

Entitás kinyerés magyar nyelvű szövegekből kétirányú LSTM-mel

Scheier Balázs – FAKK40, Mészáros Bálint – HY90XY

Budapesti Műszaki és Gazdaságtudományi Egyetem

Entitás felismerés bevezető

- NLP feladat

Entitás felismerés bevezető

- NLP feladat
- Célja, hogy felismerjük a szövegben szereplő előre meghatározott kategóriába tartozó entitásokat.

Entitás felismerés bevezető

- NLP feladat
- Célja, hogy felismerjük a szövegben szereplő előre meghatározott kategóriába tartozó entitásokat.
- Fontos alapfeladat a szövegfeldolgozásban.

Az entitásokról

A következő entitásokkal foglalkozunk:

- Helyszín (location)
- Személy (person)
- Szervezet (organization)
- Vegyes (miscellaneous)

Az entitásokról

A következő entitásokkal foglalkozunk:

- Helyszín (location)
- Személy (person)
- Szervezet (organization)
- Vegyes (miscellaneous)

Jelölések:

- B-LOC, I-LOC
- B-PER, I-PER
- B-ORG, I-ORG
- B-MISC, I-MISC
- 0

Az entitásokról

A következő entitásokkal foglalkozunk:

- Helyszín (location)
- Személy (person)
- Szervezet (organization)
- Vegyes (miscellaneous)

Jelölések:

- B-LOC, I-LOC
- B-PER, I-PER
- B-ORG, I-ORG
- B-MISC, I-MISC
- 0

Budapest egy szép város. – B-LOC

A budapesti kirándulás nagyon tetszett. – Nem entitás

- hunNERwiki

A	text	0	ART	a	0		
céljuk	text	0	NOUN<POSS<PLUR>>			cél	0
,	text	0	PUNCT	,	0		
hogy	text	0	CONJ	hogy	0		
biztosítsák	text	0	VERB<SUBJUNC-IMP><PLUR><DEF>			biztosít	0
,	text	0	PUNCT	,	0		
hogy	text	0	CONJ	hogy	0		
a	text	0	ART	a	0		
korábbi	text	0	ADJ	korai	0		
szervek	text	0	NOUN<PLUR>	szervek	0		
kilét	text	0	NOUN<POSS><CAS<ACC>>			kilét	0

Az adat struktúrája


```
'B-LOC' : 0,
'B-MISC' : 1,
'B-ORG' : 2,
'B-PER' : 3,
'I-LOC' : 4,
'I-MISC' : 5,
'I-ORG' : 6,
'I-PER' : 7,
'O' : 8,
'PAD' : 9,
'BOS' : 10,
'EOS' : 11
```

(a) A szótárunk

Model: "model"

Layer (type)	Output Shape	Param #
=====		
input_1 (InputLayer)	[(None, 28)]	0
embedding (Embedding)	(None, 28, 64)	1262464
bidirectional (Bidirectional)	(None, 28, 512)	657408
bidirectional_1 (Bidirectional)	(None, 512)	1574912
dense (Dense)	(None, 12)	6156
=====		
Total params: 3,500,940		
Trainable params: 3,500,940		
Non-trainable params: 0		

(b) A neurális háló architektúrája

- Adatok: tanító-teszt-validációs: 0.6-0.2-0.2

- Adatok: tanító-teszt-validációs: 0.6-0.2-0.2
- Optimalizáció accuracy-ra

- Adatok: tanító-teszt-validációs: 0.6-0.2-0.2
- Optimalizáció accuracy-ra
- Early stopping 5 epoch után

Hiperparaméter optimalizálás

- Rejtett rétegek mérete: 64, 128, 256, 512

Hiperparaméter optimalizálás

- Rejtett rétegek mérete: 64, 128, 256, 512
- Optimalizációs eljárások: RMSProp, Adam, SGD

Hiperparaméter optimalizálás

- Rejtett rétegek mérete: 64, 128, 256, 512
- Optimalizációs eljárások: RMSProp, Adam, SGD
- Batch méret: 32, 64, 128, 256

Hiperparaméter optimalizálás

Embed-ding réteg mérete	Első LSTM réteg mérete	Második LSTM réteg mérete	Optimalizációs algoritmus	Batch méret	Legjobb validációs accuracy
64	256	256	Adam	256	0.9905
512	256	512	RMSprop	256	0.9883
512	512	64	Adam	64	0.9881
512	512	512	RMSprop	256	0.9876
128	128	128	Adam	128	0.9875

Az 5 legjobb eredmény

A legjobb modellel (64, 256, 256, Adam, 256)

A legjobb modellel (64, 256, 256, Adam, 256)

- Tanító accuracy: 0.9973

A legjobb modellel (64, 256, 256, Adam, 256)

- Tanító accuracy: 0.9973
- Validációs accuracy: 0.9861

A legjobb modellel (64, 256, 256, Adam, 256)

- Tanító accuracy: 0.9973
- Validációs accuracy: 0.9861
- Teszt accuracy: 0.9874

- A tudomány szerint is működik az alvásmódszer, amit Salvador Dalí is használt.

- A tudomány szerint is működik az alvásmódszer, amit Salvador Dalí is használt.
 - Az eredeti címkék:
'BOS', 'O', 'O', 'O', 'O', 'O', 'O', 'O', 'O', 'O', 'B-PER', 'I-PER', 'O', 'O'

- A tudomány szerint is működik az alvásmódszer, amit Salvador Dalí is használt.
 - Az eredeti címkék:
'BOS', 'O', 'O', 'O', 'O', 'O', 'O', 'O', 'O', 'O', 'B-PER', 'I-PER', 'O', 'O'
 - A predikció:
'BOS', 'O', 'O', 'O', 'O', 'O', 'O', 'O', 'O', 'O', 'B-LOC', 'O', 'O', 'O'

- A tudomány szerint is működik az alvásmódszer, amit Salvador Dalí is használt.
 - Az eredeti címkék:
'BOS', 'O', 'O', 'O', 'O', 'O', 'O', 'O', 'O', 'O', 'B-PER', 'I-PER', 'O', 'O'
 - A predikció:
'BOS', 'O', 'O', 'O', 'O', 'O', 'O', 'O', 'O', 'O', 'B-LOC', 'O', 'O', 'O'
- Öt ok, amiért Macron Budapestre látogat.

- A tudomány szerint is működik az alvásmódszer, amit Salvador Dalí is használt.
 - Az eredeti címkék:
'BOS', 'O', 'O', 'O', 'O', 'O', 'O', 'O', 'O', 'O', 'B-PER', 'I-PER', 'O', 'O'
 - A predikció:
'BOS', 'O', 'O', 'O', 'O', 'O', 'O', 'O', 'O', 'O', 'B-LOC', 'O', 'O', 'O'
- Öt ok, amiért Macron Budapestre látogat.
 - Az eredeti címkék:
'BOS', 'O', 'O', 'O', 'B-PER', 'B-LOC', 'O'

- A tudomány szerint is működik az alvásmódszer, amit Salvador Dalí is használt.
 - Az eredeti címkék:
'BOS', 'O', 'O', 'O', 'O', 'O', 'O', 'O', 'O', 'O', 'B-PER', 'I-PER', 'O', 'O'
 - A predikció:
'BOS', 'O', 'O', 'O', 'O', 'O', 'O', 'O', 'O', 'O', 'B-LOC', 'O', 'O', 'O'
- Öt ok, amiért Macron Budapestre látogat.
 - Az eredeti címkék:
'BOS', 'O', 'O', 'O', 'B-PER', 'B-LOC', 'O'
 - A predikció:
'BOS', 'O', 'O', 'O', 'B-LOC', 'B-LOC', 'O'

- Téma: Entitás kinyerés magyar nyelvű szövegekből kétirányú LSTM-mel
- Csatatnév: Bokor, Mészáros, Scheier
- Résztvevők: Mészáros Bálint, Scheier Balázs