

Statistilise masintõlge väljundi parandamine

Mark Fišel (fishel@ut.ee)

01. märts 2007



Ettekanne sisu

- statistiline masintõlge (SMT);
- SMT seis Eesti keele puhul;
- olemasolevad parandamise meetodid;
- alternatiivne parandamise meetod



Statistiline masintõlge

$e, f, p(e), p(e|f)$

$$\begin{aligned}\hat{e} = \arg \max_e p(e|f) &= \arg \max_e \frac{p(e)p(f|e)}{p(f)} \\ &= \arg \max_e p(e)p(f|e)\end{aligned}$$

- $p(f|e)$ – tõlkimismudel
- $p(e)$ – keelemudel
- $\arg \max$ – dekodeerija



Fraasipõhine TM

Fraasipõhises tõlkimismudelis modelleeritakse erinevad etapid sõltumatult:

- fraaside tõlkimine
- fraaside järjekord

Probleemid:

- enne nägematud sõnad
 - valesti valitud sõnad/fraasid
 - vale fraaside järjekord
-



SMT + Eesti keel

- Eesti-Inglise suund
- testitud 2 korpuse peal, mõlemad – seadustekstidega
- parim tulemus:
 - BLEU: 45.22%
 - käsitsi hinnang:

Hea	15%
Rahuldav	15%
Vale	70%



Eesti-Inglise SMT probleemid

- hõredad andmed (sparse data)
 - mitu korda sisendsõna jääb tõlkimatuks – see ei esinenud treenimiskorpuses
 - seda põhjustab väike korpuste suurus, eesti keele rikas morfoloogia, liitsõnad, jne
- vale järjekord
 - järjekorra komponent ei saa hakkama



Vale järjekord

—(euroopa majandusühenduse ja šveitsi
konföderatsiooni)₁ (vaheliste kokkulepete)₂
(kohaldamisel)₃ (rakendatakse ühenduses)₄
(ühiskomitee)₅ (otsust nr 5 / 81)₆

—(for the purposes of application)₃ (of the agreements
between)₂ (the european economic community and the
swiss confederation)₁ , (decision no 5 / 81)₆ (of the
joint committee)₅ (shall apply in the community)₄

—(the european economic community and the swiss
confederation)₁ (of agreements between)₂ (the
application)₃ (of the joint committee)₅ (shall apply in
the community)₄ (decision no 5 / 81)₆



Olemasolev parandamise meetod

- kasutades teadmisi kahe keele grammatikat saab enne treenimist ja tõlkimist muuta sisend-keele järjekorda, et see oleks rohkem väljund-keelega sarnane, või
- järjekorra muutmise mustreid saab leida käsitsi või masinõppimisega
 - käsitsi tegemine nõuab väga palju inim- ja ajaresursse
 - masinõppe nõuab süntaksanalüsaatorit mõlema keele jaoks



Alternatiivne parandamise meetod

- oletame et me saame tõlkida sõnade asemel sõnaliike, ja seda väga hea täpsusega
 - S V A \Rightarrow NN VBZ JJ
- siis saab tõlkida enne lause sõnaliigid, siis lauset ennast, ja siis muuta lause fraaside järjekorda vastavalt sõnaliikidele



kommittee otsus nr 5 ->

of the committee decision nr 5

IN DT NN NN NN CD

S S S N -> NN NN CD IN DT NN

of the committee decision nr 5 ->

decision nr 5 of the committee



Sõnaliikide tõlkimine

- peaks olema lihtsam, kuna sõnavara on piiratud (eesti – 16, inglise Penn – umb. 50), mingit morfoloogiat või liitsõnu pole
- praktikas aga tavaline SMT töötab sõnaliikidega isegi halvemini kui sõnadega
 - üks teatud probleemide hulgast: peaaegu igas lauses on nii tegu- kui nimisõna, kuidas teha vahet?



Järjekorra muutmine

- brute force: $O(n!)$, kus n on sõnade (fraaside) arv
- mis võiks olla efektiivsem algoritm?
- probleemi teeb lihtsamaks see et väljund on juba jagatud fraasideks
- aga:
 - sõnaliigid võivad korduda, eriti nimisõnafraasides (lepingu allkirjastamise kuupäev – S S S)
 - tõlgitud variant ei ole ideaalne, ei pruugi sisaldada kõike “oraakli” sõnaliike



Kokkuvõte

- SMT
- Eesti SMT
- vale järjekord
- parandamise meetodid



Aitäh!

