

A Szöveglabor bemutatása

Simon Eszter

2020. szeptember 23.

PIM DBK

1. Bevezetés

2. Tervezési munkafázisok

Bevezetés

- forráskiadások
- kritikai kiadások
- born digital
- szöveglabor

a munkacsoport:

Kalcsó Gyula

Simon Eszter

Mihály Eszter

szöveglabor beszámoló:

október 7. 10:00

Tervezési munkafázisok

1. Hazai és nemzetközi jó gyakorlatok figyelembe vétele a tervezéskor
2. Szöveglabor funkciók
3. Eszközök a nyelvi feldolgozáshoz
4. A nyelvi feldolgozásra épülő eszközök a szöveglabor fejlesztéséhez
5. Feldolgozható szövegtörzsek körének meghatározása
6. DHUplába integrálható szövegtörzsek típusainak definiálása
7. Koherencia létrehozása a többi alterülettel
8. Jogi kérdések

a feldolgozható szövegek körének 3 köre:

1. Saját szövegek a dHUpa 1-3. területéről (forráskiadás, born digital, kritikai kiadás)
2. távoli szövegek (DIA, MEK stb., elsősorban konzorciumi partnerek korpuszai)
3. egyéb külső hozott szövegek a szöveglaborba

ami ehhez kell:

- minőségbiztosítás
- formai követelmények
- feltöltőfelület → ellenőrzés → elemzés, ha jó a szöveg & visszajelzés, ha nem jó

- **formátum:** kép: nem, szöveg: igen
- **nyelv:** csak magyar → a nem magyar nyelvű részeket kiszűrni és megjelölni
- **kor:** csak modern → a régi szövegekből legyen normalizált; ha nincs, akkor csak azt lehet vele kezdeni, amit az elemzetlen korpuszon lehet csinálni

- az integrálásra szánt korpusz tartalmaz más annotációkat → a formátumot igazítani kell a dHUpla elvárt bemeneti formátumához → a korpusz tulajdonosa megoldja az átalakítást, vagy lemond az annotációiról
- az integrálásra szánt korpusz már elemezve van, de más elemzővel → újra kell elemezni, különben nem lesz kompatibilis a nyelvi elemzésre épülő további elemzésekkel

- a többi alterületen is szükség lesz nyelvi feldolgozásra → ezeknél a szövegeknél meg kell határozni, hogy mely ponton mennek át a szöveglabor elemzési lépésein → minél hamarabb, annál jobb
- interoperabilitás: oda-vissza konverzió lehetséges legyen, egységes be-kimeneti formátum
- a DHupla korpuszainak összekötése a névtérrel

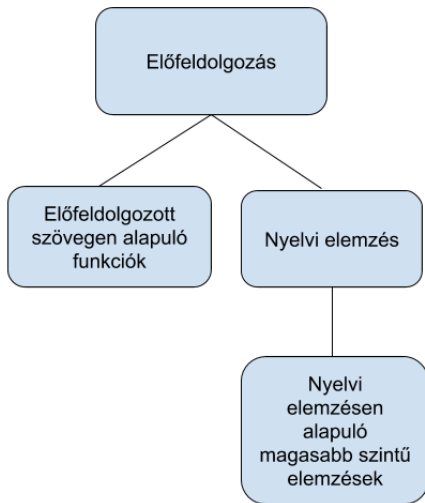
milyen a jogi státusza annak a korpusznak, amit feltöltenek a dHUplába és elemeznek a szöveglaborral?

- elvárás: legyen hozzáférhetővé téve a dHUplán a szöveg és a szöveglabor által kiadott elemzések és vizualizációk
- kérdések:
 - a dHUplán közzétett korpusz máshol is publikálható?
 - csak egy adott korpuszra kötődik az együttműködés?

- szakirodalom: kivonatok stilometriai és MSZNY cikkekből
- eszközök
 - nyelvi feldolgozó eszközök
 - nyelvi feldolgozásra épülő eszközök
- best practices: projektek, workflow-k

a szövegfeldolgozásnak 4 szintje van:

1. előfeldolgozás
2. előfeldolgozott szövegen alapuló funkciók
3. nyelvi elemzés
4. nyelvi elemzésen alapuló magasabb szintű elemzések



- bemeneti formátum ellenőrzése:
 - UTF-8 karakterkódolás
 - txt
- tokenizálás
- mondatra bontás

- w*ldcardos keresés
- [rR]eguláris kifejezések+el való keresés
- kollokációk keresése (akár wildcardosan)
- konkordancia készítése
- szógyakoriság
- szókészlet listázása
- szó n-gramok → gyakorisági lista
- központosítás vizsgálata
- magánhangzók–mássalhangzók arányának vizsgálata
- szavak és mondatok hosszának eloszlása
- az egyes szavak hány karakterből állnak
- az egyes mondatok hány szóból állnak
- verstani elemzés: metrika, rím, alliteráció automatikus detektálása

<http://clara.nytud.hu/mnsz2-dev/index.html>

- reguláris kifejezésekkel való keresés
- kollokációk keresése
- konkordancia készítése
- szógyakoriság
- szókészlet listázása
- CV-váz

.apa			
File	View	Properties	Analyse Queries Help
testével pusztai durva mén dübörög,	12	- - - - U U - U - U U U	c
és minden sors hasadt dísznő-csülökre írva.	13	- - - - U - - - U - U - U	d
PROTEUS			
A parton Proteus alakoskodik:	10	U - - - - U U - U U	a
most majdnem isten, most a lehetetlen,	11	- - U - - - U U U - U	b
most számtalan hús gyönggyé szerteröppen,	11	- - U - - - - U - U	b
most sziklává merül, most újra híg.	10	- - - - U - - - U -	a
Víz és föld határ-láncára bukik,	10	- - - U - - - U U U	a
de menny s mélység közt lakik: ő a tenger,	11	U - - - - U U - U - U	b
keblén sarat ringat, s e szerelemben	11	- - U - - - U U U - U	b
sár urává, emberré változik.	10	- U - - - - - U U	a
Bírák és bankosok areopágja	11	- - - - U U U U U - U	c
megméri, sorsa tűhegyen forog:	10	- - - U - U - U - U	d
vagy visszadobják habzó szabad árba	11	- - U - - - - U U - U	c
vagy uszonyán ember-gúzs csikorog:	10	U U U - - - - U U U	d
de minden sejtje tengert párolog:	10	U - - - - U - - - U U	d
száraz porvert tanyán nem élt hiába.	11	- - - - U - U - U - U	c
HEPHAISTOS			
Lendült, szökött az első mozdulat:	10	- - U - U - - - U U	a

Lesi Zoltán: Automatikus formai verselemzés.

In: Alkalmazott Nyelvtudomány VIII. évf. 1-2. szám 2008. 197-208.o.

- morfológiai elemzés és egyértelműsítés
- tulajdonnév-felismerés
- NP chunking
- szintaktikai elemzés
 - összetevős elemzés
 - függőségi elemzés

morfológiai információn alapuló funkciók:

a morfológiai elemző és egyértelműsítő kimenetén alapulnak

- szófaji gyakoriság
- szófajok aránya
- toldalékgyakoriság
- szógyakoriság lemmatizált szövegen
- keresések:
 - regexes keresés → bizonyos szófajú/tövé/ragozású szavak keresése
 - kollokációk keresése → bizonyos szófajú/tövé/ragozású szavak együttállása
 - konkordancia készítése → lemmatizált szövegen
 - szókészlet listázása → lemmatizált szövegen

4. Táblázat: A domének hasonlósága a szófajok eloszlása terén.

	Lenin- mesék	Lenin előtt	Lenin után	szovjet mesék	Mátyás- mesék	Rákosi- mesék	szocreál irodalom
Lenin- mesék		0,993	0,991	0,9976	0,9939	0,9951	0,9865
Lenin előtt	0,993		0,9953	0,9942	0,993	0,9917	0,9868
Lenin után	0,991	0,9953		0,9944	0,995	0,9919	0,9864
szovjet mesék	0,9976	0,9942	0,9944		0,9939	0,9939	0,9853
Mátyás- mesék	0,9939	0,993	0,995	0,9939		0,9927	0,9853
Rákosi- mesék	0,9951	0,9917	0,9919	0,9939	0,9927		0,9914
szocreál irodalom	0,9865	0,9868	0,9864	0,9853	0,9853	0,9914	

Horváth Csilla: “S azóta jól élnek, és vidám dalokat énekelnek Leninről és Sztálinról.” Szovjet propagandamesék műfaji azonosításának kísérlete.

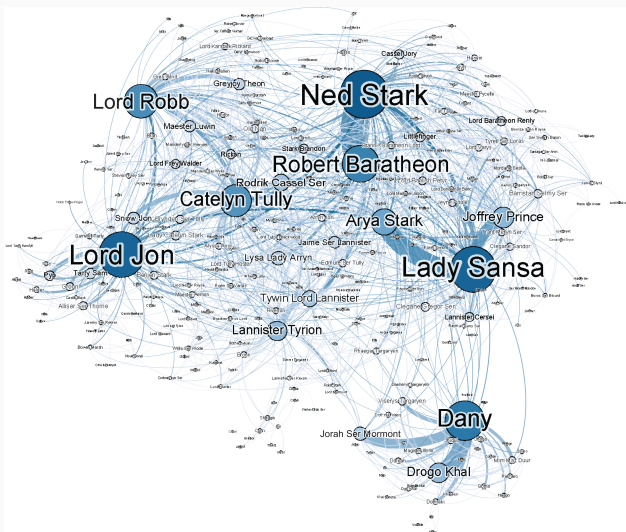
In: MSZNY 2018. 331-339.o.

szintaktikai információn alapuló funkciók:

mondatszerkezeti statisztika

tulajdonnév-felismerésen alapuló funkciók:

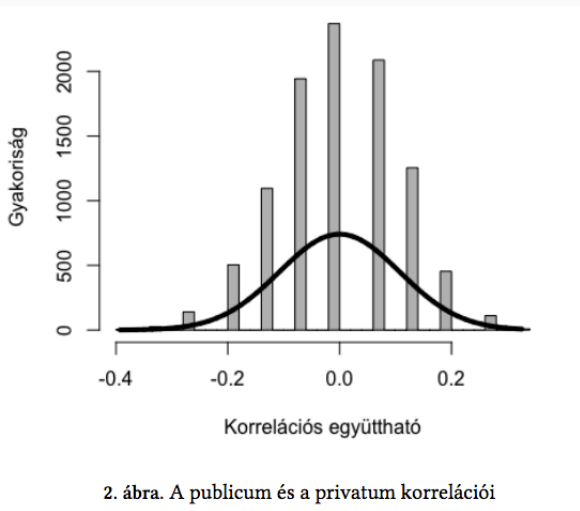
- tulajdonnév-statisztika
- egy dokumentumban előforduló személyek kapcsolati hálója



DOI: 10.7717/peerj-cs.189/fig-1

emóciógyakoriság

- nyolc kategóriás emóciószótáron alapul: öröm, bánat, düh, szeretet, félelem, undor, meglepetés, feszültség → a szótár az ezeket az emóciókat megjelenítő szavakat tartalmazza
- kell hozzá lemmatizálás
- bemenetként az egyes tokenek lemmáját kéri, és azt összehasonlítja a szótár elemeivel: szimpla listaelem-illesztés
- az egyes emóciókhoz tartozó szövegbeli előfordulásokat leszámoljuk, belőlük gyakorisági listát készítünk



Labádi Gergely: Az olvasó gép. In: Digitális Bölcsészet 1 (2018)

eszközök:

- *e-magyar*
- Magyarlánc 3.0
- UDPipe
- huspaCy
- stanza

kérdések:

- gyorsaság → offline vagy realtime elemzés?
- továbbfejleszthetőség → modularitás, nyílt forráskód
- fenntarthatóság
- kívánt funkciók
- pontosság

ezt a kívánt funkcionalitások listája fogja meghatározni

kérdések:

- meglevő eszközök adaptálása vagy új eszközök fejlesztése?
- interoperabilitás: átjárhatóság az egyes eszközök között → egységes be-kimeneti formátum, konverzió
- integrálhatóság: ha nem nyílt forráskódú, akkor hogyan integráljuk a láncba?

stilometria → Kalcsó Gyula előadása

Köszönöm a figyelmet!

`simon.eszterke@gmail.com`