactică și Pedagogică,
3. Hristea, T. (coord.) 1984. *Sinteze de limba română*, București: Editura Albatros
4. Iordan, I., 1978. *Istoria lingvisticii românești* (Coordonator: acad. Iorgu Iordan), București: Editura Științifică și Enciclopedică.
5. Iordan, I., Robu, V. 1978. *Limba română contemporană*, București: Editura Didactică și Pedagogică, București.
6. Kenneally, C., 2019. *Povestea secretă a speciei umane, Cum ne sunt modelate identitatea și viitorul de ADN și de istorie*, București: Editura Humanitas
7. Khan, R. 2013. Why Culture is Chunky and genes are Creamy, *Gene Expression*.
8. Lombard, A., 1992. Despre folosirea literelor â și î, *Limba română*, nr. 10. pp. 531-540.
9. *MacMillan Dictionary*, MacMillan, 2012
10. Macrea, D., 1943. Fizionomia lexicală a limbii române, *Dacoromania*, vol 10/2. pp. 362-373.
11. Manea, C. 1993. *Considerații statistice asupra ortografiei lui â din a și î din i*, comunicare prezentată la Sesiunea de comunicări științifice „Rolul presei și învățământului în relansarea economico-socială a României”, de la Universitatea din Sibiu, 21-22 mai 1993.
12. Manea, C. 1997. Structura etimologică a limbii române literare contemporane în lumina frecvenței cuvintelor, *Studii și cercetări lingvistice*, anul XVIII, nr. 2.

- 
13. Manea, C. 2009. *Încercare statistică asupra ortografierii cu â și î*, în vol. *Distorsionări în comunicarea lingvistică, literară și etnofolclorică românească și contextul european*, Academia Română, Institutul de Filologie Română „A. Philippide”, Iași: Editura Alfa, pp. 209-218
  14. Manea, C., 2004. *Structura etimologică a vocabularului neologic (cu specială referire la anglicismele din limba română)*, Pitești: Editura Universității din Pitești.
  15. Manea, C., 2009. *Încercare statistică asupra ortografierii cu â și î*, în vol. *Distorsionări în comunicarea lingvistică, literară și etnofolclorică românească și contextul european*, Academia Română, Institutul de Filologie Română „A. Philippide”, Ed. Alfa, Iași, pp. 209-218.
  16. Manea, C., Săvoiu, G. 2014. *Frequency Analysis of the Use of a Number of Morphological Variants in Scientific Language*, in ESMSJ, vol. IV, no. 2 (special issue), pp. 13-17.
  17. Marcu, F. 1997. *Noul dicționar de neologisme*, București: Editura Academiei Române.
  18. Marcu, F. 2000. *Dicționar uzual de neologisme*, București: Ed. Saeculum I.O
  19. O’Keeffe, A., M., McCarthy, J., and Carter, R. A. 2007. *From Corpus to Classroom*, Cambridge: Cambridge University Press.
  20. Onu, L., 1992. Oportunitatea reformei ortografice, *Limba română*, nr. 4, pp. 199-207, București.
  21. Săvoiu, G., 2003. *Statistică generală. Argumente în favoarea formării gândirii statistice*, Pitești: Editura Independența economică.
  22. Săvoiu, G. (coord.), Asandei, M., Grigorescu, R., Manole, S. 2005. *Cercetări și modelări de marketing. Metode cantitative în cercetarea pieței*, București: Editura Universitară.
  23. Săvoiu, G., 2012. *Statistică generală cu aplicații în contabilitate*, București: Editura Universitară.
  24. Săvoiu, G., Crăciuneanu, V. Țaicu, M. 2010. A New Method of Statistical Analysis of Markets’ Concentration or Diversification. *Romanian Statistical Review*, vol. 58(2), pp. 15-27.
  25. Săvoiu, G., et al. 2020. *Metode statistice și crosdisciplinaritate*, București: Editura Universitară.
  26. Sesardić, N., 2018. *Când rațiunea pleacă în vacanță. Filozofii în vacanță*, București: Editura Humanitas.
  27. \*\*\* *Biblia sau Sfânta Scriptură*. 2001. Ediție jubiliară a Sfântului Sinod, versiune redactată și adnotată de Bartolomeu Valeriu Anania (ed.), București: Editura Institutului Biblic și de Misiune al Bisericii Ortodoxe Române, p. 22. Available on - line la: <http://www.diacronia.ro/en/indexing/details/B347/pdf>. Accessed on 30.04.2020.