

Forenzní lingvistika

—

jazykověda ve službách vyšetřování

Kateřina Veselovská

veselovska@ufal.mff.cuni.cz

13. října 2017, Ústav jazyků a komunikace neslyšících

Kateřina Veselovsk



senior research associate



analza textovch dat v rmci
forezn divize

Forenzní lingvistika

= aplikovaná interdisciplinární věda na pomezí
lingvistiky, práva a kriminalistiky

Forenzní lingvistika

- v seznamu znaleckých oborů pod oborem kriminalistika
- analýza textu, fonoskopická expertiza (kriminalistická akustika)
- lingvistický materiál jako otisk prstu
- forenzní psychologie, forenzní medicína...
- nutná znalost právního řádu dané země

Forenzní lingvistika

- Ale! i srozumitelnost právních textů – pokus s právy a povinnostmi vyslýchaného
- jak oslovovat prostitutky..?

Forenzní lingvistika

- např. ztotožnění anonymního pachatele:

výhružné dopisy
udavačské dopisy
vyděračské dopisy
- dokazování trestných činů (analýza emailové komunikace)
- dopisy na rozloučenou (korpus PL)
- zasahování do smluv

Forenzní lingvistika – historie

- Německo a Anglosaské země od konce 70. let 20. st.
- u nás při Kriminalistickém ústavu v Praze od 60. let
- IAFL, International Association of Forensic Linguistics
- Malcolm Coulthard – Journal of Forensic Linguistics (od 1994)

Forenzní lingvistika

- vývoj zatím spíše skokový (case studies)
- chybí jednotná metodologie

Forenzní lingvistika

- Jan Svartvik (1968) – případ Rillington Place č. 10
- 30. 11. 1949, Notting Hill, Londýn
- vynucené přiznání Timothy Evanse k vraždě manželky a dcery (odsouzen a pověšen 1950)
- skutečným vrahem sériový vrah John Christie
- metoda: porovnání tří výpovědí, osudná pod nátlakem (delší a syntakticky složité věty)

WHAT HAPPENED TO THE WOMEN AT 10 RILLINGTON PLACE?

THE MOST
SHOCKING
STORY
OF THE
CENTURY!



VICTIM: Pretty Beryl Evans had a medical problem, when she sought help she disappeared without any trace.



10 RILLINGTON PLACE: An ordinary rooming house on an ordinary street hid an extraordinary murderer's secret.



SUSPECT: Timothy Evans' mistake was renting a room at 10 Rillington Place. He would soon be accused of ghostly murder.



VISITOR: Alice, girlfriend of Beryl Evans, Charles's attraction aroused her curiosity and almost led her to a similar fate.



NURSE CHARLIE: She knew a grim secret that was to lead her to an even more grisly death in the house of 10 Rillington Place.



EX-POUCEMAN: John Reginald Christie was a professional witness. Later he would appear in the same court as a defendant.

COLUMBIA PICTURES and FILMWAYS Present
RICHARD ATTENBOROUGH/JUDY GEESON/JOHN HURT.

A MARTIN RANSCHOFF-LESLIE LINDER PRODUCTION **10 RILLINGTON PLACE**

Screenplay by CLIVE EXTON - Associate Producer BASIL APPLEBY - Produced by MARTIN RANSCHOFF and LESLIE LINDER
Directed by RICHARD FLEISCHER • COLO •

- filmový thriller
Richarda Fleischera
(1971)



Forenzní lingvistika



- Donald Foster (1995) – případ UNABOMBER
- UNiversity and Airplane BOMBER
- Theodor Kaczynský, matematik a sériový vrah
- série bombových útoků z let 1978-1995
- charakteristika pachatele na základě manifestu zveřejněného v novinách: nadprůměrně inteligentní útočník z Chicaga, dle typických slovních spojení jej poznal jeho bratr

Stylometrie

= určení autorství (uměleckého) textu,
tj. jazyková analýza a interpretace textu,
na jejímž základě je možné určit některé
charakteristické znaky osoby, která text
vytvořila

Stylometrie

- x grafologie, x forenzní fonetika
- kvantitativní přístup, větší vzorky dat
- např. rukopisy
- Mosteller and Wallace (1963) – eseje k Americké ústavě,

Hamilton vs. Madison



Stylometrie

- statistika: např. užití obsahových a funkčních slov

TABLE 2.1. FREQUENCY DISTRIBUTION OF RATE PER THOUSAND WORDS FOR THE 48 HAMILTON AND 50 MADISON PAPERS FOR *by*, *from*, AND *to*. THE UPPER LIMIT OF A CLASS INTERVAL IS NOT INCLUDED IN THE CLASS

Rate	<i>by</i>		Rate	<i>from</i>		Rate	<i>to</i>	
	H	M		H	M		H	M
1- 3	2		1- 3	3	3	20-25		3
3- 5	7		3- 5	15	19	25-30	2	5
5- 7	12	5	5- 7	21	17	30-35	6	19
7- 9	18	7	7- 9	9	6	35-40	14	12
9-11	4	8	9-11		1	40-45	15	9
11-13	5	16	11-13		3	45-50	8	2
13-15		6	13-15		1	50-55	2	
15-17		5		—	—	55-60	1	
17-19		3	Totals	48	50	Totals	48	50
Totals	48	50						

Stylometrie

- statistika: např. užití obsahových a funkčních slov

TABLE 2.2. FREQUENCY DISTRIBUTION FOR *war*

Rate/1000	H	M
0 (exactly)	23	15
0+-2	16	13
2- 4	4	5
4- 6	2	4
6- 8	1	3
8-10	1	3
10-12	—	3
12-14	—	2
14-16	1	2
	—	—
Totals	48	50

Určované charakteristiky

- pohlaví pachatele
- stáří
- sociální a místní původ (slang, neologismy)
- vzdělání
- pravděpodobný stupeň agresivity nebo poddajnost
- vztah k autoritám
- významnost či nevýznamnost některých osob či věcí
- smyslový handicap (špatný zrak, koktavost)
- psychický handicap atd.

Metody




- větná analýza – nadprůměrně dlouhé věty atd.
- slovnědruhovú analýza – nadužívání slov atd.
- posouzení celkového stylizačního procesu
- srovnání s jinými texty vytvořenyi stejnou osobou
- obsahová analýza + postojová analýza


Metody

- větná analýza – nadprůměrně dlouhé věty atd.
- slovnědruhová analýza – nadužívání slov atd.
- posouzení celkového stylizačního procesu
- srovnání s jinými texty vytvořenými stejnou osobou
- obsahová analýza + postojová analýza

Metody

- postojová analýza = extrakce názorů a postojů z textu a řeči

facebook  Search for people, places and things   Kateřina Veselová





Miluju když jedu na melounu a odrážim se prasečima řízkama... :D
2,384 likes

✓ Liked Message * ▾

Local Business
Suggest a phone number

About – Suggest an Edit

  **2,384**

Photos Likes

Postojová analýza ve forenzní lingvistice

- výhružné dopisy / emaily
- hanobení rasy a národa / pomluvy v internetových diskusích
- terorismus
- automatická administrace internetových diskusí
- prevence sebevražd
- kyberšikana

Postojová analýza ve forenzní lingvistice

TÉMA: NÁVOD NA BEZBOLESTNOU SEBEVRAŽDU


Text: Předem Vás žádám, že nemám zájem o dotazy na důvody a rozmlouvání, důvodů mám mnoho a sílu pokračovat již nemám, zažil jsem si i na svůj nízký věk hodně a hodně věcí dokázal, bohužel jsem na totálním dně. Celý život jsem hledal smysl života a jediný smysl je pro mě konec a ať to zní jak chce už se na něj těším, život je boj... já prohrál, ikdyž jsem si to přiznal až po několika letech.



Hledám dostatečně účinný jed (či jinou látku) na 100% usmrcení, chtěl bych bezbolestnou a klidnou smrt někde v pohodlí nebo ve spánku.

Děkuji za pomoc, jinak volím variantu vlaku a tak pokud mi alespoň nepřejete konec jaký bych si přál, tak myslete na lidi co by čistili vlak.

(Potřeboval vědět i kde je látka k dostání)

13.09.2010 15:31
Uživatel: R.

Sledovat téma 

 (32 lidí)  (17 lidí)

Stránka otevřena 15248x [← Předchozí téma](#) | [Další téma →](#)

=> př. nástroj vyhledávající potenciální sebevražedné statusy na Facebooku

Postojová analýza ve forenzní lingvistice

Kyberšikana:

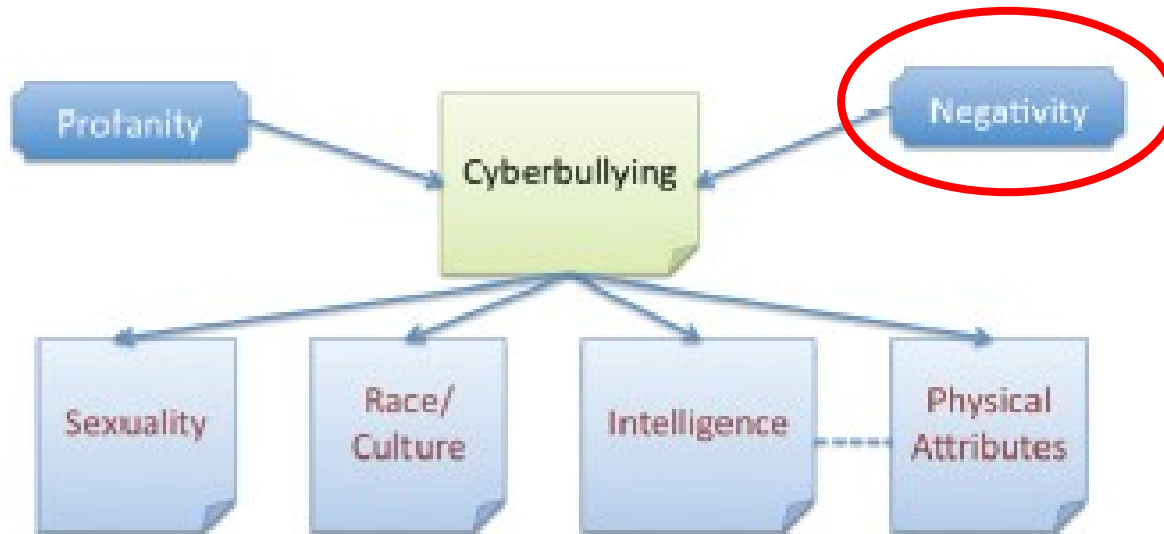
- detekce potenciálně negativních výroků na základě postojové analýzy (vulgarismy atd.)
+ klíčová slova
- Dinakar et al. 2011 – trénovací data z YouTube,
 - binární klasifikátory
 - v prototypických případech úspěšnost cca 72% (NB a SVM)

Postojová analýza ve forenzní lingvistice

Kyberšikana:

- Dinakar et al. 2011, *"Modeling the detection of Textual Cyberbullying."*
- trénovací data z YouTube
 - 4500 označovaných komentářů
 - binární klasifikátory
 - v prototypických případech úspěšnost cca 72% (NB a SVM)

Postojová analýza ve forenzní lingvistice



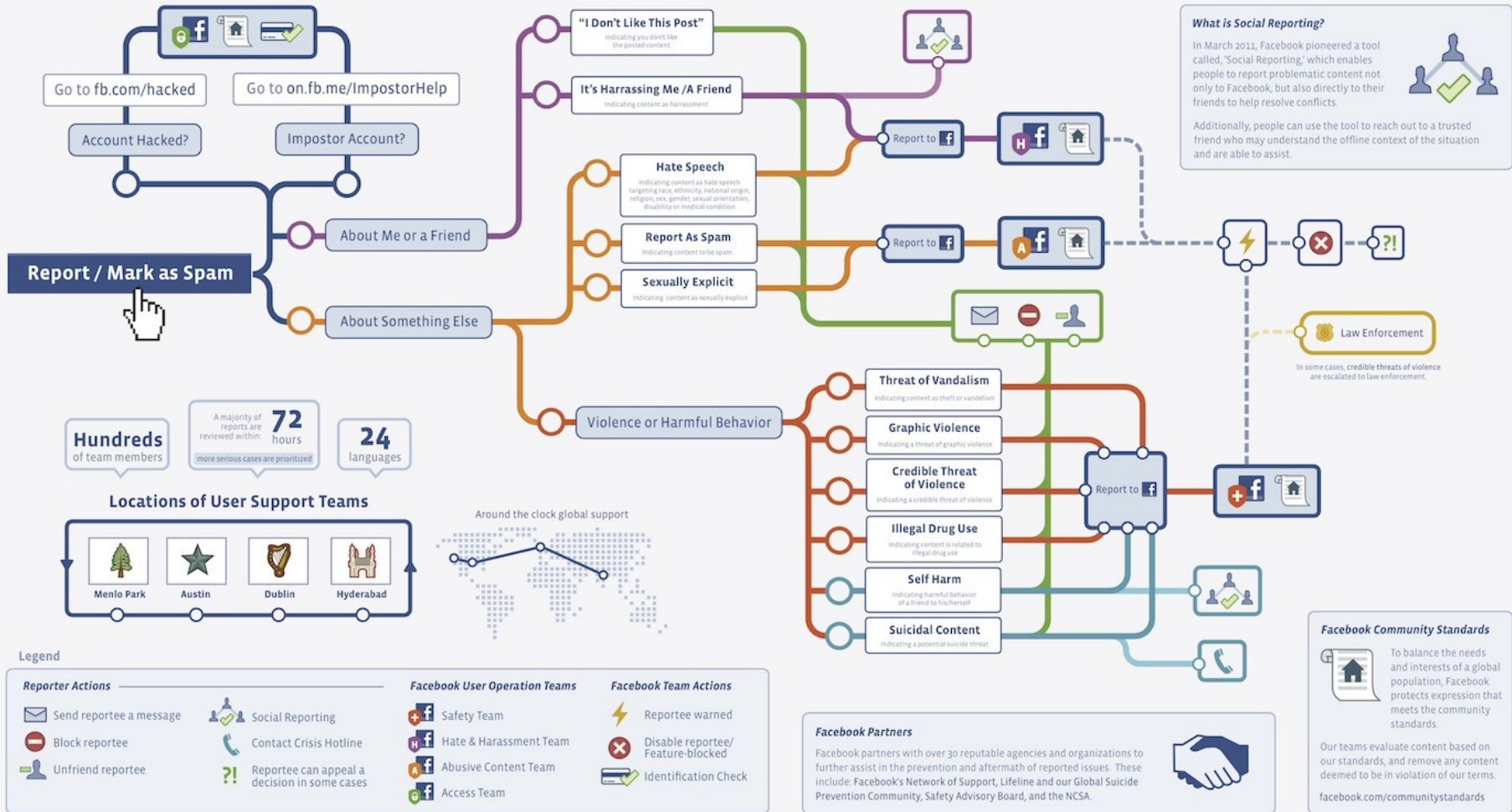
Postojová analýza ve forenzní lingvistice

- detekce **trollingu** = zasílání provokativních, hanlivých nebo irelevantních příspěvků k citlivým tématům
 - smyslem je vyprovokovat ostatní uživatele k emotivní odezvě
- detekce **hate speech** = verbální projevování nenávisti

Reporting Guide

What Happens When You Report Something?

At Facebook, nothing is more important than the safety and security of the people who use our service. With a community of over 901 million people, Facebook maintains a robust reporting infrastructure made up of dedicated teams all over the world and innovative technology systems.

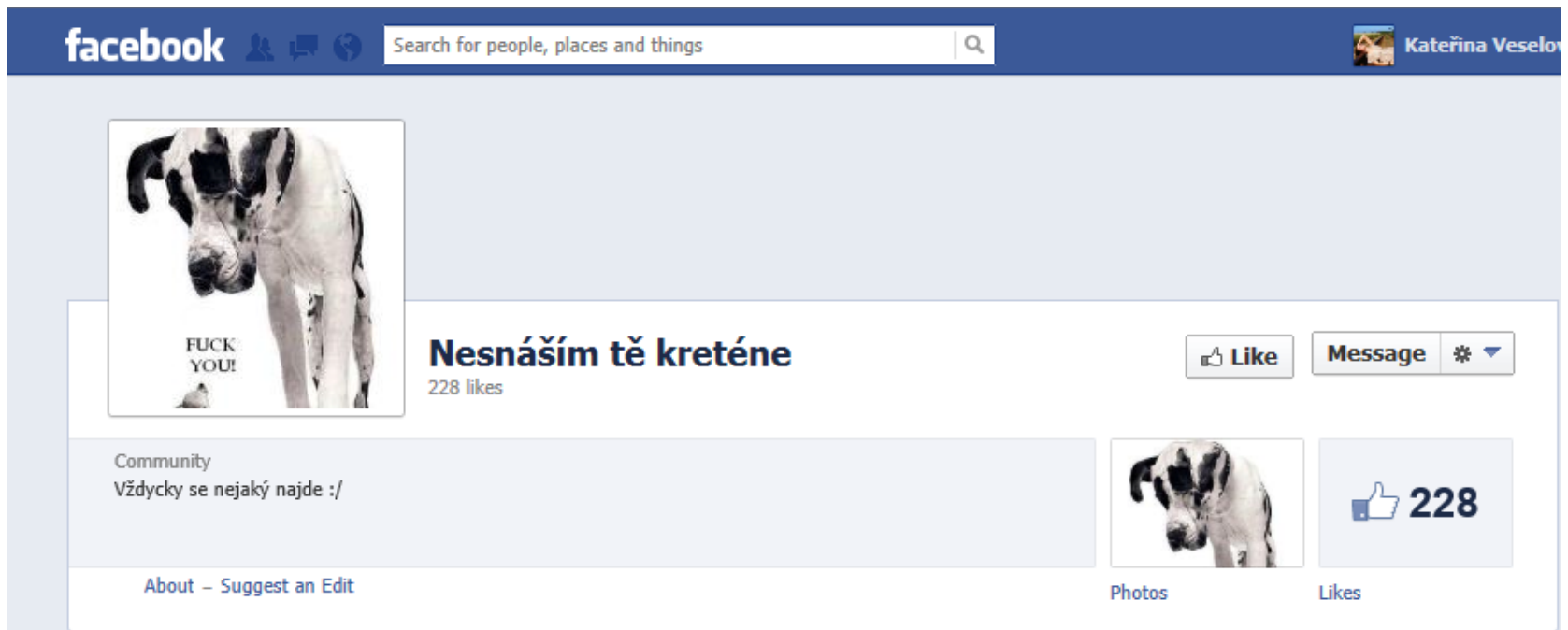


Postojová analýza ve forenzní lingvistice

= sběr velkého množství dat pro postojovou analýzu a další trénování klasifikátorů

Ztotožňování

- např. příspěvky z Facebooku vs. nesnášenlivé příspěvky v internetových diskusích



The image shows a screenshot of a Facebook post. At the top, the Facebook logo and navigation icons are visible. A search bar contains the text "Search for people, places and things". In the top right corner, the user's profile picture and name "Kateřina Veselá" are shown. The main content is a post from a community named "Nesnáším tě kreténe" (I hate you cretins), which has 228 likes. The post features a black and white photograph of a dog with the text "FUCK YOU!" overlaid. Below the post, there are buttons for "Like", "Message", and a dropdown menu. At the bottom of the post, there are links for "About" and "Suggest an Edit".

facebook

Search for people, places and things

Kateřina Veselá

FUCK YOU!

Nesnáším tě kreténe
228 likes

Like Message

Community
Vždycky se nějaký najde :/


About – Suggest an Edit

Photos Likes

Ztotožňování

- např. příspěvky z Facebooku vs. nesnášenlivé příspěvky v internetových diskusích

facebook Search for people, places and things Kateřina Veselová



Smrt důchodcům...
12 likes

Like Message

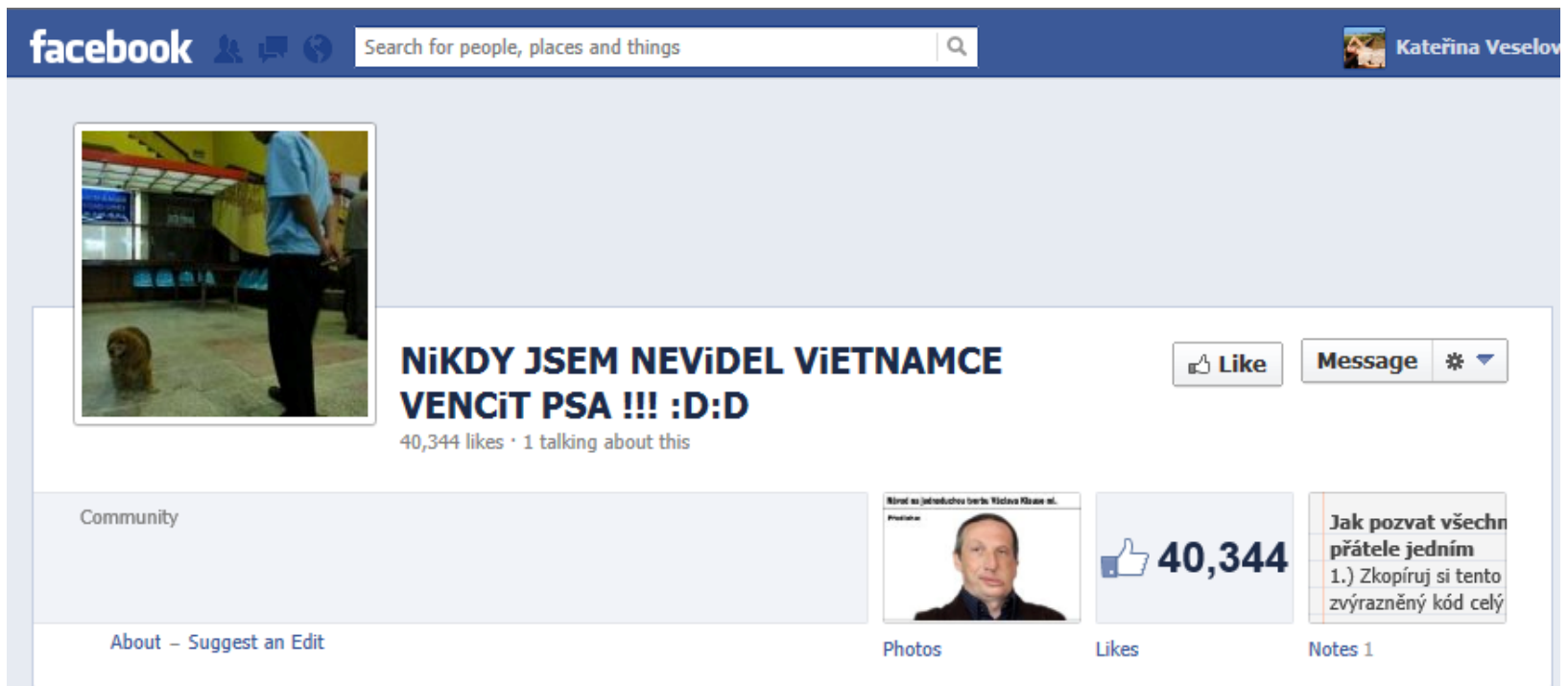
Local Business
Suggest a phone number

About – Suggest an Edit

Photos Likes


Ztotožňování

- např. příspěvky z Facebooku vs. nesnášenlivé příspěvky v internetových diskusích



The image shows a screenshot of a Facebook post. At the top, the Facebook logo and navigation icons are visible. A search bar contains the text "Search for people, places and things". In the top right corner, the name "Kateřina Veselov" is partially visible. The main post features a photograph of a man in a light blue shirt standing next to a dog. Below the photo, the text reads: "NĪKDY JSEM NEVIDEL VIETNAMCE VENCIT PSA !!! :D:D". To the right of the text are buttons for "Like", "Message", and a settings icon. Below the text, it says "40,344 likes · 1 talking about this". At the bottom of the post, there are several tabs: "Community", "Photos", "Likes", and "Notes 1". The "Likes" tab is active, showing a profile picture of a man and the number "40,344". To the right of the "Likes" tab, there is a note that says "Jak pozvat všechn přátele jedním" and "1.) Zkopíruj si tento zvýrazněný kód celý".

facebook Search for people, places and things Kateřina Veselov



NĪKDY JSEM NEVIDEL VIETNAMCE VENCIT PSA !!! :D:D

40,344 likes · 1 talking about this

Community

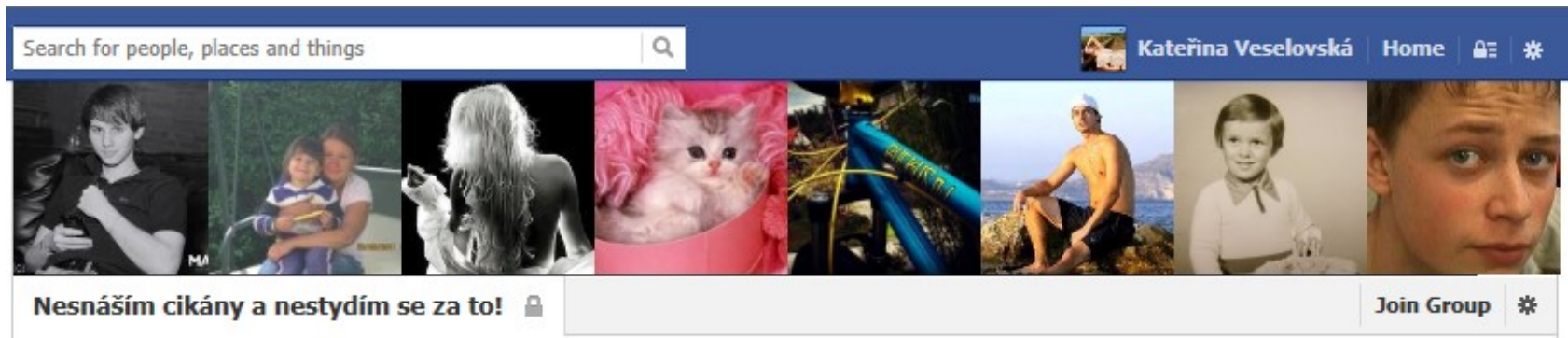
Photos Likes Notes 1

40,344

Jak pozvat všechn přátele jedním
1.) Zkopíruj si tento zvýrazněný kód celý

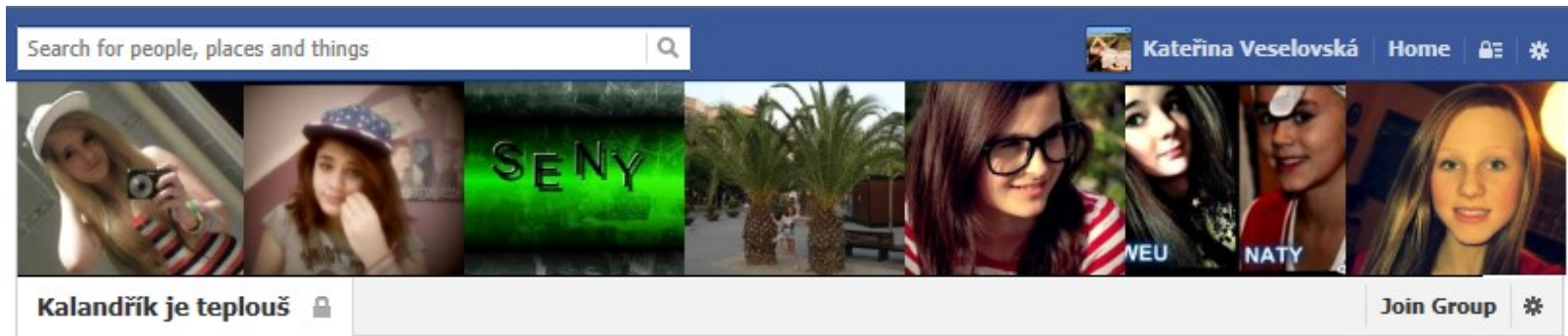
Ztotožňování

- např. příspěvky z Facebooku vs. nesnášenlivé příspěvky v internetových diskusích



Ztotožňování

- např. příspěvky z Facebooku vs. nesnášenlivé příspěvky v internetových diskusích



Postojová analýza ve forenzní lingvistice

- povaha a frekvence hodnotících výrazů (zejména těch s negativní polaritou) má zásadní vliv na odhalení nesnášenlivého textu a při existenci i krátkých referenčních pasáží (cca 500 slov) také jeho autora

= nejúčinnější pro forenzní lingvistiku:

postojová analýza + statistická analýza funkčních slov

+ objektivní aspekty (IP adresa atd.)

Detekce vulgarismů – motivace

- automatická extrakce vulgarismů v rámci postojové analýzy
- silně emocionální
- indikátory negativní polarity hodnocení v textu
- = ztotožňování hanlivých příspěvků

Vulgarismy – motivace



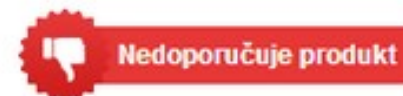
Jachyms

Přidáno: 26. listopadu 2013

10% 



+ uplný hovno

- všechno dopíči



no to si ten zasranej dr.dre ze mě už dělá prdel, větší sračku fakt udělat nemohl dopíči, mu tam ty jeho zkurvený sluchátka co hrajou jak kdybych si koupil brčko někde u čonga a strčil si to do ucha pošlu zpátky do afriky!

Utlumení vnějších ruchů:	Nedostatečné
Středky a výšky:	Zkreslené
Basy:	Dostatečné
Jsem hudební:	Profesionál
Délka kabelu:	Krátké

Je tato recenze užitečná?  [Ano](#) (14)  [Ne](#) (5)

Vulgarismy – motivace

- rozvoj webu 2.0 = rozvoj vulgarismů
- omezování pravidel u diskusí pod články
(elektronická registrace, papírová registrace,
nahrazení diskuse facebookovými thready)
- blokace diskutujících

Vulgarismy – motivace

→ nové strategie

→ současné nástroje pro lemmatizaci selhávají

Barman je @#%\$!! > barman|být|???????

Vulgarismy – teoretická východiska

- klení umožňuje expresivně vyjádřit emoční zaujetí
- expresivity je dosahováno porušováním společenských tabu

Vulgarismy – teoretická východiska

Tabu střední a východní Evropy:

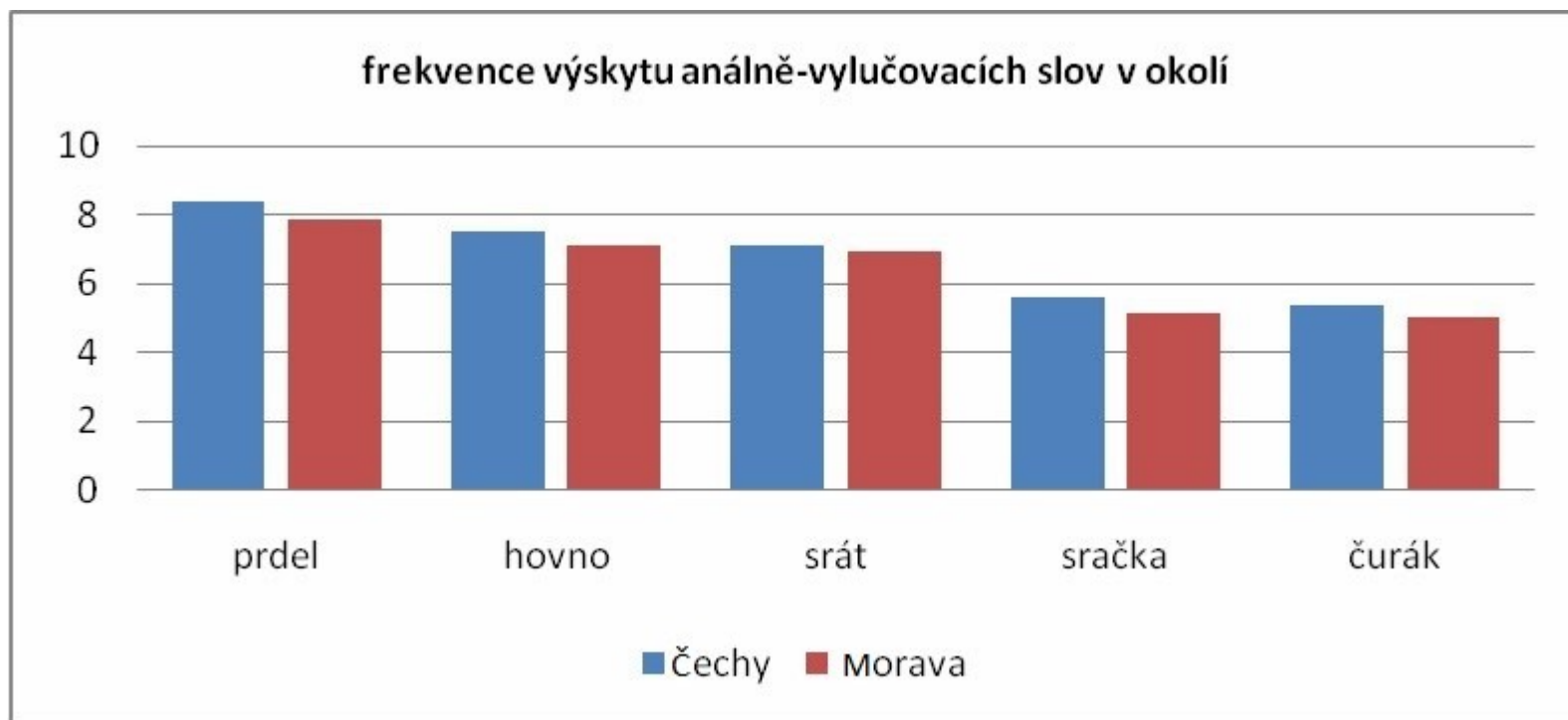
- a) sexuální chování
- b) vylučování a vylučovací orgány
- c) eschatologie (posmrtný život)
- d) náboženské hodnoty*

* *krucifix, hergot* atp. – nahrazeno výrazy z análně-vylučovací oblasti

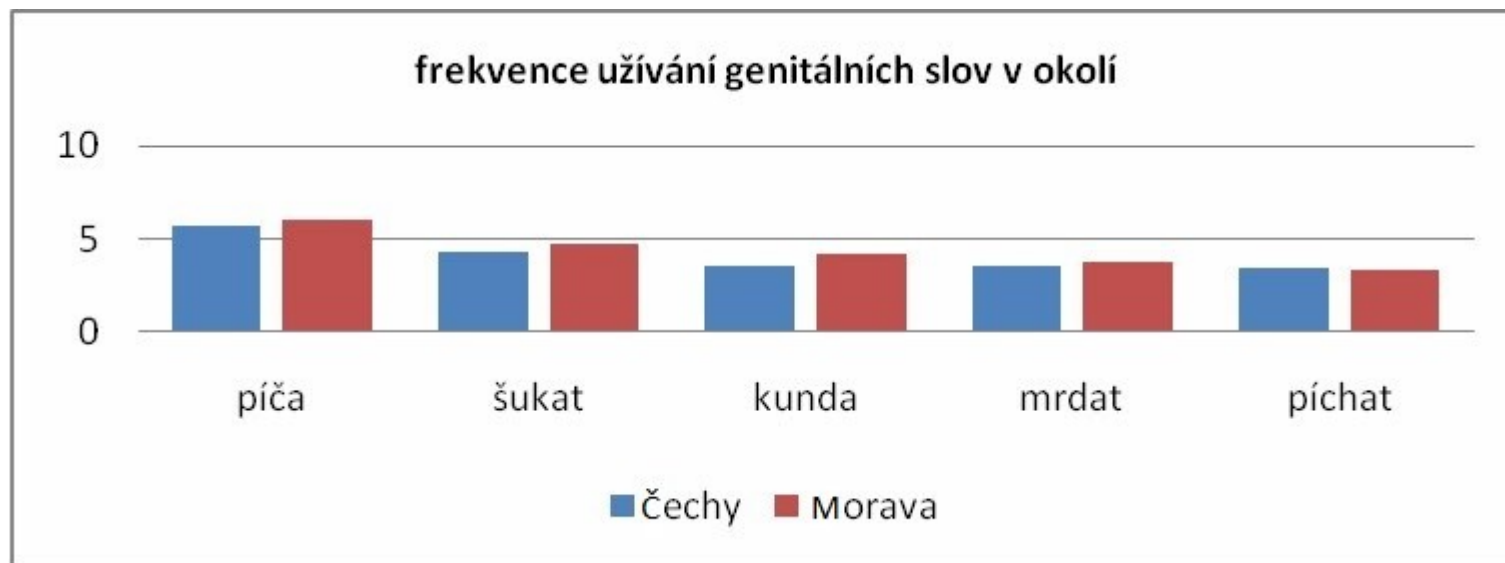
Vulgarismy – teoretická východiska

- preference typů vulgarismů dána i demograficky
- Morava genitální, Čechy anální – Franče 2009 (dotazníkové šetření)
- lokace lze vysledovat v nové podobě diskusí

Vulgarismy – teoretická východiska



Vulgarismy – teoretická východiska



Důležitost vulgarismů pro postojovou analýzu

- sémantika = co?
- formální prostředky = jak?

Vulgarismy – sémantika

- desambiguace – kolokační analýza

Dostaneš hovno. = Nedostaneš nic. (neutrální)

Stojí to za hovno. = Za nic to nestojí. (negativní)

Ale hovno. = Mýlíš se. (negace, nesouhlas)

Mele hovna. = Říká nesmysly. (nesouhlas)

Vulgarismy – sémantika

- desambiguace – kolokační analýza

Byl na sračky. = Byl opilý. (neutrální)

Ten film je sračka. = Ten film je špatný. (negativní)

Vulgarismy – sémantika

- desambiguace – kolokační analýza

Je tu tma jak v prdeli. (neutrální)

S Jardou je vždycky prdel. (pozitivní?)

A je to v prdeli. (negativní)

Vulgarismy – sémantika

- desambiguace – kolokační analýza

Svědila ji kunda. (neutrální)

Na baru je pěkná kunda. (pozitivní)

Topolánek je pěkná kunda. (negativní)

Vulgarismy – formální prostředky

- ideálně „čistý“ text



Jakub Danel · ★ Autor nejlepších komentářů · Frezař ve společnosti Kovona System

A co ja do piči ona nežila v takove díře jako ČR a jde se všeset to je piča :D Si myslí že je jedina co si nezašukala se musí jít zabít :DD Ja jebu na všechno aji na smrt :DD

Odpovědět · To se mi líbí · 6 leden v 8:34

Vulgarismy – formální prostředky

- nepísmenné znaky



Tomáš Vajčner · ★ Autor nejlepších komentářů · MěVG Klobouky u Brna
to je ale pí*a !:D

Odpovědět · To se mi líbí · 13. listopad 2013 v 3:50

Vulgarismy – formální prostředky

- nepísmenné znaky

Roman Klvač, Hustopeče nad Bečvou

Neděle, 11. května 2014, 16:24:21 | [Souhlasím](#) | [Nesouhlasím](#) | +5

pokud nejsi buzerant nebo politik,jseš ho.no....

Vulgarismy – formální prostředky

- záměna znaků

aspirinek

1.8.2013 19:24

Kundus, nemusíš tady všem předvádět, jak se doma titulujete, ty howado s modrým minimozečkem kura domácího.

+1 



Vulgarismy – formální prostředky

- jazykové hříčky

Proč ta zrádná peacha s hadíma očima konečně nezaleze do kanálu?!, Neptunes3, 30. 1. 2014 19:18:49

Nikdo na její hlody není zvědavěj, fakt že ne. Hroší kůže. Neskutečný

Reagovat



Mira UL

Malý detail:

Ad (9) — Každousek argumentuje nespravedlností daňového systému, jako by to byla výhoda. Kapitalisté neplatí daň z výrobního majetku. Daň ze zisku korporací je malá.

Vulgarismy – formální prostředky

- všeobecná znalost

Jindřich Nesrsta, Bystřice pod Hostýnem

Neděle, 11. května 2014, 12:41:00 | [Souhlasím](#) | [Nesouhlasím](#) | +27

Svět se v obrací

Vulgarismy – formální prostředky

- zkratky

 (0) 12.01.14 13:57 **Marek (164622)** **Odpověď**

Co to kua má znamenat?!

Co to kua má znamenat?!

Dosavadní doporučení: přečíst (0) Vaše doporučení: určitě přečíst přečíst nečíst

- wtf, omg, rtfm, dmnce, twl...

Vulgarismy – formální prostředky

- zkratky



Vulgarismy – možná řešení

- standartní spell-checking – porovnávání řetězců se slovníkem
- pravidlový systém – netypické kombinace znaků (bigramy, trigramy)
- Levenshteinova distance

Vulgarismy – možná řešení

- stop words list (kua, wtf, rofl...)
- devyatlátor

Vypatlátor/devypatlátor

Jednoduchá hračka pro překlad z debilštiny, kterou poslední dobou používá poměrně značný počet dětských sebevrahů.

Vstupní pole - sem napiš text který chceš přeložit, poté pouze klikni na tlačítko "DEVYPATLAT!"

Barman je howado.

↓ ↓ ↓ ↓ ↓ ↓ ↓ ↓
DEVYPATLAT
↓ ↓ ↓ ↓ ↓ ↓ ↓ ↓

Výstupní pole - zde si můžeš po devypatlání textu přečíst výsledek tohoto skromného devypatlátoru.

BARMAN JE HOVADO.

Správná detekce vulgarismů nemusí stačit...



Evelina Eva - ★ Autor nejlepších komentářů - Čalovo

Šárinka MerSí Machová ti osobo co to tady píšeš si asi velká kráva
pubertální ti sere ti u prdele začínaš semnou tak si dej bacha na mě
holka jedna nevychovaná co si osobě myslíš ?? ja nejsem žádná
pubertačka za prve a podruhy sem starší osoba tak se tak chovej osobo
jedna ..ti tvoje sexi hrátky ze zamylovanosti nech jo sem muslimka tak
sed pohov krávo ,,

Odpověď - To se mi líbí - 13. listopad 2013 v 2:54

Forenzní psycholinguvistika

„Psycholinguvistické profilování případu série výhrůžných anonymních dopisů“ (Veselá a Musilová, 2012)

- demonstrace a psychologická interpretace zrodu a vývoje přípravy úmyslného násilného aktu
- forenzní psychoanalýza

Děkuji za pozornost.