

Kell ide egy
cím

E. Bordi,
T. Both,
Sz. Pável,
Cs. Sándor,
A. Szász

Motiváció

Létező
megoldások

Létező
megoldások

Kutatási terv

Megtenni
szándékozott
lépések
Csapat tagjainak
hozzájárulása
Megvalósítandó
anyagok
Munkabeosztás
és ütemterv

Kell ide egy cím

A statisztikus gépi fordítás hatékonyságának jelentés egyértelműsítéssel történő javítása

E. Bordi, T. Both, Sz. Pável, Cs. Sándor, A. Szász

Babeş-Bolyai Tudományegyetem, Matematika és Informatika Kar, Kolozsvár

2016 április 15.

Tartalom

1 Motiváció

2 Létező megoldások

3 Létező megoldások

4 Kutatási terv

Megtenni szándékozott lépések
Csoport tagjainak hozzájárulása
Megvalósítandó anyagok
Munkabeosztás és ütemterv

Motiváció

Létező megoldások

Létező megoldások

Kutatási terv

Megtenni
szándékozott
lépések
Csoport tagjainak
hozzájárulása
Megvalósítandó
anyagok
Munkabeosztás
és ütemterv

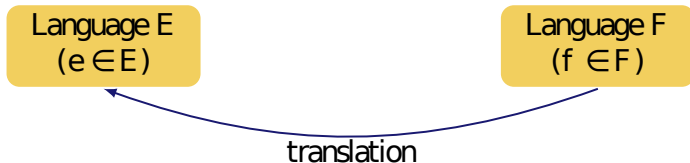
- Szövegek fordítása a számítógép által
- Statisztikus gépi fordítás (SMT)
- Filmfeliratok használata tanulási adatként
- Hatékonyságon javítani
- Jelentés egyértelműsítése (WSD)

Létező megoldások

Statistical machine translation

e: Good morning!

f: Bon jour!



$$P(e|f) = \frac{P(e)P(f|e)}{P(f)} \quad (1)$$

$$T(e) = \hat{e} = \operatorname{argmax}_e P(e|f) = \operatorname{argmax}_e P(e)P(f|e) \quad (2)$$

Létező megoldások

Statistical machine translation

Motiváció

Létező megoldások

Létező megoldások

Kutatási terv

Megtenni
szándékozott
lépések
Csapat tagjainak
hozzájárulása
Megvalósítandó
anyagok
Munkabeosztás
és ütemterv

$$T(e) = \hat{e} = \operatorname{argmax}_e P(e)P(f|e)$$

- nyelvi modell: $P(e)$
 - folytonosságot biztosít a célnyelvben
- fordítási modell: $P(f|e)$
 - lexikális megfeleltetés a nyelvek között
 - szó alapú (word based) modellek [?] [?]
 - kifejezés alapú (phrase based) modellek [?] [?]
- argmax : argmax_e
 - keresés

Létező megoldások

Előnyök, hátrányok

Szó alapú vs kifejezés alapú modellek:

- szó alapú modell: nehéz a tokeninzálás [?]
- kifejezés alapú modell pontosabb komplex nyelveknél [?]

Hátrányai ezeknek a rendszereknek:

- Még mindig nem elég pontosak az SMT rendszerek a WSD-hez képest [?]

Kell ide egy
cím

E. Bordi,
T. Both,
Sz. Pável,
Cs. Sándor,
A. Szász

Motiváció

Létező
megoldások

Létező
megoldások

Kutatási terv

Megtenni
szándékozott
lépések
Csapat tagjainak
hozzájárulása
Megvalósítandó
anyagok
Munkabeosztás
és ütemterv

Létező megoldások

Aktuális próbálkozások a javításra

- WSD beépítése az SMT-be: [?] [?]

Megtenni szándékozott lépések

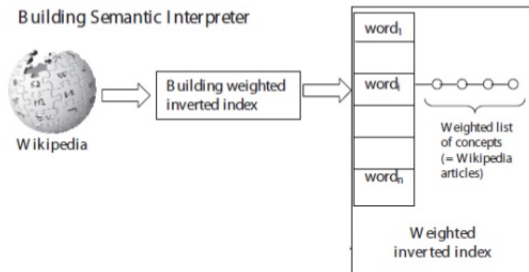
- Wikipedia alapú WSD rendszer felépítése
 - ① Egyik legnagyobb létező adathalmaz
 - ② Több nyelven elérhető
 - ③ A cikkek címei felhasználhatóak, mint conceptek
- WSD rendszer integrálása a létező SMT rendszerekbe
 - ① Statisztikai modellek nem veszik figyelembe többértelmű szavakat
 - ② A fordítás pontossága nagyban növelhető lenne ezen problémák megoldásával

Megtenni szándékozott lépések

Wikipedia alapú WSD rendszer

1 Invertált index felépítése minden nyelv számára

- Minden szóhoz hozzárendelünk concepteket és a hozzájuk tartozó súlyt
- A conceptek az angol Wikipedia címek
- Súlyozáshoz használhatjuk például a tf-idf súlyozást



Megtenni szándékozott lépések

Wikipedia alapú WSD rendszer

② Súlyozott concept vektor hozzárendelése a fordítandó és SMT által fordított szöveghez

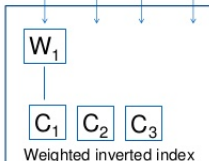
- Szöveg ábrázolása, mint szavak vektora
- Conceptek megfelelő súllyal történő hozzárendelése minden szóhoz
- Összesített súlyozott concept lista felépítése a szöveghez

Text fragment



Word vector

W_1 W_2 W_3 W_3



Concept vector

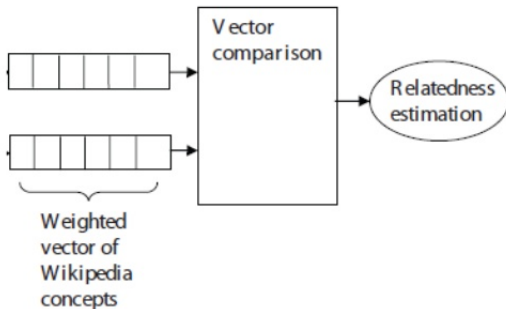
C_1 C_2 C_3 C_4

Megtenni szándékozott lépések

Wikipedia alapú WSD rendszer

③ Szemantikai hasonlóság meghatározza a fordítandó és célszöveg között

- A fordítandó és célszöveghez felépítjük a concept vektorokat
- A hasonlóság vizsgálatát a vektorok összehasonlítása jelenti
- A hasonlóság metrika lehet például a gyakran használt cos távolság



Megtenni szándékozott lépések

WSD súlyok integrálása

Motiváció

Létező
megoldások

Létező
megoldások

Kutatási terv

Megtenni
szándékozott
lépések

Csapat tagjainak
hozzájárulása

Megvalósítandó
anyagok

Munkabeosztás
és ütemterv

4 Legjobb N találat újrarangsorolása

- SMT rendszer alapján meghatározni a legjobb N fordítást
- WSD rendszer alapján súly hozzárendelése a találatokhoz
- A találatok újrarangsorolása az SMT és WSD együttes eredményei alapján

Kell ide egy
cím

E. Bordi,
T. Both,
Sz. Pável,
Cs. Sándor,
A. Szász

Motiváció

Létező
megoldások

Létező
megoldások

Kutatási terv

Megtenni
szándékozott
lépések

Csapat tagjainak
hozzájárulása

Megvalósítandó
anyagok

Munkabeosztás
és ütemterv

Csapat tagjainak hozzájárulása

Feladatkör	Csapattag
Megvalósítandó anyagok	Bordi Eszter
Munkabeosztás	
Gantt diagram	
Motiváció	Both Tibor
Megtenni szándékozott lépések	Pável Szabolcs
Létező megoldások	Sándor Csanád
Csapat tagjainak hozzájárulása	Szász Adorján

Kell ide egy
cím

E. Bordi,
T. Both,
Sz. Pável,
Cs. Sándor,
A. Szász

Motiváció

Létező
megoldások

Létező
megoldások

Kutatási terv

Megtenni
szándékozott
lépések
Csapat tagjainak
hozzájárulása

**Megvalósítandó
anyagok**

Munkabeosztás
és ütemterv

Megvalósítandó anyagok

E. Bordi, T. Both,
Sz. Pável,
Cs. Sándor,
A. Szász

Motiváció

Létező
megoldások

Létező
megoldások

Kutatási terv

Megtenni szándékozott
lépések

Csapat tagjainak
hozzájárulása

Megvalósítandó anyagok

Munkabeosztás és
ütemterv

Munkaegységek:

- A kutatás beindítása, adatgyűjtés
- Az adathalmaz összeállítása
- Algoritmus kidolgozása
- Az új és régi módszerek összevetése
- Az algoritmus optimalizálása, prototípus-fejlesztés

Munkaegységek és részfeladatai:

- A kutatás beindítása, adatgyűjtés
 - 1 Kutatási módszertani lehetőségek felvázolása
 - 2 Kutatási terv kidolgozása
 - 3 Meglévő módszerek feltérképezése
 - 4 Prototípus első koncepciójának kidolgozása
- Az adathalmaz összeállítása
 - 1 Adatok gyűjtése
 - 2 Standard formátumra való alakítás
 - 3 Adathalmaz validálása
- Algoritmus kidolgozása
 - 1 Algoritmus leírása
 - 2 Teszt jegyzőkönyvek, a statisztikus gépi fordítás hatékonyságával kapcsolatos hatások vizsgálata

Munkaegységek és részfeladatai:

- A kutatás beindítása, adatgyűjtés
 - 1 Kutatási módszertani lehetőségek felvázolása
 - 2 Kutatási terv kidolgozása
 - 3 Meglévő módszerek feltérképezése
 - 4 Prototípus első koncepciójának kidolgozása
- Az adathalmaz összeállítása
 - 1 Adatok gyűjtése
 - 2 Standard formátumra való alakítás
 - 3 Adathalmaz validálása
- Algoritmus kidolgozása
 - 1 Algoritmus leírása
 - 2 Teszt jegyzőkönyvek, a statisztikus gépi fordítás hatékonyságával kapcsolatos hatások vizsgálata

Munkaegységek és részfeladatai:

- A kutatás beindítása, adatgyűjtés
 - 1 Kutatási módszertani lehetőségek felvázolása
 - 2 Kutatási terv kidolgozása
 - 3 Meglévő módszerek feltérképezése
 - 4 Prototípus első koncepciójának kidolgozása
- Az adathalmaz összeállítása
 - 1 Adatok gyűjtése
 - 2 Standard formátumra való alakítás
 - 3 Adathalmaz validálása
- Algoritmus kidolgozása
 - 1 Algoritmus leírása
 - 2 Teszt jegyzőkönyvek, a statisztikus gépi fordítás hatékonyságával kapcsolatos hatások vizsgálata

E. Bordi, T. Both,
Sz. Pável,
Cs. Sándor,
A. Szász

Motiváció

Létező
megoldások

Létező
megoldások

Kutatási terv

Megtenni szándékozott
lépések

Csapat tagjainak
hozzájárulása

Megvalósítandó anyagok

Munkabeosztás és
ütemterv

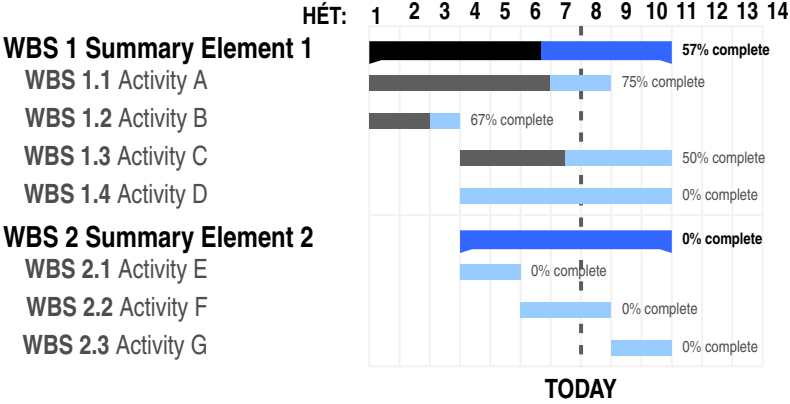
Munkaegységek és részfeladatai:

- Az új és régi módszerek összevetése
 - ① Tesztfuttatások különböző algoritmusokkal és adatokkal, a statisztikus gépi fordítás hatékonyságával kapcsolatos hatások vizsgálata
 - ② Az eltérés statisztikus validációja
 - ③ Szintézis
- Az algoritmus optimalizálása, prototípus-fejlesztés
 - ① Performanciaoptimalizált algoritmus
 - ② A gépi fordító prototípusának tesztelése és tesztjegyzőkönyvek
 - ③ Nyilvánosságra hozandó eredmények dokumentációja
 - ④ A projekt lezárása, dokumentáció

Munkaegységek és részfeladatai:

- Az új és régi módszerek összevetése
 - 1 Tesztfuttatások különböző algoritmusokkal és adatokkal, a statisztikus gépi fordítás hatékonyságával kapcsolatos hatások vizsgálata
 - 2 Az eltérés statisztikus validációja
 - 3 Szintézis
- Az algoritmus optimalizálása, prototípus-fejlesztés
 - 1 Performanciaoptimalizált algoritmus
 - 2 A gépi fordító prototípusának tesztelése és tesztjegyzőkönyvek
 - 3 Nyilvánosságra hozandó eredmények dokumentációja
 - 4 A projekt lezárása, dokumentáció

Munkabeosztás és ütemterv



Irodalomjegyzék I



A. Author.

Handbook of Everything.

Some Press, 1990.



S. Someone.

On this and that.

Journal of This and That, 2(1):50–100, 2000.