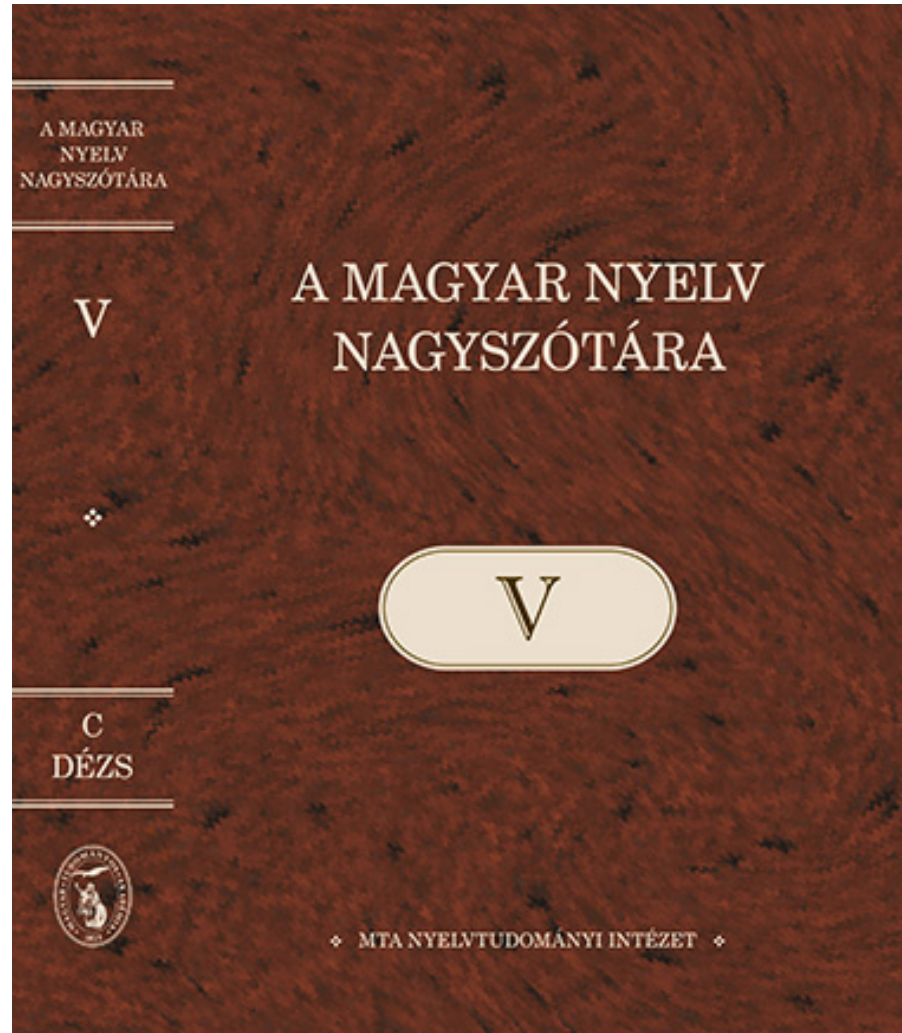


Korpusznyelvészet

2016. április 18., ELTE

Sass Bálint
MTA Nyelvtudományi Intézet
`sass.balint@nytud.mta.hu`



<http://nszt.nytud.hu/nszt.html>

Mi mindent kell csinálni ahhoz,
hogy sima szövegből ilyen korpuszlekérdezőfelület legyen?
Ezt fogjuk most megnézni lépésről lépésre.

Példa:

/ Ám de vizont hallá, hogy majd a' Trójai vérből /

1.

Az Mtsz építése

Karakterkódolás

1 karakter = 1 byte

kitûnõ és idõszerû

1 karakter > 1 byte: Unicode

helytakarékos kódolás: UTF-8

HÃ©tfõn találkozunk.

Megvan az egységes UTF-8 kódolású szövegünk.

é ✓ ő ✓ f ✓ ö ✓

XML

Kenjük a vaját

a késsel

a kenyérre.

XML – tagek

<recept>

Kenjük a **<hozzavaló>**vajat**</hozzavaló>**

a **<eszköz>**késsel**</eszköz>**

a **<hozzavaló>**kenyérre**</hozzavaló>**.

</recept>

XML – attribútumok

```
<recept nev="vajas kenyér">
```

```
Kenjük a <hozzavalo id="41">vaját</hozzavalo>
```

```
a <eszkoz id="5">késsel</eszkoz>
```

```
a <hozzavalo id="12">kenyérre</hozzavalo>.
```

```
</recept>
```


Mtsz XML

```
<section>
  <head>
    <id>7021030</id>
    <author>Baróti Szabó Dávid</author>
    <wdate>1808</wdate>
    ... egyéb adatok
  </head>
  <text>
    <page>
      <p>18</p>
      <par>
        Ám de vizont hallá, hogy majd a' Trójai vérből<br/>
        Nemzet ered, melly e' várat valahára le-dönti;<br/>
        ... további sorok
      </par>
      ... további bekezdések
    </page>
    ... további oldalak
  </text>
</section>
```

TEI

„The Text Encoding Initiative (TEI) is a consortium which collectively develops and maintains a standard for the representation of texts in digital form.”

Korpuszkezelő

NoSketchEngine (NoSkE)

<https://nlp.fi.muni.cz/trac/noske>

NoSkE XML

```
<doc mtsz_id="7021030" author="Baróti Szabó Dávid" wdate="1808" ...>
  <oldal oldalszam="18">
    <par>
      Ám de vízfönt hallá, hogy majd a' Trójai vérből<br/>
      Nemzet ered, melly e' várat valahára le-dönti;<br/>
      ... további sorok
    </par>
    ... további bekezdések
  </oldal>
  ... további oldalak
</doc>
```

NoSkE XML ← XSLT ← Mtsz XML

```
<doc mtsz_id="7021030" author="Baróti Szabó Dávid" wdate="1808" ...>
  <oldal oldalszam="18">
    <par>
      Ám de vízfönt hallá, hogy majd a' Trójai vérből<br/>
      Nemzet ered, melly e' várat valahára le-dönti;<br/>
      ... további sorok
    </par>
    ... további bekezdések
  </oldal>
  ... további oldalak
</doc>
```

Találatok időrendben

Hogy lehet a találatok idő szerinti rendezését megvalósítani?

1. Lekérdezés után az eredményt mindig rendezzük.
2. Előre rendezzük és azt kérdezzük le.

Megoldás (erre is): XSLT

„Tokenizálás”

```
<doc mtsz_id="7021030" author="Baróti Szabó Dávid" wdate="1808" ...>
  <oldal oldalszam="18">
    <par>
      Ám
      de
      vifzont
      hallá
      <g/>
      '
      hogy
      ...
      <br/>
      ... további tokenek
    </par>
    ... további bekezdések
  </oldal>
  ... további oldalak
</doc>
```

„XML+TAB” formátum

```
<doc mtsz_id="7021030" author="Baróti Szabó Dávid" wdate="1808" ...>
  <oldal oldalszam="18">
    <par>
      Ám      ám      KOT
      de      de      KOT
      vifzont vifzont KOT
      hallá   hall   V.Ipf.S3.Def
    </g/>
      ,      ,      WPUNCT
      hogy    hogy    KOT
      ...
    <br/>
      ... további tokenek
    </par>
      ... további bekezdések
  </oldal>
  ... további oldalak
</doc>
```

2.

Az Mtsz használata

Reguláris kifejezések

Bizonyos tulajdonságú karaktersorozatok megadására.

- . tetszőleges karakter
- * a megelőző karakterből 0 vagy több
- + a megelőző karakterből 1 vagy több
- ? a megelőző karakterből 0 vagy 1
- [ab] 'a' vagy 'b' karakter
- [^ab] nem 'a' és nem is 'b' karakter
- r|s 'r' vagy 's' reguláris kifejezés
- (..) egybefoglalás

Példák:

- | | |
|---------------|------------------|
| 1. alma | 5. .* |
| 2. tejf.l | 6. .*bb |
| 3. mentők? | 7. alma almá.* |
| 4. nélk[üúű]l | 8. mondjá(to)?k |

CQL (Corpus Query Language)

[..] egy tokenre vonatkozó megkötések

$x="y"$ x attrib értéke legyen y – Mtsz: csak *word* attrib van

$x!="y"$ x attrib értéke *ne* legyen y

& és kapcsolat megkötések között

Példák:

1. [] []

2. [word="majd"]

3. "majd"

4. [word!="a.*"]

5. []?

6. [word="nem"] [word="kellett"] [word="volna"]? [word=".*ni"]

Mtsz példalekérdezés

Feladat. Keressünk olyet: tárgyesetű szó + múltidejű E/3 ige!

Mtsz példalekérdés

Feladat. Keressünk ilyen: tárgyestű szó + múltidejű E/3 ige!

" . * t " " . * . . t t "

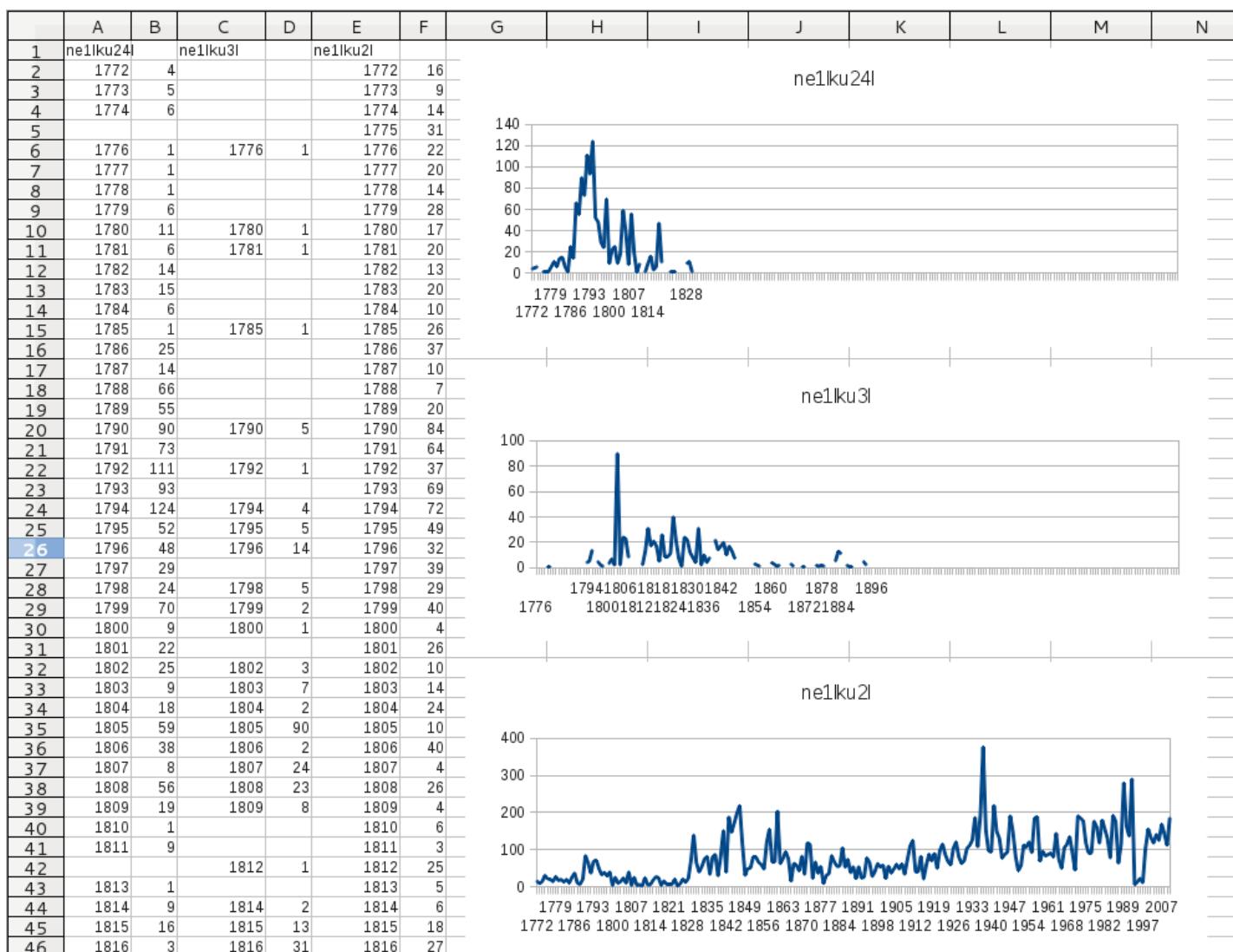
Mtsz példalekérdezés

Feladat. Keressünk olyet: tárgyesetű szó + múltidejű E/3 ige!

"*.t" "*.tt"

[word="*.t"] [word="*.tt" & word!="alatt" & word!="azelőtt"]

Diakrón vizsgálat: *nélkül* helyesírása



3.

Korpuszok

Korpuszok

MNSZ2 – elemzett, nagy méretű

- * körülültük, felszededegettük, elsimítottuk, végigcsináltuk
- * forrásokban, fellegekben, falvakban, fejekben (bazi lassú)
- * cél, csal, csaj, csel, dzsal

Mazsola – igék és bővítmények

reprezentáció: A lány vállat vont → ige=von alany=lány tárgy=váll

BUSZI – beszélt de írott

... bizonyos dógokban □ mmm tát, hogy ööö lustább annál, mint amilyennek elképzeltem, ...

Ómagyar korpusz – normalizálás, ómagyar morfológia

NKP (Nemzeti Korpuszportál)

<http://corpus.nytud.hu/nkp>

4.

Korpuszvezérelt kutatás

Korpuszalapú és korpuszvezérelt

„A korpusz segédeszköz, ami empirikus adataival támogatja az intuíciót, mérhetővé teszi a nyelvi jelenségeket, meglévő elméleteket bizonyít/cáfol.”

„A korpusz maga szolgáltatja az «elméletet», a nyelvész minden feltevés és elvárás nélkül fordul az adatokhoz. Minden következtetést kizárólag korpuszmegfigyelésekből von le.”

serendipity principle:

lényeges jelenség véletlen felfedezése

pl.: Tognini-Bonelli: *Corpus Linguistics at Work* (2001)

Braille-rövidírás bővítése

ban/ben → b (⠠⠠⠠ → ⠠) hoG → h (⠠⠠⠠ → ⠠)

Alapötlet: a maximális rövidítési képességgel bíró ideális rövidírás a magyar nyelv korpuszgyakorisági adatai alapján **korpuszvezérelt** módon, automatikusan kialakítható.

Elv: a lehető leggyakoribb *betűkapcsolatokat* kell a lehető legrövidebbre rövidíteni.

Eredmény: 33 új szabály

meg → mg (⠠⠠⠠ → ⠠⠠) maGar → mG (⠠⠠⠠⠠⠠ → ⠠⠠)

Tanulság:

A korpuszvezérelt módon létrehozott rendszer még úgy is kétszeres teljesítményre képes az intuíció illetve hagyomány talaján álló rendszerrel szemben, hogy már eleve jelentősen rövidített szövegen kell dolgoznia.

Ha valamit meglévő (korpuszgyakorisági) adatokból automatikusan származtatni tudunk, akkor nem érdemes intuitív megközelítést alkalmazni.

Igei szerkezetek felfedezése

vet [15728]

vet -nAk VÉG-t [1463] *vessen véget az erőszaknak*

vet SZEM-A-rA -t [805] *hasonló diszkriminációkat vetnek az albán hatóságok szemére*

vet -rA PILLANTÁS-t [708] *vess egy pillantást a térképre*

vet -t [703] *vetem a magot*

vet -rA -t [380] *a humanista könyveket máglyára vetették*

vet PAPÍR-rA -t [377] *vesse papírra az új problémákat*

vet SZÁM-t -vAl [297] *ez rossz fényt vet az edzők nevelő munkájára*

vet -rA FÉNY-t [267] *vessünk számot eddigi politikánkkal*

vet -bA -t [252] *a tó vizébe vetette magát*

csóvál [1078]

csóvál FEJ-A-t [754] *csóválta a fejét*

Korpusznyelvészet

2016. április 18., ELTE

Sass Bálint
MTA Nyelvtudományi Intézet
`sass.balint@nytud.mta.hu`