

Számítógépes szintaxis

DÖMÖTÖR ANDREA

DIGITÁLIS ÖRÖKSÉG NEMZETI LABORATÓRIUM
ELTE BTK TI DIGITÁLIS BÖLCSÉSZET TANSZÉK

2022. március 28.

Az előadás témakörei

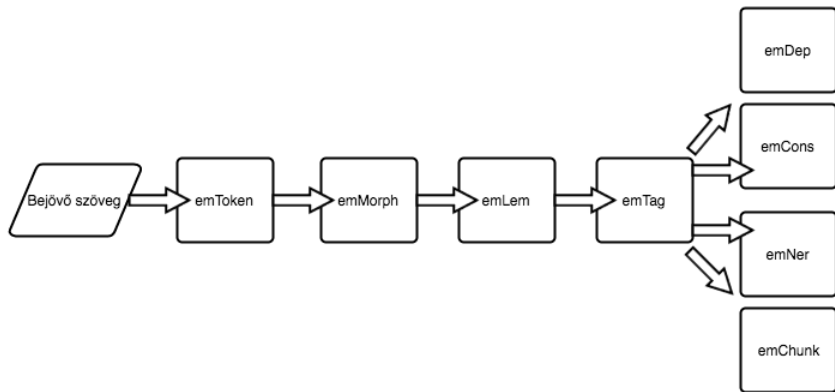
- A (számítógépes) nyelvi elemzés szintjei
- A mondat szerkezet megközelítései a számítógépes nyelvészetben
- Treebankek
- Elemző eszközök

A nyelvi elemzés lépései (szintjei)

- Szegmentálás, tokenizálás
- Szótövesítés (lemmatizálás)
- Morfológiai elemzés
- Szófaji egyértelműsítés
- Szintaktikai elemzés

Az egyes szintek **elemzési láncot** alkotnak. Az alacsonyabb szintek kimenete adja a magasabb szintek bemenetét (tehát pl. a szintaktikai elemzésnek feltétele a morfológiai elemzés és szófaji egyértelműsítés).

e-magyar pipeline



Szegmentálás, tokenizálás

A **szegmentálás** során a szöveget bekezdés vagy/és mondat egységekre bontjuk.

A **tokenizálás** az elemzési egységekre bontást jelenti.

- Token: “nyelvtechnológiai értelemben vett” szó
 - Az írásjelek külön tokenek
- Miért fontos ez?
 - A mondatelemzéshez nyilván meg kell határozni a mondatot...
 - Az *újság* és az *újság?* nem két különböző szó
 - Nem akarunk *újság?*-szerű szavakat

Szótövesítés (lemmatizálás)

A **szótövesítés** során a szótőt és a szófajt állapítjuk meg.

Miért fontos ez?

- Intelligens keresés
- A magyarban rengeteg szóalak tartozik egy szóhoz
 - → kezelhetetlen szótárméret
- Szófajspecifikus jegyek taníthatósága

Morfológiai elemzés

A **morfológiai elemző** a szavakat szóelemekre bontja, megadja a szóalakok lehetséges elemzéseit, a szó morfológiai jegyeit.

Fontos, hogy ezen a szinten **lehetséges elemzéseket** kapunk (az összeset, sőt...), és még nem tudjuk, melyik a jó.

Miért fontos ez?

- Elemzés nélkül túl sok szóalakunk lenne
- Az elemzett korpuszon tanított mesterséges intelligencia képes lehet:
 - ismeretlen szavak szóelemeit is felismerni
 - új szóalakokat generálni
- A szintaktikai elemzéshez szükség van a morfológiai jegyekre

Szófaji egyértelműsítés

A morfológiai elemző minden szóalakhoz megadja az összes lehetséges elemzést. A **szófaji egyértelműsítő** feladata, hogy ezek közül kiválassza a helyeset (a kontextus alapján legvalószínűbbet).

Miért fontos ez?

- Homonímiák feloldása
 - pl. *meg*: kötőszó vagy igekötő?
- Nyelv gazdaságossága → sok azonos alakú szó
 - Az emberi nyelvfeldolgozás is a kontextus alapján különbözteti meg őket.

Szintaktikai elemzés

A **szintaktikai elemző** a mondat szerkezetet határozza meg.

Ennek két fő megközelítése van:

- **Közvetlen összetevős elemzés**

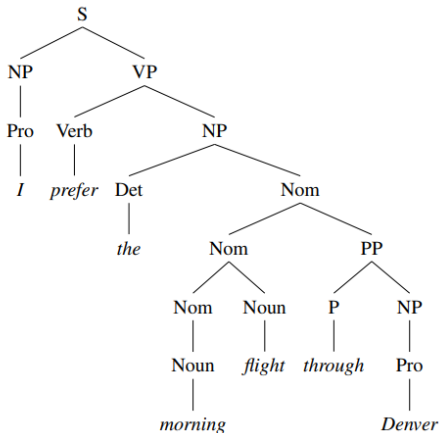
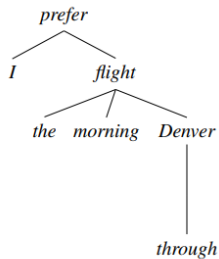
- A szavak frázisokat alkotnak (NP, VP, stb.)
- A frázisokból épül fel a mondat
- A szintaktikai szabályok a mondat és a frázisok lehetséges szerkezeteit írják le

- **Függőségi elemzés**

- A szavak közötti függőségeket és a relációk típusát adja meg
- Minden szónak van egy feje
- A szó és a feje közti viszony írja le a szó szintaktikai szerepét (pl. subj, obj...)
- A teljes elemzés mindig egy gyökércsomópontból indul ki, és fa struktúrában ábrázolható

Közvetlen összetevős elemzés

I prefer the morning flight to Denver.



Lexikon és szabályok

Lexikon	
Pro	<i>I, Denver</i>
Verb	<i>prefer</i>
Det	<i>the</i>
Noun	<i>morning, flight</i>
P	<i>through</i>

Szabályok
S → NP VP
NP → Pro
NP → Det Nom
VP → Verb NP
Nom → Nom PP
Nom → Nom Noun
Nom → Noun
PP → P NP

Top-down és bottom-up elemzés

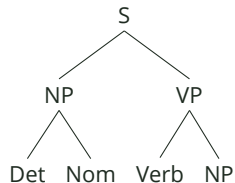
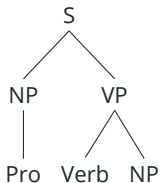
- **Bottom-up elemzés**

- A levelektől indul, és felfelé haladva építi a fát
- Ha eljutunk a mondatszimbólumig (S-ig), akkor a mondat grammatikus
- Előnyei:
 - Nem töltjük az időt a bemenettel nem kompatibilis szabályokkal
 - Minden részszerkezet kompatibilis a bemenet valamelyik részével

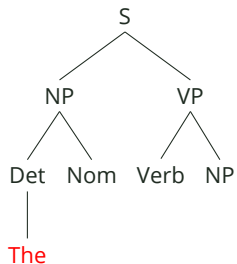
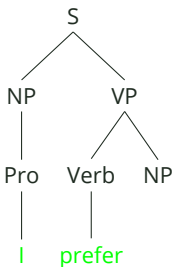
- **Top-down elemzés**

- A mondatszimbólumtól indul, és lefelé haladva építi a fát
- Ha eljutunk a levelekig, akkor grammatikus a mondat
- Előnyei:
 - Nem töltjük az időt nem mondatszimbólummal végződő részelemzésekkel
 - Minden részszerkezetnek helye lesz a fában

Top-down elemzés – példa



Top-down elemzés – példa

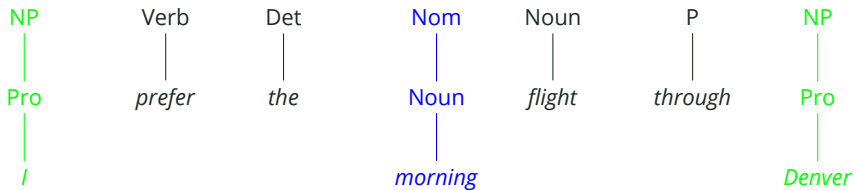


A jobb oldali elemzést elvetjük, mert nem kompatibilis a bemenettel. A bal oldalt folytatjuk az NP-vel az előzőekhez hasonlóan, amíg meg nem kapjuk az összes levelet.

Bottom-up elemzés – példa

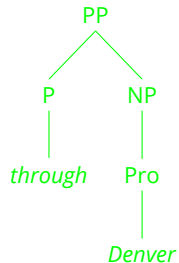
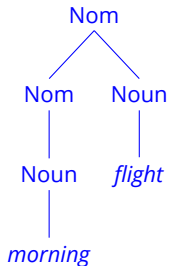


- Nom → Noun
- NP → Pro



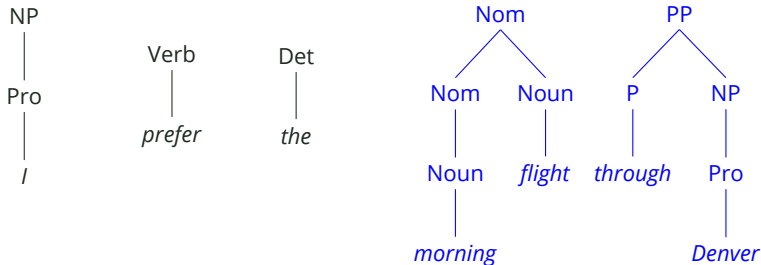
Bottom-up elemzés – példa

- Nom → Nom Noun
- PP → P NP



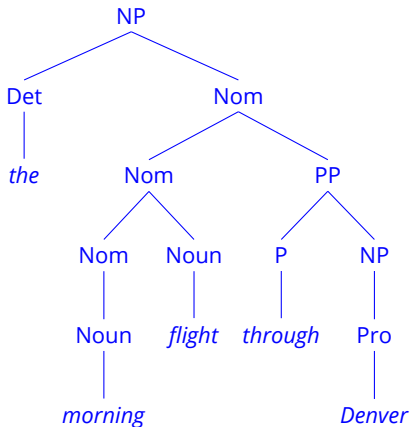
Bottom-up elemzés – példa

- Nom → Nom PP

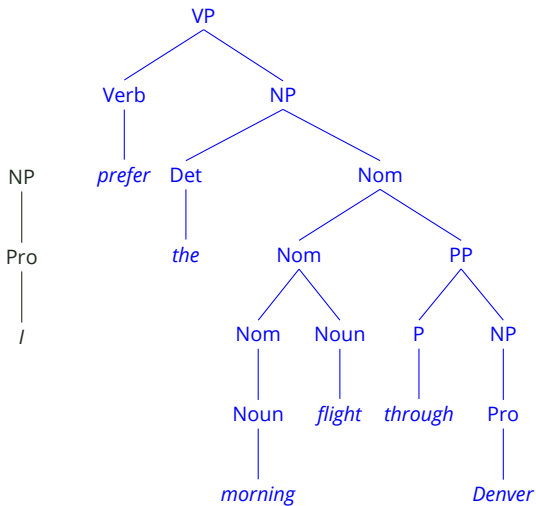


Bottom-up elemzés – példa

- NP → Det Nom

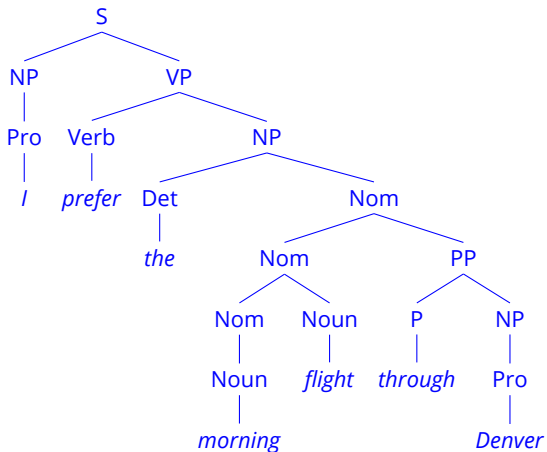


Bottom-up elemzés – példa



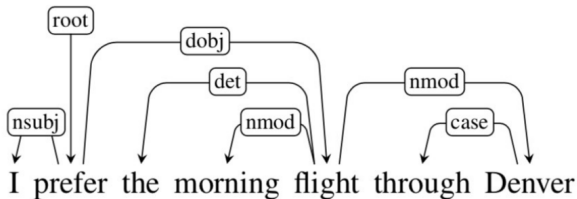
- VP → Verb NP

Bottom-up elemzés – példa



- $S \rightarrow NP VP$

Függőségi elemzés



- Minden szónak pontosan egy feje van
- A egy virtuális gyökércsomópontból (root) indul ki
- A szavak közötti élek címkéi a nyelvtani viszonyt írják le

Relációtípusok – példák

Clausal Argument Relations	Description
NSUBJ	Nominal subject
DOBJ	Direct object
IOBJ	Indirect object
CCOMP	Clausal complement
XCOMP	Open clausal complement
Nominal Modifier Relations	Description
NMOD	Nominal modifier
AMOD	Adjectival modifier
NUMMOD	Numeric modifier
APPOS	Appositional modifier
DET	Determiner
CASE	Prepositions, postpositions and other case markers
Other Notable Relations	Description
CONJ	Conjunct
CC	Coordinating conjunction

Összetevős vs. függőségi elemzés

Összetevős elemzés	Függőségi elemzés
mi mit előzhet meg	mi mitől függ
szavak → frázisok	szavak → bináris relációk
magasabb szintű kategóriák lehetősége	szintaktikai viszonyok kifejezhetősége
érzékeny a szórendre	nem érzékeny a szórendre

- A gazdag morfológiájú és/vagy szabad szórendű nyelvek számára előnyösebb a függőségi elemzés
- Általában is ez az uralkodó megközelítés

Universal Dependencies

A függőségi annotációt is tartalmazó korpuszokat **treebank**nek nevezzük. Magyar nyelvre jelenleg a legnagyobb ilyen korpusz a [Szeged Treebank](#).

A treebankek annotációjának egységesítésére jött létre a [Universal Dependencies](#) projekt.

- Nemzetközi szabvány
- Célja: a különböző nyelvű treebankek annotációja kompatibilis legyen egymással
- → univerzális nyelvtant feltételez
- Nyílt, közösségi projekt
 - ~200 résztvevő
 - >100 nyelv (köztük a magyar)

CoNLL-U formátum

CoNLL: Conference on Natural Language Learning

A CoNLL fájlformátumok a CoNLL konferencia "shared task"-jaihoz tartozó fájlok formátumai.

- Vertikális fájl (tsv)
 - Minden szó (token) külön sorban
 - A szó tulajdonságai külön oszlopokban
- Többféle verziója létezik
- **CoNLL-U:** a Universal Dependencies által használt formátum

CoNLL-U fájlok mezői

- **ID:** egyedi azonosító
- **Form:** eredeti szóalak (token)
- **Lemma:** szótő
- **Upos:** szófajcímke
- **Xpos:** nyelvspecifikus szófajcímke
- **Feats:** morfológiai jegyek
- **Head:** a szó feje (id)
- **Deprel:** a függőségi viszony típusa
- **DepS:** továbbfejlesztett függőségi gráf
- **Misc:** egyéb annotációk

CoNLL-U példa

Szöveg: Kosztolányi: Esti Kornél

Elemző: UDPipe

```
# sent_id = 5278
# text = Meg volt gyözdöve, hogy maradandó alkotása, s évek múlva is boldogan gondol erre a téli estére, amikor a semmiből létrehívta.
1 Meg meg PART _ _ 2 compound:preverb _ _ TokenRange=369083:369086
2 volt van VERB _ _ Definite=Ind|Mood=Ind|Number=Sing|Person=3|Tense=Past|VerbForm=Fin|Voice=Act 0 root _ _ TokenRange=369087:369091
3 gyözdöve gyözdöve ADV _ _ VerbForm=Conv 2 advmod:mode _ _ SpaceAfter=No|TokenRange=369092:369100
4 , , PUNCT _ _ 2 punct _ _ TokenRange=369100:369101
5 hogy hogy SCONJ _ _ 2 punct 7 mark _ _ TokenRange=369102:369106
6 maradandó maradandó ADJ _ _ Case=Nom|Number=Sing|Number[psd]=None|Number[psor]=None|Person[psor]=None|VerbForm=PartFut 7 amod:att _ _ TokenRange=369107:369116
7 alkotása alkotás NOUN _ _ Case=Nom|Number=Sing|Number[psd]=None|Number[psor]=Sing|Person[psor]=3 2 nsubj _ _ SpaceAfter=No|TokenRange=369117:369125
8 , , PUNCT _ _ 14 punct _ _ TokenRange=369125:369126
9 s s CCONJ _ _ 14 cc _ _ TokenRange=369127:369128
10 évek év NOUN _ _ Case=Nom|Number=Plur|Number[psd]=None|Number[psor]=None|Person[psor]=None 14 obl _ _ TokenRange=369129:369133
11 múlva múlva ADP _ _ 18 case _ _ TokenRange=369134:369139
12 is is CCONJ _ _ 18 cc _ _ TokenRange=369140:369142
13 boldogan boldog ADJ _ _ Case=Ess|Degree=Pos|Number=Sing|Number[psd]=None|Number[psor]=None|Person[psor]=None 14 amod:mode _ _ TokenRange=369143:369151
14 gondol gondol VERB _ _ Definite=Ind|Mood=Ind|Number=Sing|Person=3|Tense=Pres|VerbForm=Fin|Voice=Act 2 conj _ _ TokenRange=369152:369158
15 erre ez PRON _ _ Case=Sub|Number=Sing|Number[psd]=None|Number[psor]=None|Person=3|Person[psor]=None|PronType=Dem 18 nmod:obl _ _ TokenRange=369159:369163
16 a a DET _ _ Definite=Def|PronType=Art 18 det _ _ TokenRange=369164:369165
17 téli téli ADJ _ _ Case=Nom|Degree=Pos|Number=Sing|Number[psd]=None|Number[psor]=None|Person[psor]=None 18 amod:att _ _ TokenRange=369166:369178
18 estére est NOUN _ _ Case=Sub|Number=Sing|Number[psd]=None|Number[psor]=Sing|Person[psor]=3 14 nmod:obl _ _ SpaceAfter=No|TokenRange=369171:369177
19 , , PUNCT _ _ 14 punct _ _ TokenRange=369177:369178
20 amikor amikor ADV _ _ PronType=Rel 23 advmod:tlocy _ _ TokenRange=369179:369185
21 a a DET _ _ Definite=Def|PronType=Art 22 det _ _ TokenRange=369186:369187
22 semiből semmi PRON _ _ Case=El|Number=Sing|Number[psd]=None|Number[psor]=None|Person[psor]=None 23 nmod:obl _ _ TokenRange=369188:369196
23 létrehívta létrehív VERB _ _ Definite=Def|Mood=Ind|Number=Sing|Person=3|Tense=Past|VerbForm=Fin|Voice=Act 14 advcl _ _ SpaceAfter=No|TokenRange=369197:369207
24 . . PUNCT _ _ 2 punct _ _ TokenRange=369207:369208
```

Elemző eszközök

e-magyar

- Integrált elemzőlánc
 - Tokenizálás, morfológiai elemzés, összetevős elemzés, függőségi elemzés, főnévcsoport-felismerés, névelem-felismerés
- Webes felület
 - Vagy: [emtsv](#) parancssori alkalmazás/Python csomag sok egyéb modullal
- Saját kódkészlet ([emMorph](#))
 - UD-re konvertálható
- tsv, xml kimenet

Elemző eszközök

UDPipe

- Cseh fejlesztés
- Több nyelvre (köztük magyarra) tanított elemzőlánc
 - Tokenizálás, morfológiai elemzés, függőségi elemzés
- Könnyen kezelhető webes felület
- CoNLL-U (tsv) kimenet

Stanza

- Stanford NLP fejlesztés
- Szövegfeldolgozásra alkalmas Python-csomag
- Magyar nyelvre is elérhető

Keresés a függőségi annotációban – demo

CoNLL vonzatkereső **demo** (colab notebook)

Feladatok

1. **Összetevős elemzés** papíros

- dravida_fak.pdf

2. **Függőségi elemzés** programozós

Csinálj valamilyen összehasonlító kimutatást a vonzatkereső demo segítségével (pl.: a *kap* ige gyakori tárgyai irodalmi szövegekben és hírportál szövegekben)!

- Elemezd a szövegeket UDPipe-pal!
- Futtasd le a demo kódot a szükséges paraméterekkel!
- Írj egy rövid összefoglalót a kapott eredményekről!

3. **Szorgalmi** extra programozós

Írj egy függvényt a demo kódba, amely a paraméterben megadott lemma fejeiről készít gyakorisági listát! Csinálj ezzel is valamilyen összehasonlítást (pl. milyen fejek tartoznak az *orosz* és az *ukrán* szavakhoz)!

Köszönöm a figyelmet!

domotor.andrea@btk.elte.hu

<https://dh-lab.hu/>

<https://elte-dh.hu/>

<https://github.com/elte-dh>

<https://zenodo.org/communities/elte-dh>