

Miért van szükség a nyelvtechnológiára?

Ács Judit

BME AUT

természetes nyelvfeldolgozás
natural language processing
számítógépes nyelvészet
NLP
nyelvtechnológia

AND THE DUMBEST THING ABOUT
EMO KIDS IS THAT... I...
YOU KNOW, I'M SICK OF EASY TARGETS.
ANYONE CAN MAKE FUN OF EMO KIDS.
YOU KNOW WHO'S HAD IT TOO EASY?
COMPUTATIONAL LINGUISTS.



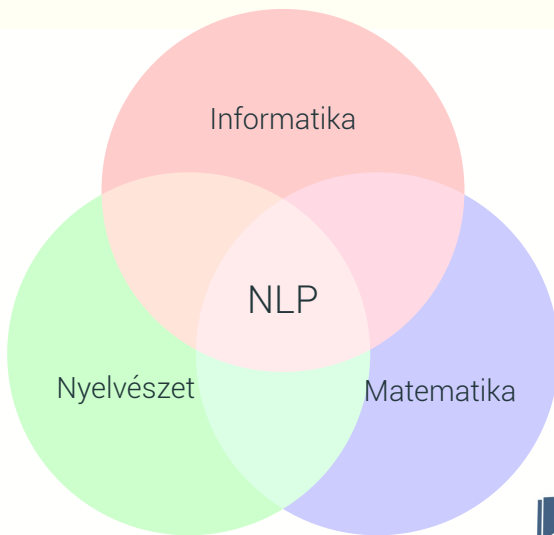
"OOH, LOOK AT ME!
MY FIELD IS SO IL-DEFINED
I CAN SUBSCRIBE TO ANY OF
DOZENS OF CONTRADICTIONARY
MODELS AND STILL BE
TAKEN SERIOUSLY!"



<https://xkcd.com/114/>



Schönherz



Definíció?



Schönherz

Definíció?

- az input és/vagy az output természetes nyelv



Definíció?

- az input és/vagy az output természetes nyelv
- formája sokféle lehet: írott, beszélt, jelnyelv, Braille stb.





- Fordítsuk le az alábbi mondatokat magyarra!

He went back to New York and left Mary behind.
She was devastated.

Fordítás szavanként

He went back to New York and left Mary behind. She was devastated.



Schönerz

Fordítás szavanként

He went back to New York and left Mary behind. She was devastated.

Ő ment hát nak/nek Új York és van hátra Mary mögött. Ő volt elpusztított.



Schönherz

Fordítás szavanként

He went back to New York and left Mary behind. She was devastated.

Ő ment hát nak/nek Új York és van hátra Mary mögött. Ő volt elpusztított.

- a **back** szó kétértelmű (vissza, hát),



Fordítás szavanként

He went back to New York and left Mary behind. She was devastated.

Ő ment hát nak/nek Új York és van hátra Mary mögött. Ő volt elpusztított.

- a *back* szó kétértelmű (vissza, hát),
- *left, devastated* is többértelműek



Schönherz

Fordítás szavanként

He went back to New York and left Mary behind. She was devastated.

Ő ment hát nak/nek Új York és van hátra Mary mögött. Ő volt elpusztított.

- a *back* szó kétértelmű (vissza, hát),
 - *left, devastated* is többértelműek
- *New York* egy tulajdonnév,

Fordítás szavanként

He went back to New York and left Mary behind. She was devastated.

Ő ment hát nak/nek Új York és van hátra Mary mögött. Ő volt elpusztított.

- a **back** szó kétértelmű (vissza, hát),
 - **left, devastated** is többértelműek
- **New York** egy tulajdonnév,
- **She** kire vonatkozik? Angolul egyértelmű, magyarul nem,

Fordítás szavanként

He went back to New York and left Mary behind. She was devastated.

Ő ment hát nak/nek Új York és van hátra Mary mögött. Ő volt elpusztított.

- a **back** szó kétértelmű (vissza, hát),
 - **left, devastated** is többértelműek
- **New York** egy tulajdonnév,
- **She** kire vonatkozik? Angolul egyértelmű, magyarul nem,
- **went back** együtt fordítandó
- szórend, toldalékok

Fordítás szavanként

He went back to New York and left Mary behind. She was devastated.

Ő ment hát nak/nek Új York és van hátra Mary mögött. Ő volt elpusztított.

- a **back** szó kétértelmű (vissza, hát), **word sense disambiguation**
 - *left, devastated* is többértelműek
- **New York** egy tulajdonnév,
- **She** kire vonatkozik? Angolul egyértelmű, magyarul nem,
- **went back** együtt fordítandó
- szórend, toldalékok

Fordítás szavanként

He went back to New York and left Mary behind. She was devastated.

Ő ment hát nak/nek Új York és van hátra Mary mögött. Ő volt elpusztított.

- a **back** szó kétértelmű (vissza, hát), **word sense disambiguation**
 - **left, devastated** is többértelműek
- **New York** egy tulajdonnév, **named entity recognition**
- **She** kire vonatkozik? Angolul egyértelmű, magyarul nem,
- **went back** együtt fordítandó
- szórend, toldalékok

Fordítás szavanként

He went back to New York and left Mary behind. She was devastated.

Ő ment hát nak/nek Új York és van hátra Mary mögött. Ő volt elpusztított.

- a **back** szó kétértelmű (vissza, hát), **word sense disambiguation**
 - **left, devastated** is többértelműek
- **New York** egy tulajdonnév, **named entity recognition**
- **She** kire vonatkozik? Angolul egyértelmű, magyarul nem, **anaphora resolution**
- **went back** együtt fordítandó
- szórend, toldalékok

He went back to New York and left Mary behind. She was devastated.



He went back to New York and left Mary behind. She was devastated.

Visszament New York-i és a bal Mary mögött. Ő volt elpusztított.



He went back to New York and left Mary behind. She was devastated.

Visszament New York-i és a bal Mary mögött. Ő volt elpusztított.

- *New York-i* lett a *to New Yorkból*,



He went back to New York and left Mary behind. She was devastated.

Visszament New York-i és a bal Mary mögött. Ő volt elpusztított.

- *New York-i* lett a *to New Yorkból*,
- *left* rossz jelentését fordítja,

He went back to New York and left Mary behind. She was devastated.

Visszament New York-i és a bal Mary mögött. Ő volt elpusztított.

- *New York-i* lett a *to New York*ból,
- *left* rossz jelentését fordítja,
- *left behind*-ot szó szerint értelmezi?

He went back to New York and left Mary behind. She was devastated.

Visszament New York-i és a bal Mary mögött. Ő volt elpusztított.

- *New York-i* lett a *to New York*ból,
- *left* rossz jelentését fordítja,
- *left behind*-ot szó szerint értelmezi?
- a 2. mondatról inkább ne beszéljünk.



tokenizálás szöveg kisebb egységekre bontása (mondat, szó).



tokenizálás szöveg kisebb egységekre bontása (mondat, szó).

- *New York-i* hány token?
- elválasztás
- dátumok kezelése stb.



tokenizálás szöveg kisebb egységekre bontása (mondat, szó).

- *New York-i* hány token?
- elválasztás
- dátumok kezelése stb.

szótövezés toldalékok (képző, rag) eltávolítása. *házaink* - ház



tokenizálás szöveg kisebb egységekre bontása (mondat, szó).

- *New York-i* hány token?
- elválasztás
- dátumok kezelése stb.

szótövezés toldalékok (képző, rag) eltávolítása. *házaink* - ház

- *kályhák, ettem*

tokenizálás szöveg kisebb egységekre bontása (mondat, szó).

- *New York-i* hány token?
- elválasztás
- dátumok kezelése stb.

szótövezés toldalékok (képző, rag) eltávolítása. *házaink* - ház

- *kályhák, ettem*

morfológiai elemzés

- *zúzalékával*



tokenizálás szöveg kisebb egységekre bontása (mondat, szó).

- *New York-i* hány token?
- elválasztás
- dátumok kezelése stb.

szótövezés toldalékok (képző, rag) eltávolítása. *házaink* - ház

- *kályhák, ettem*

morfológiai elemzés • *zúzalékával*

névelemazonosítás named entity recognition

tokenizálás szöveg kisebb egységekre bontása (mondat, szó).

- *New York-i* hány token?
- elválasztás
- dátumok kezelése stb.

szótövezés toldalékok (képző, rag) eltávolítása. *házaink - ház*

- *kályhák, ettem*

morfológiai elemzés • *zúzalékával*

névelemazonosítás named entity recognition

- *Catholic* angolul, *katolikus* magyarul



szófaji elemzés part-of-speech tagging



szófaji elemzés part-of-speech tagging

„nyelvtani elemzés” nyelvtani szerkezet feltárása. Különböző nyelvészeti elméletek különböző elemzéseket adnak.

szófaji elemzés part-of-speech tagging

„nyelvtani elemzés” nyelvtani szerkezet feltárása. Különböző nyelvészeti elméletek különböző elemzéseket adnak.

egyértelműsítés word sense disambiguation



gépi fordítás rule-based MT, statistical MT, neural MT



Magasabb szintű NLP feladatok

gépi fordítás rule-based MT, statistical MT, neural MT

koreferencia feloldás mik utalnak ugyanarra az entitásra? *New York – Big Apple*



Schönherz

Magasabb szintű NLP feladatok

gépi fordítás rule-based MT, statistical MT, neural MT

koreferencia feloldás mik utalnak ugyanarra az entitásra? *New York – Big Apple*

szentiment analízis vélemény, érzelem kinyerése



Schönherz

Magasabb szintű NLP feladatok

gépi fordítás rule-based MT, statistical MT, neural MT

koreferencia feloldás mik utalnak ugyanarra az entitásra? *New York – Big Apple*

szeniment analízis vélemény, érzelem kinyerése

kapcsolatkinyerés relation extraction, ki kinek a rokona?

Magasabb szintű NLP feladatok

gépi fordítás rule-based MT, statistical MT, neural MT

koreferencia feloldás mik utalnak ugyanarra az entitásra? *New York – Big Apple*

szeniment analízis vélemény, érzelem kinyerése

kapcsolatkinyerés relation extraction, ki kinek a rokona?

Magasabb szintű NLP feladatok

gépi fordítás rule-based MT, statistical MT, neural MT

koreferencia feloldás mik utalnak ugyanarra az entitásra? *New York – Big Apple*

szentiment analízis vélemény, érzelem kinyerése

kapcsolatkinyerés relation extraction, ki kinek a rokona?

kérdésmegválaszolás

szövegösszegzés text summarization

Magasabb szintű NLP feladatok

gépi fordítás rule-based MT, statistical MT, neural MT

koreferencia feloldás mik utalnak ugyanarra az entitásra? *New York – Big Apple*

szeniment analízis vélemény, érzelem kinyerése

kapcsolatkinyerés relation extraction, ki kinek a rokona?

kérdésmegválaszolás

szövegösszegzés text summarization

helyesírás-ellenőrzés

Hogyan oldjuk meg ezeket a feladatokat?



Schönherz

Hogyan oldjuk meg ezeket a feladatokat?

1. Szabályalapú



Schönherz

Hogyan oldjuk meg ezeket a feladatokat?

1. Szabályalapú

- kézzel írt szabályokat keresünk



Hogyan oldjuk meg ezeket a feladatokat?

1. Szabályalapú

- kézzel írt szabályokat keresünk
- nagyon időigényes



Schönherz

Hogyan oldjuk meg ezeket a feladatokat?

1. Szabályalapú

- kézzel írt szabályokat keresünk
- nagyon időigényes
- nem skálázódik



Schönherz

Hogyan oldjuk meg ezeket a feladatokat?

1. Szabályalapú

- kézzel írt szabályokat keresünk
- nagyon időigényes
- nem skálázódik
- jó pontosság, de milyen a fedés?



Schönherz

Hogyan oldjuk meg ezeket a feladatokat?

1. Szabályalapú

- kézzel írt szabályokat keresünk
- nagyon időigényes
- nem skálázódik
- jó pontosság, de milyen a fedés?

2. Gépi tanulás



Hogyan oldjuk meg ezeket a feladatokat?

1. Szabályalapú

- kézzel írt szabályokat keresünk
- nagyon időigényes
- nem skálázódik
- jó pontosság, de milyen a fedés?

2. Gépi tanulás

- felügyelt vs. felügyeletlen



Schönherz

Hogyan oldjuk meg ezeket a feladatokat?

1. Szabályalapú

- kézzel írt szabályokat keresünk
- nagyon időigényes
- nem skálázódik
- jó pontosság, de milyen a fedés?

2. Gépi tanulás

- felügyelt vs. felügyeletlen
- a tanítóadat a legtöbb feladathoz nagyon költséges

Hogyan oldjuk meg ezeket a feladatokat?

1. Szabályalapú

- kézzel írt szabályokat keresünk
- nagyon időigényes
- nem skálázódik
- jó pontosság, de milyen a fedés?

2. Gépi tanulás

- felügyelt vs. felügyeletlen
- a tanítóadat a legtöbb feladathoz nagyon költséges
- nem felügyelt módszerekhez sok adatunk van

Hogyan oldjuk meg ezeket a feladatokat?

1. Szabályalapú

- kézzel írt szabályokat keresünk
- nagyon időigényes
- nem skálázódik
- jó pontosság, de milyen a fedés?

2. Gépi tanulás

- felügyelt vs. felügyeletlen
- a tanítóadat a legtöbb feladathoz nagyon költséges
- nem felügyelt módszerekhez sok adatunk van
- ***zaj!***

Hogyan oldjuk meg ezeket a feladatokat?

1. Szabályalapú

- kézzel írt szabályokat keresünk
- nagyon időigényes
- nem skálázódik
- jó pontosság, de milyen a fedés?

2. Gépi tanulás

- felügyelt vs. felügyeletlen
- a tanítóadat a legtöbb feladathoz nagyon költséges
- nem felügyelt módszerekhez sok adatunk van
- **zaj!**
- itt pezseg jobban az élet

Hogyan oldjuk meg ezeket a feladatokat?

1. Szabályalapú

- kézzel írt szabályokat keresünk
- nagyon időigényes
- nem skálázódik
- jó pontosság, de milyen a fedés?

2. Gépi tanulás

- felügyelt vs. felügyeletlen
- a tanítóadat a legtöbb feladathoz nagyon költséges
- nem felügyelt módszerekhez sok adatunk van
- **zaj!**
- itt pezseg jobban az élet
- deep learning



- címkézett



- címkézett
 - névelemekre:



- címkézett
 - névelemekre:
 - New York - B-LOC E-LOC
 - Google - ORG



- címkézett
 - névelemekre:
 - New York - B-LOC E-LOC
 - Google - ORG
 - gold standard, silver standard



- címkézett
 - névelemekre:
 - New York - B-LOC E-LOC
 - Google - ORG
 - gold standard, silver standard
- címkézetlen



- címkézett
 - névelemekre:
 - New York - B-LOC E-LOC
 - Google - ORG
 - gold standard, silver standard
- címkézetlen
 - nyers szöveg
 - kontextusból nyerhetünk információt

- címkézett
 - névelemekre:
 - New York - B-LOC E-LOC
 - Google - ORG
 - gold standard, silver standard
- címkézetlen
 - nyers szöveg
 - kontextusból nyerhetünk információt
- rengeteg standard adat

- címkézett
 - névelemekre:
 - New York - B-LOC E-LOC
 - Google - ORG
 - gold standard, silver standard
- címkézetlen
 - nyers szöveg
 - kontextusból nyerhetünk információt
- rengeteg standard adat
- de sose elég, 6000+ nyelv :(



- Python, Java



- Python, Java
 - NLTK (Natural Language Toolkit, Python)
 - Stanford CoreNLP, Apache OpenNLP (Java)



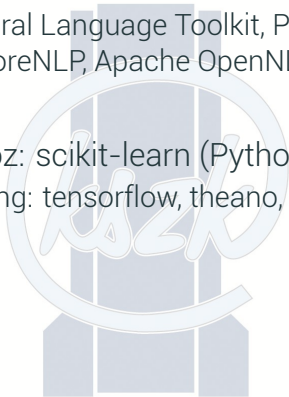
- Python, Java
 - NLTK (Natural Language Toolkit, Python)
 - Stanford CoreNLP, Apache OpenNLP (Java)
- Linux (OSX)



- Python, Java
 - NLTK (Natural Language Toolkit, Python)
 - Stanford CoreNLP, Apache OpenNLP (Java)
- Linux (OSX)
- gépi tanuláshoz: scikit-learn (Python), Weka (Java)



- Python, Java
 - NLTK (Natural Language Toolkit, Python)
 - Stanford CoreNLP, Apache OpenNLP (Java)
- Linux (OSX)
- gépi tanuláshoz: scikit-learn (Python), Weka (Java)
 - deep learning: tensorflow, theano, torch, keras stb.



- Python, Java
 - NLTK (Natural Language Toolkit, Python)
 - Stanford CoreNLP, Apache OpenNLP (Java)
- Linux (OSX)
- gépi tanuláshoz: scikit-learn (Python), Weka (Java)
 - deep learning: tensorflow, theano, torch, keras stb.
- plain text, TSV, esetleg XML

- Python, Java
 - NLTK (Natural Language Toolkit, Python)
 - Stanford CoreNLP, Apache OpenNLP (Java)
- Linux (OSX)
- gépi tanuláshoz: scikit-learn (Python), Weka (Java)
 - deep learning: tensorflow, theano, torch, keras stb.
- plain text, TSV, esetleg XML
 - Linux CL toolokkal is elvégezhető sok egyszerű feladat (awk, sed)

Hogyan készítenénk helyesírásellenőrzőt?



Schönherz

Hogyan készítenénk helyesírássellenőrzőt?

Hogyan sorolnánk fel az összes magyar szót?



Hogyan készítenénk helyesírásellenőrzőt?

Hogyan sorolnánk fel az összes magyar szót?

Demo



Köszönöm a figyelmet

judit@aut.bme.hu

Demo: [https://gist.github.com/juditacs/
4435129e6f79015ba98fba13f1736b84](https://gist.github.com/juditacs/4435129e6f79015ba98fba13f1736b84)

Önlab, szakdolgozat, diplomaterv, TDK lehetőség



Schönherz