

Bevezetés az e-magyar programcsomag használatába

Vadász Noémi

2019. május 2.

MTA Nyelvtudományi Intézet
vadasz.noemi@nytud.mta.hu

1. szövegelemzés számítógéppel
 - elemzési lépések
 - elemzőláncok
 - az **e-magyar** működése
2. az **e-magyar** használata

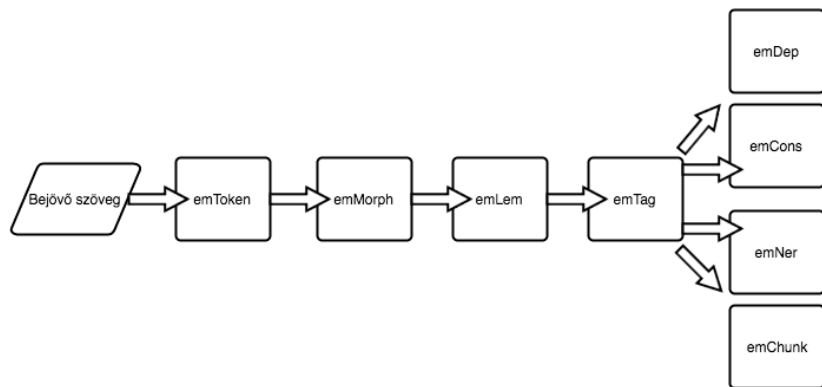
Szövegfeldolgozó lépések

1. szöveg felbontása kezelhető egységekre: **mondatszegmentálás** és **tokenizálás**
2. a szavak elemzése: **morfológiai elemzés**
3. az elemzések **egyértelműsítése** a mondatban
4. **főnévi csoportok, tulajdonnevek** meghatározása
5. **függőségi** és **összetevős** mondatelemzés

A lépéseket **láncba** köthetjük, ahol az alsóbb szintű elemzési lépés **kimenete** a felsőbb szintű elemzés **bemenetét** képezi.

Elemzőláncok magyarra: magyarlanc és **e-magyar**

Az e-magyar elemzőlánc felépítése



Két feladat:

1. mondatszegmentálás

- rövidítések, római számok, sorszámok (pont van a végén)
- mondatok idézetben vagy zárójelben

2. tokenizálás

- detokenizálhatóság
- páros írásjelek (zárójel, idézőjel, aposztróf) kezelése
- kérdőpartikula kezelése

Miért kell tokenizálni?

- a szintaktikai elemző mondatokon dolgozik
- a morfológiai elemző tokeneken dolgozik

- szabályalapú
- a leggyakoribb rövidítéseket ismeri (kb. 450 rövidítés)
- a kimenetben `xml`-tagek jelölik a szó- és mondathatárokat

A kutya váratlanul ugatni kezdett. Ettől úgy megijedt dr. Thorotzkay Alfréd, hogy hanyatt esett az aszfalton. Felesége, aki egyébként a BKV Zrt.-nél dolgozik, egyből rohant hozzá, amint ezt megtudta.

<code><s> ... </s></code>	mondat
<code><w> ... </w></code>	szó
<code><c> ... </c></code>	írásjel
<code><ws> ... </ws></code>	whitespace

<S><W>A</W><WS> </WS><W>kutya</W><WS></WS><W>váratlanul</W><WS>
</WS><W>ugatni</W><WS></WS><W>kezdett</W><C>.</C></S><WS> </WS>
<S><W>Ettől</W><WS> </WS><W>ügy</W><WS> </WS><W>megijedt</W><WS>
</WS><W>dr.</W><WS> </WS><W>Thorotzkay</W><WS>
</WS><W>Alfréd</W><C>,</C><WS> </WS><W>hogy</W><WS>
</WS><W>hanyatt</W><WS> </WS><W>esett</W><WS> </WS><W>az</W><WS>
</WS><W>aszfalt</W><C>.</C></S><WS>
</WS><S><W>Felesége</W><C>,</C><WS> </WS><W>aki</W><WS>
</WS><W>egyébként</W><WS> </WS><W>a</W><WS>
</WS><W>BKV</W><WS></WS><W>Zrt.-nél</W><WS>
</WS><W>dolgozik</W><C>,</C><WS></WS><W>egyből</W><WS>
</WS><W>rohant</W><WS> </WS><W>hozzá</W><C>,</C><WS>
</WS><W>amint</W><WS> </WS><W>ezt</W><WS>
</WS><W>megtudta</W><C>.</C></S>

A
kutya
váratlanul
ugatni
kezdett
.

Ettől
úgy
megijedt
dr.
Thorotzkay
Alfréd
,
hogy
hanyatt
esett

- `tsv` fájl
- minden sor egy token
- mondatok között üres sor

- a tokenhez hozzárendeljük **az összes lehetséges elemzést**
- különböző címkekészletek léteznek
- a címkékben lehetnek:
 - morfoszintaktikai információk
 - a jelentésre vonatkozó információk
 - a hangalakra vonatkozó információk (allomorfia)
 - szófajkód
 - lemma
 - morfológiai szegmentumok

- tőtár, toldaléktár és nyelvtan
- szabályalapú morfológiai elemző (véges állapotú transzdúcer)
- szabályalapú tövesítés a tőalkotó morfémák figyelembevételével

szemét

1. szem[/N]=szem+e[Poss.3Sg]=é+t[Acc]=t
2. szem[/N]=szem+é[AnP]=é+t[Acc]=t
3. szemét[/N]=szemét+[Nom]=

Két morfológiai címkekészlet

emMorph

- nem tükrözi a morfológiai jelöltséget
- derivációt, szóösszetételt is kódol
- tartalmazza a lemmát, a morfológiai szegmentumokat, allomorfort
- `adtad ad[/V]=ad+tad[Pst.Def.2Sg]=tad`

UD (Universal Dependencies)

- külön szófajkód
- linearizált jegy-érték struktúra
- nem tükrözi a morfológiai jelöltséget
- képzést, derivációt nem kódol, csak inflexiót
- nem tartalmazza a lemmát, a morfológiai szegmentumokat
- `adtad ad VERB`
`Definite=Def | Mood=Ind | Number=Sing | Person=2 |`
`Tense=Past | VerbForm=Fin | Voice=Act`

Morfológiai egyértelműsítés

- minden tokenhez megvan az összes lehetséges elemzése
- el kell dönteni, hogy **az aktuális mondatban** melyik elemzést válasszuk
- tehát az egyértelműsítő már mondatokat néz
- a lehetséges elemzések legvalószínűbb láncát keresi

Péter Marira vár.

Péter a vár előtt áll.

Statisztikai alapú egyértelműsítés:

- sok, már elemzett és emberek által egyértelműsített mondatot megnézett a program
- ebből megtanulja, hogy mik a gyakori mintázatok a szövegben
- nem egyszerűen azt nézi, hogy a vár gyakrabban ige, mint főnév
- hanem azt, hogy egy névelőt sokkal gyakrabban követ főnév, mint ige
- és azt, hogy a -RA esetragú főnév után valószínűbb az ige, mint a főnév

A kastély nem vár.

A	a	[/Det art.Def]
kastély	kastély	[/N][Nom]
nem	nem	[/Adv]
vár	vár	[/N][Nom]
.	.	[/PUNCT]

A kastély nem vár senkire.

A	a	[/Det art.Def]
kastély	kastély	[/N][Nom]
nem	nem	[/Adv]
vár	vár	[/V][Prs.NDef.3Sg]
senkire	senki	[/N Pro][Subl]
.	.	[/PUNCT]

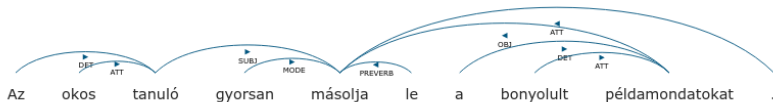
- bemenete a mondat egyértelműsített tokenekkel
- kimenete a függőségi fa
- függőség: a szavak közötti reláció
- a fa minden csomópontja egy szó
- a csomópontok közötti él a gyerek (függő) csomóponttól mutat a szülő felé
- az élcímkék a függőségi viszony típusai

Az **e-magyarba** épített függőségi mondatelemző statisztikai megközelítésű.

Az exkatonát kórházba szállították, ahol két műtétet is végrehajtottak rajta.

1	Az	2	DET
2	exkatonát	4	OBJ
3	kórházba	4	OBL
4	szállították	0	ROOT
5	,	4	PUNCT
6	ahol	10	LOCY
7	két	8	ATT
8	műtétet	10	OBJ
9	is	8	CONJ
10	végrehajtottak	4	ATT
11	rajta	10	OBL
12	.	0	PUNCT

Az okos tanuló gyorsan másolja le a bonyolult példamondatokat.



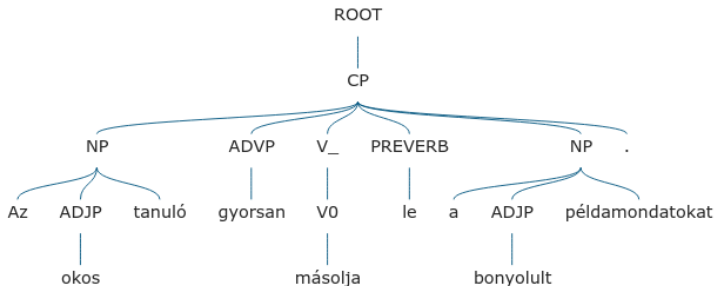
- bemenete a mondat egyértelműsített tokenekkel
- kimenete az összetevős fa
- a mondatot frázisokra bontja
- a fa levelei a mondat szavai (terminálisok)
- a fa nem-levél csomópontjai a frázisok (nem-terminálisok)
- az élek címkézetlenek

Az **e-magyarba** épített összetevős mondatelemző statisztikai megközelítésű.

Az exkatonát kórházba szállították, ahol két műtétet is végrehajtottak rajta.

1	Az	(ROOT(CP(NP*
2	exkatonát	*)
3	kórházba	(NP*)
4	szállították	(V_(V0**))
5	,	*
6	ahol	(ADVP*)
7	két	(NP*
8	műtétet	*)
9	is	(C0*)
10	végrehajtottak	(V_(V0**))
11	rajta	(NP*)
12	.	**))

Az okos tanuló gyorsan másolja le a bonyolult példamondatokat.



A főnévi csoportok és a tulajdonnevek felismerése

- ún. szekvenciális címkézési feladatok
- főnévi csoport keresése: **maximális NP** (olyan főnévi csoport, amely nem része főnévi csoportnak)
- **tulajdonnevek** keresése és névkategóriákba sorolása (személynév, intézménynév, földrajzi név, egyéb)
- statisztikai gépi tanuláson alapul, a HunTag3 program végzi
- sok szöveget látott, amelyben be vannak jelölve az NP-k/tulajdonnevek
- és hogy milyen tulajdonságokra figyeljen a tanulóskor (nagybetűvel kezdődött, mondat elején van, szófaj)
- összefüggéseket talál a tulajdonságok és a címkék között, amelyből valószínűségi modellt épít

A szállásunk egy Balaton melletti kis üdülőfaluban, Zamárdiban volt.

A	B-NP
szállásunk	E-NP
egy	B-NP
Balaton	I-NP
melletti	I-NP
kis	I-NP
üdülőfaluban	I-NP
,	I-NP
Zamárdiban	E-NP
volt	O
.	O

[...] közölte Wolf László, az OTP Bank vezérigazgató-helyettese az MTI érdeklődésére.

közölte	0
Wolf	B-PER
László	E-PER
,	0
az	0
OTP	B-ORG
Bank	E-ORG
vezérigazgató-helyettese	0
az	0
MTI	1-ORG
érdeklődésére	0
.	0

A nyelvi elemző korlátai

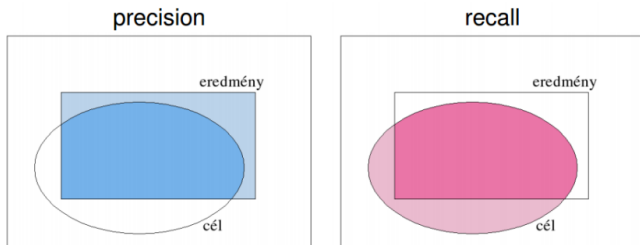
gold standard: kézzel címkézett anyag, amihez hasonlítunk

	cél pozitív	cél negatív
eredmény pozitív	True positive (TP)	False positive (FP)
eredmény negatív	False negative (FN)	True negative (TN)

Egy példa: egy kép gépi feldolgozása

- egy gép eldönti 1 000 képről, hogy van-e rajta kismacska
- az 1 000 képet ember is kiértékeli (ez lesz a gold standard)
- összevetjük az ítéleteket
- összeszámoljuk a TP, TN, FN találatokat
- a számok segítségével értékeljük a gép teljesítményét

A nyelvi elemző korlátai



- **pontosság** (precision): $TP/(TP+FP)$
A macskásnak ítélt képek közül hány tartalmaz ember szerint is kismacskát?
- **fedés** (recall): $TP/(TP+FN)$
Megtalálta-e a gép az összes macskás felvételt?
- **F-mérték**: a pontosság és a fedés harmonikus közepe

1. GATE

- nyelvfeldolgozó keretrendszer
- grafikus felület
- platformfüggetlen
- egységes annotációs modell
- más elemzőeszközök is beilleszthetők a láncba
- jól dokumentált
- soklépéses, viszonylag bonyolult telepítés

2. emTSV

- terminálból működtethető
- egységes adatformátum a modularitásért
- az elemzőlánc lépései között ki- és beléphetünk
- a kimenet egy könnyen feldolgozható fejléces **tsv**
- kétféle morfológiai címkékészlet: emMorph és UD

3. honlap

- a modulok és a lánc kipróbálására
- max. 6 000 karakter hosszú szövegek elemzésére
- kétféle formátumban letölthető kimenet: **xml** és **tsv**

File Options Tools Help

GATE

- Applications
 - Pipeline 0003
 - Language Resources
 - maistantulum_1.txt_00012
 - Processing Resources
 - HU 4. 'emEnter' Named Entity Reco
 - HU 4. 'emDep' Dependency Parser
 - HU 4. 'emCons' Consistency Parser
 - HU 4. 'emChunk' NP Chunker
 - HU 3. 'emTag' POS Tagger and Lemmatizer
 - HU 2. 'emMorph-em Lem' Morphological Analyzer
 - HU 1. 'emToken' Sentence Splitter
 - Datastores

Messages Pipeline 00013 maistantulum_1...

Annotation Sets Annotations List Annotations Stack Co-reference Editor Text

Nyaralásra néha repülővel utazó európai turistákért hozzászokhattunk, hogy a reptér egyenlő egy nagy épülettel, ami mellett óriási sík terület húzódik, ahol a gépek kényelmesen le- és felszállhatnak, mellől vagy utánán jól kigúrtulak magukat a hosszú pályán.

Na, ha egyszer megadatik önök, hogy ellátogathat Buzsábra, akkor mindent gyorsan felejtse el, különben elég nagy meglepetés fogja érni. Az ország első és a nagyjából még az egyetlen nemzetközi repülőtér, a Paró reptér bevetése ugyanis minden, csak nem kényelmes rutinfeladat.

A Paró érteke csak a bhután nemzeti légitársaság a Druk Air hasznáta is az első menetrend szerinti járat 1983-ban állt szolgálatba -, 2011 óta már néhány más cég is tart fenn járatokat. A leszállás ugyanakkor az nemén főtársaj adottsága miatt olyan nehéz, hogy a néhány évtel elköltő adatok szerint csak nyolc pilótának volt engedélye a világon, hogy landolhasson ott.

A reptér 2200 méteres magasságban fekszik, egy mély völgyben, amelyet 5000 méteres hegycsúcsok vesznek körbe, emiatt sokat és erős szélben kell manőverezni, miközben a gép a hegoldalaktól, majd a reptér körüli fűzektől is csak pár száz méteres távolságra. A landoláshoz pedig egy keskeny, keszletben mindössze 1,2, ma 2 kilométer hosszú pályát kell pontosan elkészíteni. Nem csoda, hogy csak napnál is jó látsá viszonyok között használható a reptér.

Type	Set	Start	End	Id
Token	0	10	4	(NER-BIO1-O, NP-BIO-I-NP, anas-{{ana-nyaral[V]-nyaral-ás[Ger/N]-ás[raSubj]-ra, feats-{{N[Subj], lemma-nyaralás}, {ana-nyaralás[N]-nyaralás+ra[Subj]-
Token	11	15	6	(NER-BIO1-O, NP-BIO-I-NP, anas-{{ana-neha[Adv]-neha, feats-{{Adv, lemma-neha}}, cons-{{ADVP}}, feature-SubPOS-xDeg-none Num-none Per-none, hfstat
Token	16	25	8	(ana-repulo[V]-repulo-ol, impPtcp[Adj]-o-vo[Ins]-vel, feats-{{Adj[Ins], lemma-repulo}, {ana-repulo[N]-repulo-vo[Ins]-vel, f
Token	26	31	10	(NER-BIO1-O, NP-BIO-I-NP, anas-{{ana-utaz[VF]-utaz-ol, impPtcp[Adj]-o-ol, ImpPtcp[Adj]-o-ol, ImpPtcp[Adj]-o-ol, ImpPtcp[Adj]-o-ol, Imp
Token	32	39	12	(NER-BIO1-O, NP-BIO-I-NP, anas-{{ana-turista[N]-turista-i, Adjz[Adj]-i-[Nomi], feats-{{Adj[Nomi], lemma-európai}, {ana-európai[Adj]-európai-[Nomi], feats
Token	40	51	14	(NER-BIO1-O, NP-BIO-E-NP, anas-{{ana-turista[N]-turista+kent[EssFor:kent]-kent, feats-{{N[EssFor:kent], lemma-turista}}, cons-}, feature-SubPOS-c[Num-s
Token	52	68	16	(NER-BIO1-O, NP-BIO-O, anas-{{ana-hozza[PreV]-hozza-zsokk[VF]-szok-hat[Mod/V]-hat-tunk[Pat.NDef.IPl]-tunk, feats-{{V[Mod/V]Pat.NDef.IPl], lemma-h
Token	68	69	17	(NER-BIO1-O, NP-BIO-O, anas-}, cons-}, feature-}, hfstana-OTHER, kind-punctuation, lemma-}, length-1, pos-, string-}
Token	70	74	19	(NER-BIO1-O, NP-BIO-O, anas-{{ana-abogy[AdvPro:tel]-hogy, feats-{{AdvPro:tel, lemma-abogy}, {ana-hogy[AdvPro:tel]-hogy, feats-{{AdvPro:tel, lemma-
Token	75	76	21	(NER-BIO1-O, NP-BIO-I-NP, anas-{{ana-a[DetPro:Post]-a, feats-{{DetPro:Post}, lemma-a}, {ana-a[Det:art.Def]-a, feats-{{Det:art.Def}, lemma-a}, {ana-
Token	77	83	23	(NER-BIO1-O, NP-BIO-E-NP, anas-{{ana-repter[N]-repter-[Nomi], feats-{{N[Nomi], lemma-repter}}, cons-}, feature-SubPOS-s[Cas-n Num-p none Per-p
Token	84	91	25	(NER-BIO1-O, NP-BIO-I-NP, anas-{{ana-egyenlo[Adj]-egyenlo-[Nomi], feats-{{Adj[Nomi], lemma-egyenlo}}, cons-{{NP}}, feature-SubPOS-c[Deg-p Num-s Cas-n
Token	92	95	27	(NER-BIO1-O, NP-BIO-I-NP, anas-{{ana-egy[Det:art.NDef]-egy, feats-{{Det:art.NDef, lemma-egy}, {ana-egy[Num]-egy-[Nomi], feats-{{Num[Nomi], lemma-
Token	96	100	29	(NER-BIO1-O, NP-BIO-I-NP, anas-{{ana-nagy[Adj]-nagy-[Nomi], feats-{{Adj[Nomi], lemma-nagy}, {ana-nagy[AdvAdMod]-nagy, feats-{{AdvAdMod, lemma-n
Token	101	110	31	(NER-BIO1-O, NP-BIO-I-NP, anas-{{ana-epulet[VF]-epulet-tel[Ins]-tel, feats-{{N[Ins], lemma-epulet}}, cons-}, feature-SubPOS-c[Num-s Cas-i Num-p none
Token	110	111	32	(NER-BIO1-O, NP-BIO-E-NP, anas-}, cons-}, feature-}, hfstana-OTHER, kind-punctuation, lemma-}, length-1, pos-, string-}
Token	112	115	34	(NER-BIO1-O, NP-BIO-O, anas-{{ana-ami[DetPro:Rel]-ami, feats-{{DetPro:Rel, lemma-ami}, {ana-ami[Pro:Rel]-ami-[Nomi], feats-{{Pro:Rel[Nomi], lemma-
Token	116	123	36	(NER-BIO1-O, NP-BIO-O, anas-{{ana-mellet[Post]-mellett, feats-{{Post, lemma-mellet}}, cons-}, feature-SubPOS-t-, hfstana-[Post], kind-word, lemma-m
Token	124	130	38	(NER-BIO1-O, NP-BIO-I-NP, anas-{{ana-erules[Adj]-erules-i, Adjz:ia[Adj]-i-[Nomi], feats-{{Adj[Nomi], lemma-erules}, {ana-erules[Adj]-erules-i-[Nomi], feats-{{Adj[N
Token	131	134	40	(NER-BIO1-O, NP-BIO-I-NP, anas-{{ana-sik[Adj]-sik-[Nomi], feats-{{N[Adj], lemma-sik}, {ana-sik[Adj]-sik-[Nomi], feats-{{Adj[Nomi], lemma-sik}, {ana-
Token	135	142	42	(NER-BIO1-O, NP-BIO-I-NP, anas-{{ana-terulet[N]-terulet-[Nomi], feats-{{N[Nomi], lemma-terulet}}, cons-}, feature-SubPOS-c[Num-s Cas-n Num-p none Per
Token	143	154	44	(NER-BIO1-O, NP-BIO-O, anas-{{ana-huz[V]-huz-od[Mod.MedPass.V]-od-ik[Prs.NDef.3sg]-ik, feats-{{N[Prs.NDef.3sg], lemma-huzodik}, {ana-huzodik[VF]-huzo
Token	150	151	45	(NER-BIO1-O, NP-BIO-O, anas-}, cons-}, feature-}, hfstana-OTHER, kind-punctuation, lemma-}, length-1, pos-, string-}
Token	152	156	47	(NER-BIO1-O, NP-BIO-O, anas-{{ana-ahol[AdvPro:Rel]-ahol, feats-{{AdvPro:Rel, lemma-ahol}}, cons-{{CPADVP}}, feature-SubPOS-xDeg-none Num-none Per
Token	157	158	49	(NER-BIO1-O, NP-BIO-I-NP, anas-{{ana-a[DetPro:Post]-a, feats-{{DetPro:Post}, lemma-a}, {ana-a[Det:art.Def]-a, feats-{{Det:art.Def}, lemma-a}, {ana-

233 Annotations (0 selected) Select:

Document Editor Initialisation Parameters Relation Viewer

- Sentence
- SpaceToken
- Token
- Original markups

New

e-magyar honlap

<http://e-magyar.hu>

e-magyar GATE

<https://github.com/dlt-rilmta/hunlp-GATE>

emTSV

<https://github.com/dlt-rilmta/emtsv>

magyarlanc

<http://rgai.inf.u-szeged.hu/index.php?lang=en&page=magyarlanc>

morfológiai címkekészletek

<https://github.com/dlt-rilmta/panmorph>