

Data mining in the age of social web

Ákos Bardóczi

Downloads: <https://bardoczi.net>

Web: <http://genetics.dote.hu/bardoczi>

Email: bardoczi@med.unideb.hu

Tel: +1 202 4700-790

Fax: +36 1 9997-993

Comments: akos@cryptolab.net



*If You Steal From One Author, It's Plagiarism;
If You Steal From Many, It's Research*

unknown author



Highlighted topics

- a part of open source intelligence
- simple data extraction/mining techniques from web 2.0 services
- authorship attribution/authorship profiling
- plagiarism detection



Old fashioned sources in nutshell

- WHOIS-records
- email long headers – which sent via web, commercial providers
- DNS-records
- patterns in cookies and HTTP-response streams



Old fashioned sources – WHOIS

- in case gTLD – WHOIS masking
- the fake contact details similars to each other between suspected domains - frequent
- private registrar companies owned by domain owner



Old fashioned sources email long header

- The popular misbelief about Gmail, AOL, YahooMail and others



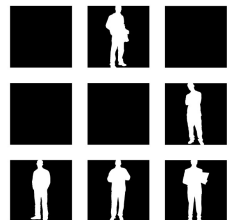
Old fashioned sources email long header in fact

Return-path: <bardoczi@gmail.com>
Received: from st11p00mm-smtpin012.mac.com ([17.172.87.212]) by ms03531.mac.com (Oracle Communications Messaging Server 7u4-27.08 (7.0.4.27.7) 64bit (built Aug 22 2013)) with ESMTP id <0MUC007OQVDVDMG0@ms03531.mac.com> for bardoczi@me.com; Tue, 08 Oct 2013 15:08:19 +0000 (GMT)
Original-recipient: rfc822;bardoczi@me.com
Received: from mail-we0-f182.google.com ([74.125.82.182]) by st11p00mm-smtpin012.mac.com (Oracle Communications Messaging Server 7u4-27.08(7.0.4.27.7) 64bit (built Aug 22 2013)) with ESMTP id <0MUC00MUCVDUUTJ1@st11p00mm-smtpin012.mac.com> for bardoczi@me.com (ORCPT bardoczi@me.com); Tue, 08 Oct 2013 15:08:19 +0000 (GMT)
Received-SPF: pass (st11p00mm-smtpin012.mac.com domain of bardoczi@gmail.com designates 74.125.82.182 as permitted sender) helo=mail-we0-f182.google.com; envelope-from=bardoczi@gmail.com; x-software=spfmlter 0.97 http://www.acme.com/software/spfmlter/ with libspf-unknown;
Received: by mail-we0-f182.google.com with SMTP id t61so8978535wes.13 for <bardoczi@me.com>; Tue, 08 Oct 2013 08:08:17 -0700 (PDT)
DKIM-Signature: v=1; a=rsa-sha256; c=relaxed/relaxed; d=gmail.com; s=20120113; h=mime-version:reply-to:date:message-id:subject:to:content-type; bh=F7h8uU2A3blCvcj6XbYiBDxawl1FFa98TBwa/d1NMpM=; b=zE5NNhuRYbJ1pF6FVzzV3sUlwYShaPj/b5aFpEAY4geEe0Avi+nYh45jZI66gaFZg/FWm/WrVXOypzB6BDklqvQzXnYSq8XbP6BIJY4W4calkByv0Gi6RS5RSECV+goLCul4DTRHfcKyKNCBoNSSL1IQdijZaRNL/BrK0oNm5rhChSiepYurr4uqnkeUiHk5hpr6W4+Et8ZC3J1vpa7l4Ln4fWviG6p63olwZFMMrmvghOuH3lenzRPZCI636z5PvwXYLorwWJc7gDkhshPztIjKimUtWclm3Wgc50we+VEre3QKy+4rd8bEhSTlrXzJ2MUJtef27M3cbfZZPr pWYA==
MIME-version: 1.0
X-Received: by 10.180.93.166 with SMTP id cv6mr2026676wib.37.1381244897844; Tue, 08 Oct 2013 08:08:17 -0700 (PDT)
Received: by 10.180.93.166 with SMTP id cv6mr2026676wib.37.1381244897844; Tue, 08 Oct 2013 08:08:17 -0700 (PDT)
Reply-to: akos@secure.bardoczi.net
Date: Tue, 08 Oct 2013 17:08:17 +0200
Message-id: <CAFMsfx1G4uXMG-B_1P0ihvYBhMRe3Zc8diQxbkoRXEtDM4TSEw@mail.gmail.com>
Subject: Nem mondom meg, hogy hol vagyok
From: ?ISO-8859-1?Q?E1kos_bard=F3czi?= <bardoczi@gmail.com>
To: bardoczi@me.com
Content-type: multipart/alternative; boundary=f46d043c7fbc7befc104e83c23ec
Authentication-results: st11p00mm-smtpin012.mac.com dkim=pass reason="2048-bit key" header.d=gmail.com header.i=@gmail.com
Character encoding: UTF-8
X-icloud-spam-score: 33322 f=mail.com e=gmail.com pp=ham sof=pass dkim=pass wl=absent pwl=absent
X-Proofpoint-Spam-Details: rule=notspam policy=default score=1 spamscore=1 suspectscore=3 phishscore=0 adultscore=0 bulkscore=0 classifier=spam adjust=0 reason=mlx scancount=1 engine=7.0.1-1308280000 definitions=main-131008006

the server which first received

sender's timezone

character encoding – may indicates
the senders language and nationality



HACKTIVITY

Old fashioned sources

„client IP”

Tracing route to mail-we0-f182.google.com [74.125.82.182] over a maximum of 30 hops:

1	17 ms	2 ms	*	192.168.4.1
2	17 ms	17 ms	17 ms	192.168.0.1
3				
4	11 ms	26 ms	17 ms	catv-89-135-222-30.catv.broadband.hu [89.135.222.30]
5	17 ms	17 ms	17 ms	84.116.240.78
6	17 ms	17 ms	40 ms	84.116.240.1 - UPC Austria GmbH
7	17 ms	17 ms	14 ms	72.14.217.146 - Google Inc
8	27 ms	17 ms	41 ms	209.85.243.121
9	47 ms	48 ms	48 ms	72.14.234.11
10	40 ms	55 ms	63 ms	209.85.241.228
11	52 ms	62 ms	56 ms	209.85.240.221
12	62 ms	58 ms	39 ms	209.85.252.83
13	*	*	*	Request timed out.
14	50 ms	58 ms	43 ms	mail-we0-f182.google.com [74.125.82.182]

Trace complete.



Old fashioned sources - emailhipotetical interpretation

- the Google doesn't make transatlantic traffic if not necessary!
- typical ping times between Europe and USA :)
- „de novo” signal
- syllogism: the server and the sender both located near Vienna +/- 200 km
- (dedicated connection across Google and ISP?)



Old fashioned sources

DNS-records

- non-masked A-record – note:
entire IP-block may owned by one
server/owner as „hosting provider”
- masked A-record – eg.

Cloudflare:

- by default adding direct-
connect subdom.
- cannot masking MX – reverse
lookup!



Old fashioned sources cookies and HTTP-responses

- same owner - same developer - same admin - same CMS - frequent
- consequence – similar HTTP-stream patterns and cookie-patterns between suspected websites



Different layers, different methods

...and cognitive/linguistical/behavioural

OSI LAYERS	EXAMPLE PROTOCOLS
APPLICATION LAYER	HTTP, FTP, IRC, SSH, DNS
PRESENTATION LAYER	SSL, FTP, IMAP, SSH
SESSION LAYER	VARIOUS API'S, SOCKETS
TRANSPORT LAYER	TCP, UDP, ECN, SCTP, DCCP
NETWORK LAYER	IP, IPSec, ICMP, IGMP
DATA-LINK LAYER	Ethernet, SLIP, PPP, FDDI
PHYSICAL LAYER	Coax, Fiber, Wireless

Question:
how can examine...
...and why?

Why we loves the Facebook API?

Facebook Graph API (2013. sept.)

- Graph API object set: Achievement, Album, Application, Checkin, **Comment**, Domain, Errors, Event, FriendList, Group, Insights, Link, Message, **Note**, Offer, Order, Page, Payment, Photo, Pictures, **Post**, Question, QuestionOption, Review, **Status message**, Thread, User, Video

- The fields of the User object:

id, name, first_name, middle_name, last_name, **gender**, locale, languages, link, username, **age_range**, third_party_id, installed, timezone, updated_time, verified, **bio**, birthday, cover, currency, devices, education, email, hometown, interested_in, location, political, payment_pricepoints, favorite_athletes, favorite_teams, picture, quotes, relationship_status, religion, security_settings, significant_other, video_upload_limits, website, work



Why we loves social sources?

Some methods not needs special permissions, API keys, etc. Eg. this simple URL:

graph.facebook.com/https://facebook.com/hacktivity?metadata=1

• In real e.g.

<http://www.weknowwhatyouredoing.com/>

We know what you're doing...

a social networking privacy experiment



32,064 people like this. Be the first of your friends.



6,936



Follow @callumhaywood

Public Facebook statuses - Status Search - Foursquare location finder - Facebook friend checkins - Contact

[About this tool](#)

Who wants to get fired?



Jesse Alber W.

I fucking hate being off work but have to listen to my boss get my eyes examed this freaking blows big time about 58 minutes ago, no people like this, posted from Facebook for iPhone, report

Who's hungover?



Deb C.

Pouring rain, a hungover passenger, not even close to Minot yet...I am breaking all the rules and hitting Starbucks as soon as I pull into Minot! about 58 minutes ago, no people like this, posted from Facebook for Android, report

Who's taking drugs?



Deonte J.

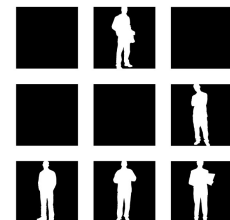
Woke up wit ah headache smoke ah blunt headache gone#loveweed#jakefrosteBlunt about 1 hour ago, no people like this, posted from Facebook for Android, report

Who's got a new phone number?



Matthew T.

Got a new phone, same number (07xx611x4x). I'll try gets whatsapp up and running asap :) about 2 hours ago, no people like this, posted from web, report



HACKTIVITY

The rabbit hole

one of my fav. studies

Private traits and attributes are predictable from digital records of human behavior (Michal Kosinskia, David Stillwella, Thore Graepelb, 2012)

- input: 58000 volunteers likes, some demographical and psychometrical data
- Processed with simple linear regression and logistics regression



The found, successfully like-correlated properties:

“sexual orientation,” “ethnic origin,” “political views,” “religion,” “personality,” “intelligence,” “satisfaction with life” (SWL), substance use (“alcohol,” “drugs,” “cigarettes”), “whether an individual’s parents stayed together until the individual was 21 y old,” and basic demographic attributes such as “age,” “gender,” “relationship status,” and “size and density of the friendship network.”



Private traits and attributes are predictable from digital records of human behavior (Michal Kosinskia, David Stillwella, Thore Graepelb, 2012)

The researcher...

Trait		Selected most predictive Likes		
IQ	High	The Godfather	Jason Aldean	Low
		Mozart	Tyler Perry	
		Thunderstorms	Sephora	
		The Colbert Report	Chiq	
		Morgan Freemans Voice	Bret Michaels	
		The Daily Show	Clark Griswold	
		Lord Of The Rings	Bebe	
		To Kill A Mockingbird	I Love Being A Mom	
		Science	Harley Davidson	
		Curly Fries	Lady Antebellum	
Satisfaction With Life	Satisfied	Sarah Palin	Hawthorne Heights	Dissatisfied
		Glenn Beck	Kickass	
		Proud To Be Christian	Atreyu (Metal Band)	
		Indiana Jones	Lamb Of God	
		Swimming	Gorillaz	
		Jesus Christ	Science	
		Bible	Quote Portal	
		Jesus	Stewie Griffin	
		Being Conservative	Killswitch Engage	
		Pride And Prejudice	Ipod	



<http://www.pnas.org/content/suppl/2013/03/07/1218772110.DCSupplemental/st01.pdf>

Private traits and attributes are predictable from digital records of human behavior (Michal Kosinskia, David Stillwella, Thore Graepelb, 2012)

...can...

Openness	<i>Liberal & Artistic</i>	Oscar Wilde Charles Bukowski Sylvia Plath Leonardo Da Vinci Bauhaus Dmt The Spirit Molecule American Gods John Waters Plato Leonard Cohen	NASCAR Austin Collie Monster-In-Law I don't read Justin Moore ESPN2 Farmlandia The Bachelor Oklahoma State University Teen Mom 2	<i>Conservative</i>
Conscientiousness	<i>Well Organized</i>	Law Officer National Law Enforcement Lowfares.Com Accounting Foursquare Emergency Medical Services Sunday Best Kaplan University Glock Inc Mycalendar 2010	Wes Anderson Bandit Nation Omegle Vocaloid Serial Killer Screamo Anime Vamplets Join If Ur Fat Not Dying	<i>Spontaneous</i>
Extraversion	<i>Outgoing & Active</i>	Beerpong Michael Jordan Dancing Socializing Chris Tucker I Feel Better Tan Modeling Cheerleading Theatre Flip Cup	RPGs Fanfiction.Net Programming Anime Manga Video Games Role Playing Games Minecraft Voltaire Terry Pratchet	<i>Shy & Reserved</i>



<http://www.pnas.org/content/suppl/2013/03/07/1218772110.DCSupplemental/st01.pdf>

Private traits and attributes are predictable from digital records of human behavior (Michal Kosinskia, David Stillwella, Thore Graepelb, 2012)

...read between...

Agreeableness	<i>Cooperative</i>	Compassion International Logan Utah Jon Foreman Redeeming Love Pornography Harms The Book Of Mormon Circles Of Prayer Go To Church Christianity Marianne Williamson	I Hate Everyone I Hate You I Hate Police Friedrich Nietzsche Timmy South Park Atheism / Satanism Prada Sun Tzu Julius Caesar Knives	<i>Competitive</i>
		Sometimes I Hate Myself Emo Girl Interrupted So So Happy The Addams Family Vocaloid Sixbillionsecrets.com Vampires Everywhere Kurt Donald Cobain Dot Dot Curve	Business Administration Getting Money Parkour Track & Field Skydiving Mountain Biking Soccer Climbing Physics / Engineering 48 Laws Of Power	
Emotional Stability	<i>Neurotic</i>			<i>Calm & Relaxed</i>
Gender	<i>Female</i>	Tv Fanatic Chiq Gillette Venus Shoedazzle Bebe Proud To Be A Mom Covergirl Wet Seal Aerie By American Eagle Mall World	Modern Warfare 2 ESPN Sportscenter Band Of Brothers Starcraft Deadliest Warrior Dos Equis Red Vs Blue X Games Bruce Lee	<i>Male</i>



<http://www.pnas.org/content/suppl/2013/03/07/1218772110.DCSupplemental/st01.pdf>

Private traits and attributes are predictable from digital records of human behavior (Michal Kosinskia, David Stillwella, Thore Graepelb, 2012)

...the lines likes

Sexual Orientation	Politics	Republican George W Bush John McCain Conservative Rush Limbaugh Sean Hannity Bill Oreilly Positively Republican Sarah Palin Ronald Reagan Glenn Beck	Democrat Joe Biden Speaker Nancy Pelosi Health Care Reform The White House Democrats Barbara Boxer Anthony Weiner Being Liberal Left Action Barack Obama2012 Ted Kennedy
	Homosexual Males	No H8 Campaign Kathy Griffin Kurt Hummel Glee Human Rights Campaign Mac Cosmetics Adam Lambert Ellen DeGeneres Juicy Couture Sue Sylvester Glee Wicked The Musical	Heterosexual Males X Games Nike Basketball Bungie WWE Sportsnation Wu-Tang Clan Foot Locker Shaq Bruce Lee Being Confused After Waking Up From Naps
	Homosexual Females	Girls Who Like Boys Who Like Boys Rupauls Drag Race No H8 Campaign Gay Marriage Human Rights Campaign The L Word Sometimes I Just Lay In Bed And Think About Life Not Being Pregnant Gay Marriage Tegan And Sara	Heterosexual Females Lipton Brisk Yahoo Adidas Originals Foot Locker WWE Inbox 1 Makes Me Nervous Thinking Of Something And Laughing Alone I Just Realized Immature Spells I'm Mature Did You Get A Haircut No It Grew Shorter



<http://www.pnas.org/content/suppl/2013/03/07/1218772110.DCSupplemental/st01.pdf>

Private traits and attributes are predictable from digital records of human behavior (Michal Kosinskia, David Stillwella, Thore Graepelb, 2012)

Twitter-users closer

- **problem:** e.g. linkfarm identification
- **hipothesis:** Let's measure the influence and popularity.



Twitter

Similis simili gaudet

- popular users have more follower than followed
- deliberate connections
- identification of problematic accounts [bots/linkfarms]: measure the popularity and the followers dispersion. Abnormal dispersion may indicate problematic user.
- (cluster analysis)



Twitter

Measure the most popular followers of the suspected user

- An elegant implementation in Python
- sample result:

droidPisti 12

FakeFeri 11

HirhozoHugo 9

Blogcsokor 9

Wannabehirportal 8....



Twitter

Entity analysis in the stream

- *Extracting entities from tweets and performing simple frequency analysis - (Mining the Social Web)*
- Abnormal dispersion simple result:

#kiralywebhely 2834

#legujabbpinamagazin 2830

#ittmindenwareztmegtalasz 2820

#torrentekmindenmennyisegben 2811



Possibilities of IF-IDF measures



About digital humanities and computational linguistics

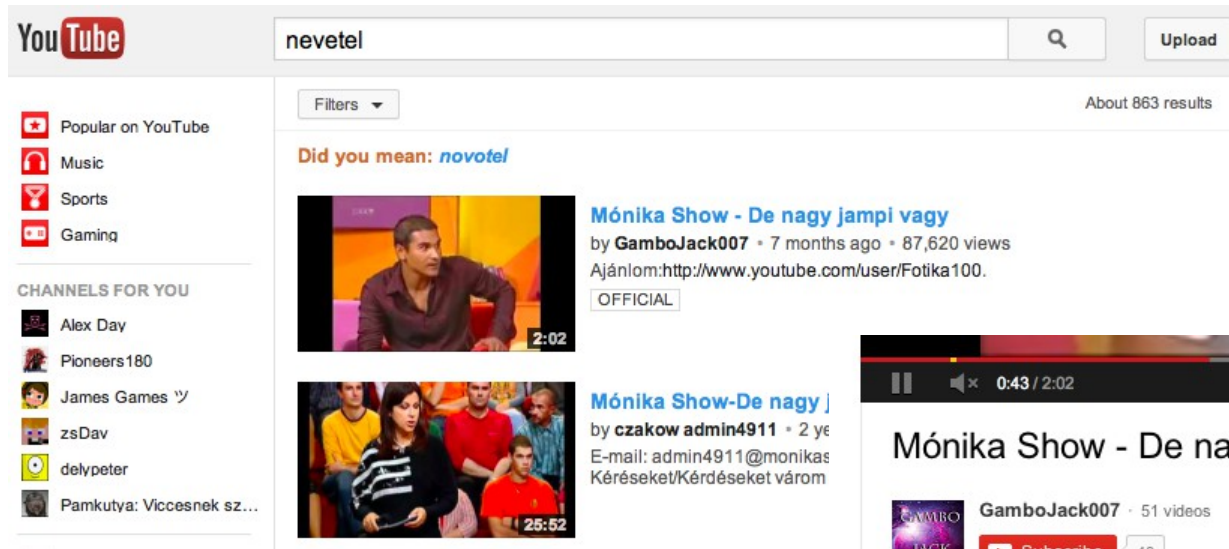
Example: linguistics and the fine tuned search

„To be, or not to be.”

Shakespeare W. (soap opera addict), cca. 1603



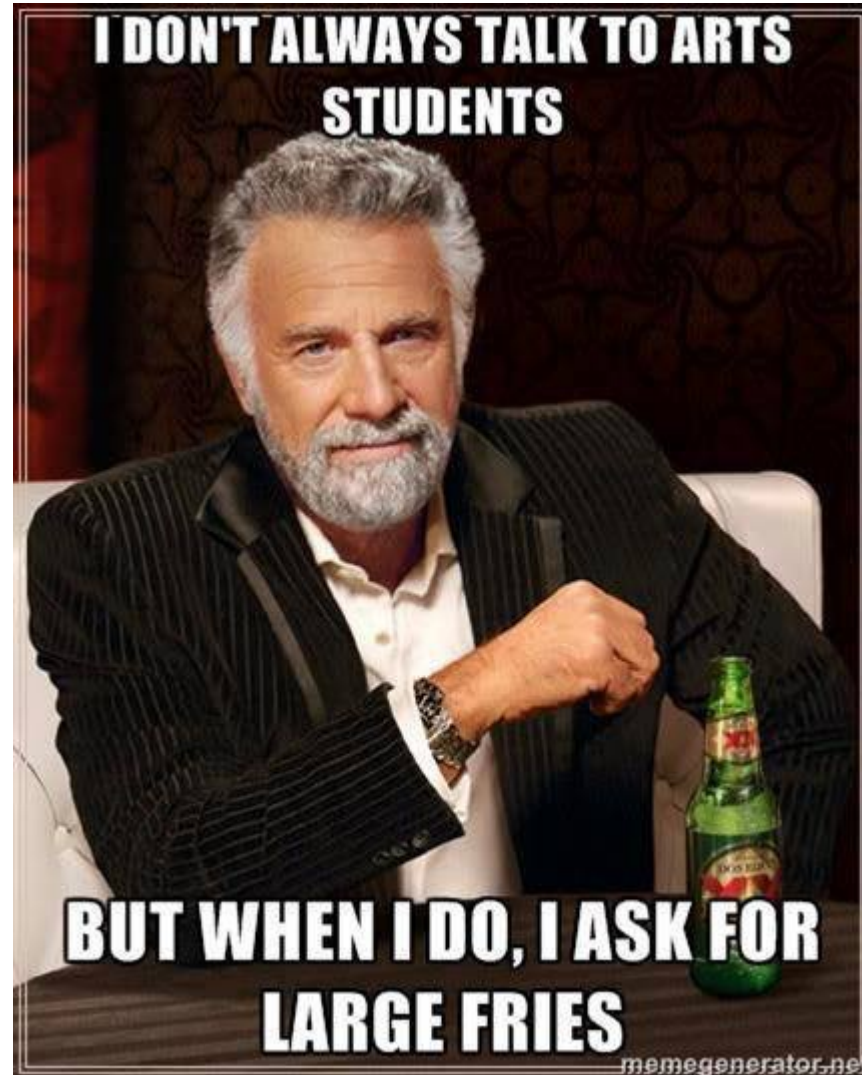
An identified former YT-hapax legomenon and the „did you mean”: *nevetel*-novotel



In this case, the titles of the first two results, don't contains „nevetel”, only the comments (with significant weight)

„Jampi” Zsolti, 2009,
Maunika Erdélyi et al.

Who they are really? - popular misbelief

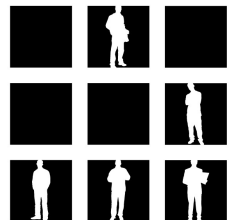
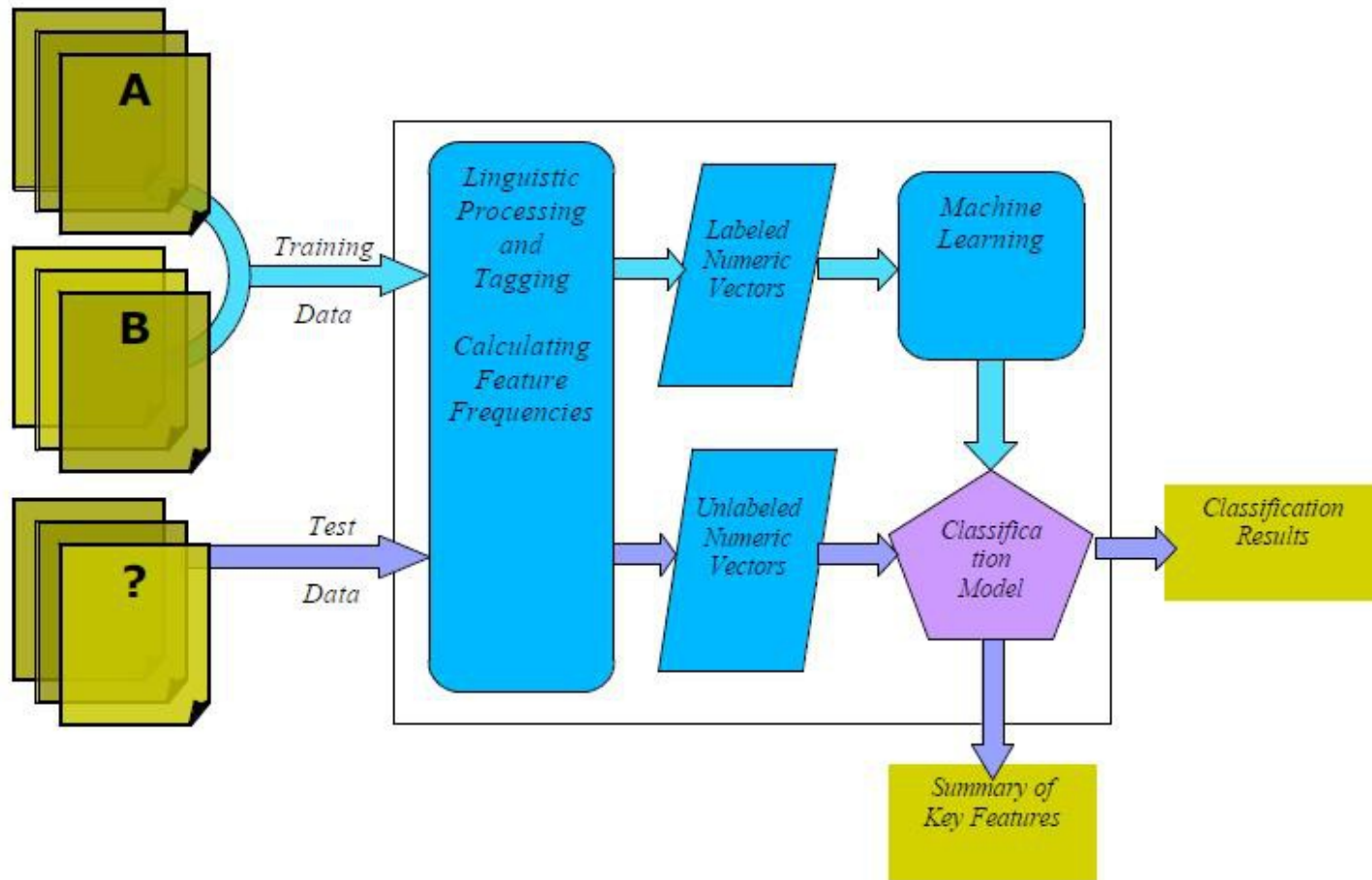


Computational linguists:
they are not simply BAs/MAs.
They must have experiences in...
strong mathematical and computer
science skills, such as

- **classifiers**
- **machine learning**
- **declarative programming**
- **fields of AI**
- **data mining**



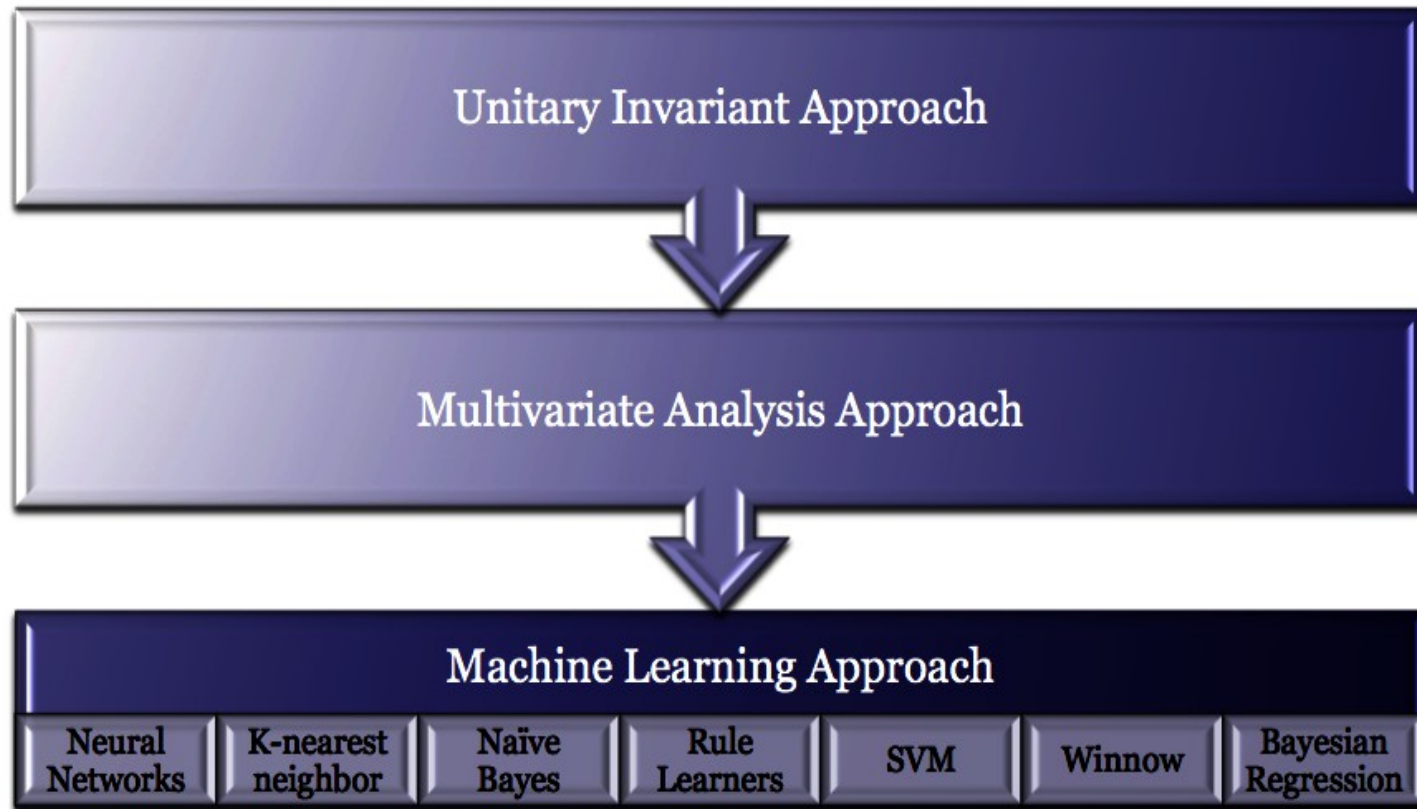
One principle of the many



HACKTIVITY

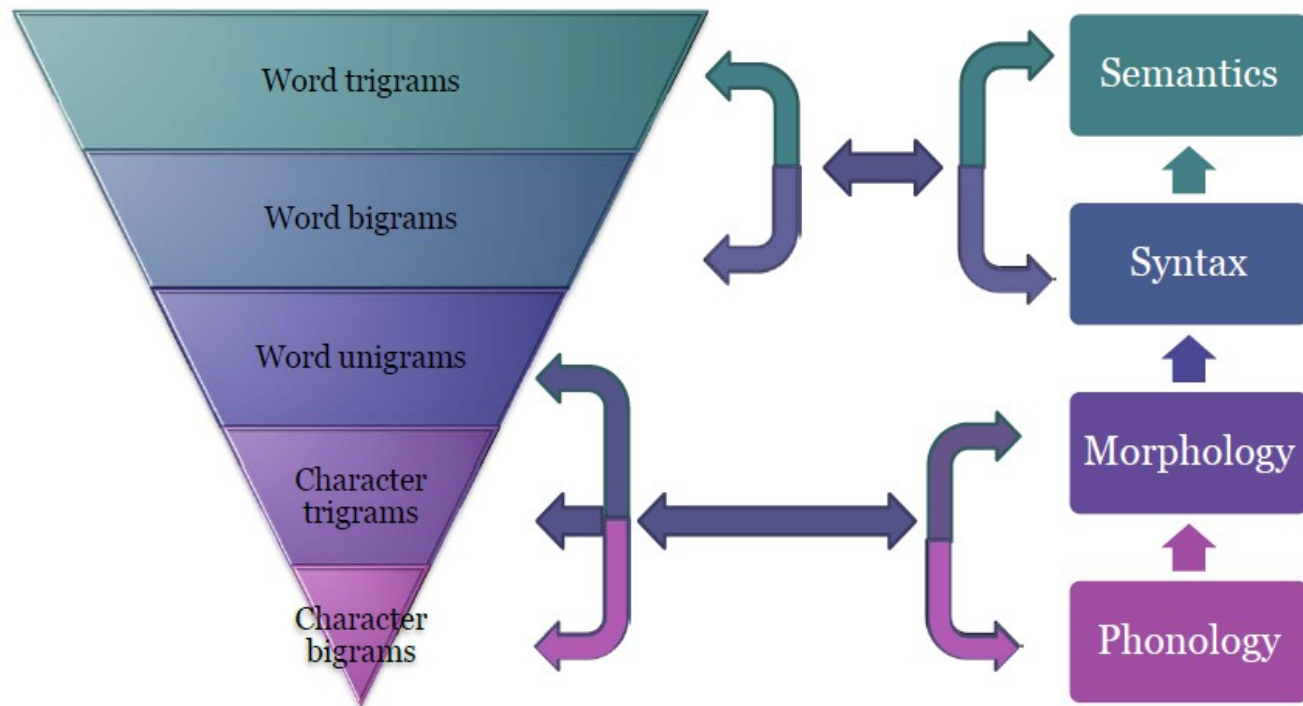
An quick historical overview

History Of Analysis Methods



Authorship identification on big sets

An hierarchical representation of features and related linguistic levels



Authorship identification in large email collections: Experiments using features that belong to different linguistic levels (George K. Mikros & Kostas Perifanos, 2011)

About Natural Language Processing

Python NLTK

Typical:

1. End of Sentence (EOS) Detection
2. Tokenization
3. Part-of-Speech Tagging
4. Chunking
5. Extraction

```
nltk.tokenize.sent_tokenize(valtozo)
tokenek = [nltk.tokenize.word_tokenize(s) for s in sentences]
pos_tagged_tokens = [nltk.pos_tag(t) for t in tokenek]
pos_tagged_tokens
chunking
mondattani elemzés -
ne_chunks = nltk.batch_ne_chunk(pos_tagged_tokens)
ne_chunks
```



Some selected linguistics weapons in authorship identification and plagiarism detection

- frequency and dispersion of suffixed words
- Syntactic complexity
 - Syntactic tree depth
 - Distance between syntactically dependent words
- cohesion

A multitude of linguistically-rich features for authorship attribution (Ludovic Tanguy, Assaf Urieli, Basilio Calderone, Nabil Hathout, Franck Sajous, 2011)



The style of the sentences likes the authors fingerprint!

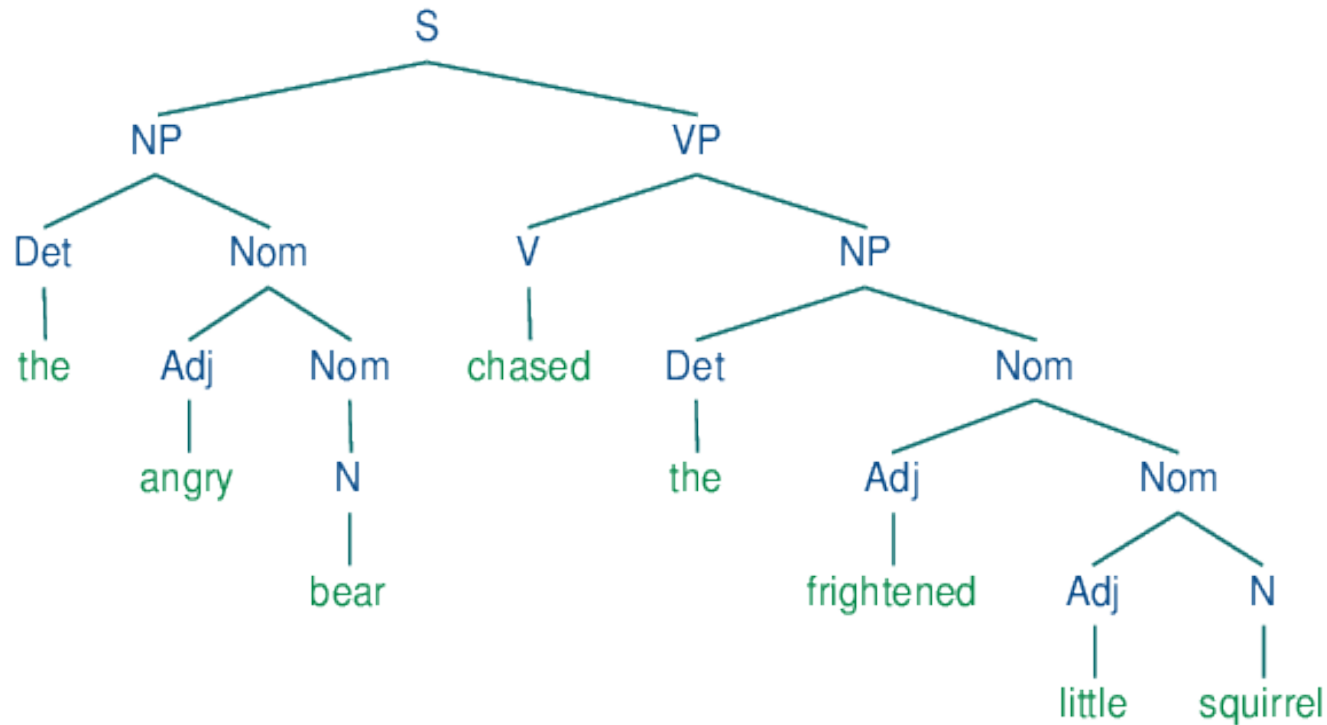


Fig: <http://nltk.org/book/ch08.html>



Quick explanation of a plagiarism detection technique [translation plagiarism]

core problems:

- non-friendly language pairs:
 - a) fixed/non-fixed word order, b) conjugation, c) major grammar differences
- former, partially solutions

Principles of new algo:

- sentence as translation unit
- sentence as context
- sentence, because
 1. represent an unit in mind
 2. easy to find the borders
 3. enough unique
- similarity metrics



Quick explanation of a plagiarism detection technique [translation plagiarism] – algo in nutshell (similarity metrics)

Let's $w_x \in S_x$ and $w_y \in S_y$ $\tilde{\text{Sim}}(x,y) = \alpha \cdot |S_x \cap S_y| - \beta \cdot |S_x \setminus S_y|$
will show the found and not-found words (with different weights).

After a bit fine tuning $\text{Sim}(x,y) = \min (\alpha \cdot |S_x \cap S_y| - \beta \cdot |S_x \setminus S_y| ,$
 $\alpha \cdot |S_y \cap S_x| - \beta \cdot |S_y \setminus S_x|)$

Some lemmas $\text{trans}(w_x) = W_y$ where $w_y \in W_y$

$\text{trans}(w_y) = W_x$ where $w_x \in W_x$

$w_x \in \text{trans}(w_y)$ therefore $w_y \in \text{trans}(w_x)$

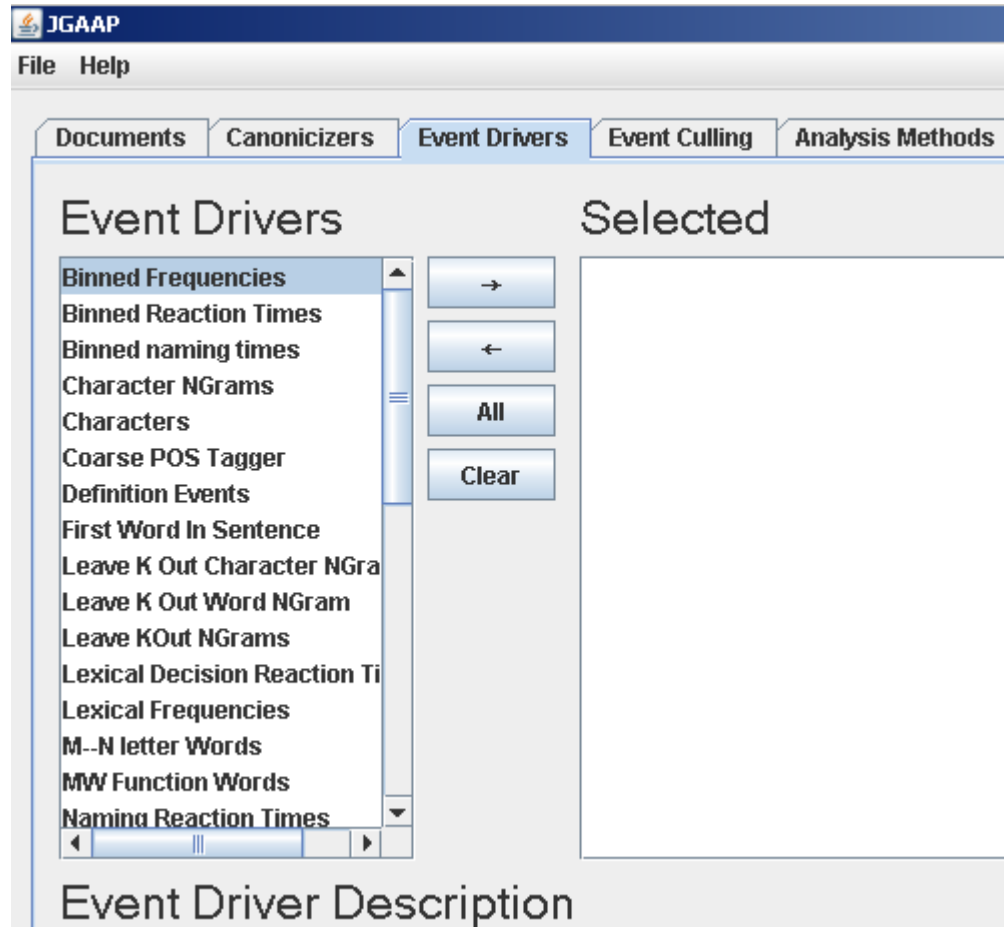
$w_y \in \text{trans}(w_x)$ therefore $w_x \equiv w_y$

Signs: S – sentence, n, w – number of words in sentence and words

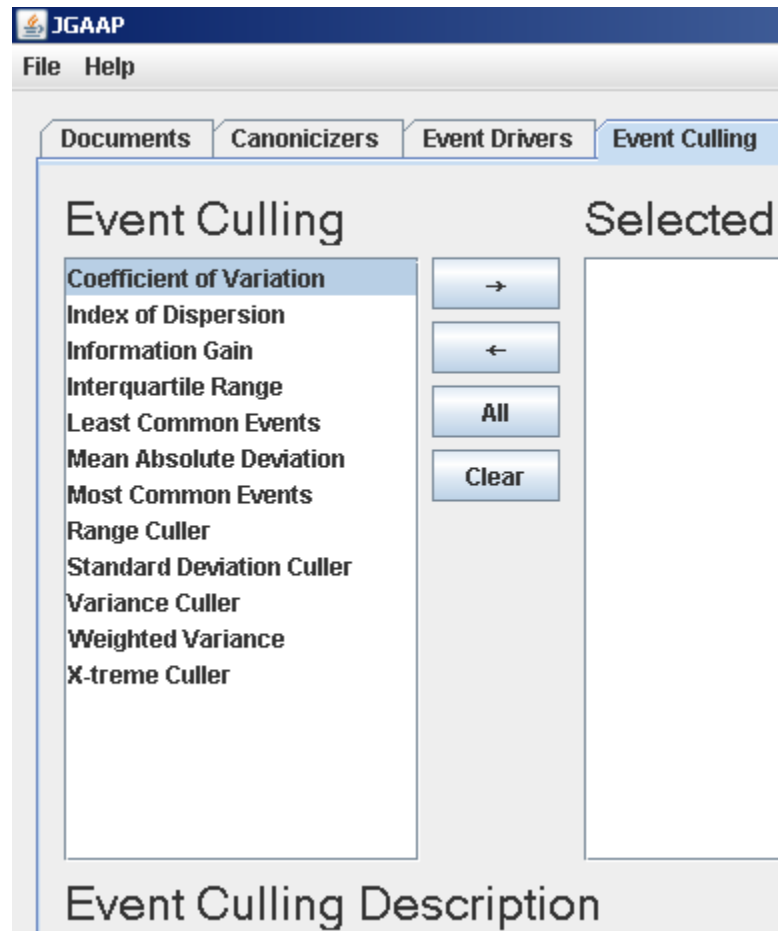


Many-many sources, mostly from Máté Pataki (MTA SZTAKI)

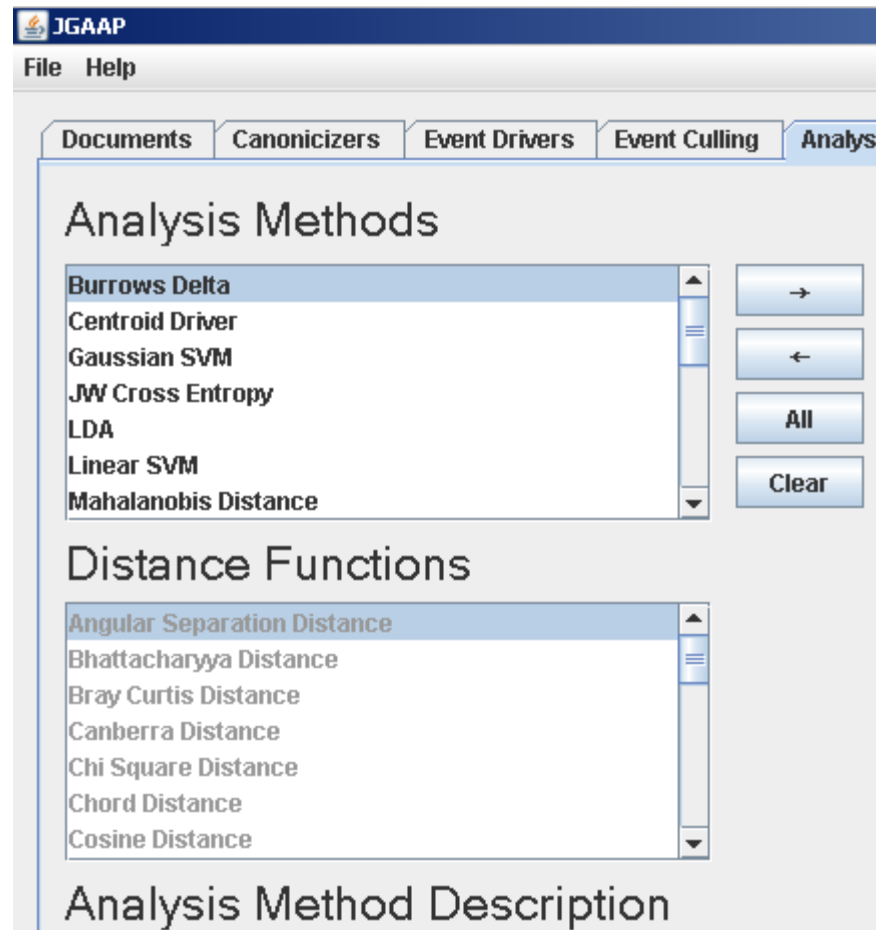
Go deep – JGAAP, event drivers



Go deep – JGAAP, event culling



Go deep – JGAAP, analysis methods



See also

- Text Mining: Applications and Theory (Michael W. Berry, Jacob Kogan)
- Natural Language Processing with Python (Steven Bird, Ewan Klein, Edward Loper)
- Mining the Social Web (Matthew A. Russel)
- Forensic Linguistics - Advances in Forensic Stylistics (Gerald McMenamin et al.)
- The Routledge Handbook of Forensic Linguistics (ed. Malcolm Coulthard, Alison Johnson)

PAN CLEF papers

Recommended blogs in hungarian:

<http://kereses.blog.hu/>

<http://kep-es-kod.blogspot.ch/>

<http://szamitogepesnyelveszet.blogspot.ch/>

- Digital humanities MA in Hungary:

<http://www.btk.unideb.hu/>

<https://btk.ppke.hu/>

