

T-61.5020 Luonnollisten kielten tilastollinen käsittely

Vastaukset 9, to 28.3.2007, 12:15–14:00 — Tilastollinen konekääntäminen

Versio 1.0

1. Etsitään todennäköisin käänös \hat{e} ruotsinkieliselle lauseelle r :

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

Käytetään kirjassa esitettyä mallia todennäköisyydelle $P(r|e)$:

$$P(r|e) = \frac{1}{Z} \sum_{a_1=0}^l \cdots \sum_{a_m=0}^l \prod_{j=1}^m P(r_j|e_{a_j})$$

missä m on alkuperäisen ruotsinkielisen lauseen pituus ja l on käännetyn englanninkielisen lauseen pituus. Lasketaan kummallekin vaihtoehdolle:

$$P(r|e_1) = 1.0 \cdot 0.7 \cdot 0.9 \cdot 1.0 \cdot 1.0 \cdot 1.0 \cdot 0.1 = 0.063$$

$$P(r|e_2) = 1.0 \cdot 0.7 \cdot 1.0 \cdot 1.0 \cdot 1.0 \cdot 1.0 \cdot 1.0 \cdot 1.0 = 0.7$$

Tässä siis kokeillaan kaikkia käänössääntöjä jokaiselle ruotsinkielisen lauseen sanalle (huomioimatta sanajärjestystä). Koska säännöstö on hyvin harva, päästään näin yksinkertaiseen laskutoimitukseen.

Prioritodennäköisyys $P(e)$ saadaan kielimallista. Lasketaan se kummallekin lauseelle:

$$P(e_1) = \prod_{i=1}^l P(w_i) = 0.18 \cdot 0.05 \cdot 0.01 \cdot 0.13 \cdot 0.1 \cdot 0.12 \cdot 0.02 = 2.8 \cdot 10^{-9}$$

$$P(e_2) = 0.18 \cdot 0.07 \cdot 0.11 \cdot 0.21 \cdot 0.01 \cdot 0.13 \cdot 0.1 \cdot 0.01 = 3.8 \cdot 10^{-10}$$

Kertomalla priori ja käänöstodennäköisyys huomataan, että jälkimmäinen käänös on todennäköisempi:

$$P(e_1)P(r|e_1) = 0.063 \cdot 2.8 \cdot 10^{-9} = 1.8 \cdot 10^{-10}$$

$$P(e_2)P(r|e_2) = 2.6 \cdot 10^{-10}$$

Huomattavaa on, että käänösmalli ei ota mitään kantaa sanajärjestykseen. Koska myöskään kielimalli (unigrammimalli) ei tätä tee, mallin mielestä sanajärjestyksellä ei ole mitään merkitystä. Jos mallilta kysytään todennäköisintä lausetta (ei testata eri vaihtoehtoja) huomataan, että todennäköisimpään lauseeseen ei voi tulla artikkeleja eikä sanaa “into”. Tämä johtuu siitä, että niiden lisääminen ei muuta käänöstodennäköisyyttä, mutta tiputtaa aina kielimallitodennäköisyyttä. Kielimalli siis suosii muutenkin lyhyempiä lauseita. Kasvattamalla kielimallin konteksti trigrammiksi, saisi ehkä artikkelit ja sanajärjestyksen paremmin kohdalleen.

Yleisesti tällä menetelmällä tarvitaan heuristiikkaa valitsemaan käänökset, joita tutkitaan. Kaikkien vaihtoehtojen läpikäynti on käytännössä mahdotonta.

2. Käytetään esimerkkinä sanaa $f = \text{“tosiasia”}$. Se on esiintynyt korpuksessa 989 kertaa. Normalisointia varten aineistosta pitää laskea erikseen kaikkien englannin sanojen esiintymismäärät.

a-b) Allaolevassa taulukossa on 20 suurinta arvoa saanutta englannin sanaa yhteisesiintymien frekvenssin sekä englannin sanan esiintymismäärällä normalisoidun frekvenssin mukaisesti. Nähdään että kumpikaan menetelmä ei anna toivottuja tuloksia. Ensimmäisessä ovat ongelmana hyvin yleiset sanat, joita esiintyy lähes jokaisessa lauseessa ja siten myös yhtäaikaa f :n kanssa. Jälkimmäisessä taas ongelmaksi muodostuvat hyvin harvinaiset sanat: Jos kerran esiintynyt sana sattuu esiintymään yhtä aikaa f :n kanssa, se saa suurimman mahdollisen arvon.

e	$C(e, f)$	e	$\frac{C(e, f)}{C(e)}$
the	2563	winkler	1.0000
of	1128	visarequired	1.0000
that	1086	visaexempt	1.0000
is	1040	veiling	1.0000
and	829	valuejudgment	1.0000
to	823	undisputable	1.0000
in	726	stayers	1.0000
a	716	semipermeable	1.0000
fact	654	rulingout	1.0000
it	385	roentgen	1.0000
this	376	residuarity	1.0000
we	316	regionallevel	1.0000
are	295	redhaired	1.0000
not	280	poorlyfounded	1.0000
for	274	philippic	1.0000
have	253	pemelin	1.0000
be	234	paiania	1.0000
which	230	overcultivation	1.0000
on	224	outturns	1.0000
has	212	onesixth	1.0000

c) Edellisissä menetelmissä oli ongelmana että ne eivät ottaneet huomioon käännöksen molempia suuntia: Jotta e olisi todennäköinen käännös f :lle, e :n pitäisi esiintyä niissä lauseissa joissa f esiintyi, sekä f :n pitäisi esiintyä niissä lauseissa joissa e esiintyi. Tällöin todennäköisyyksien $P(e|f) = \frac{C(e, f)}{C(f)}$ sekä $P(f|e) = \frac{C(e, f)}{C(e)}$ pitäisi kummankin olla suuria. Kokeillaan seuraavana painotuksena näiden todennäköisyyksien tuloa.

Tulokset ovat seuraavan sivun vasemmanpuoleisessa taulukossa. Tällä kertaa löydettiin oikea käännös, ja myös toinen merkitykseltään läheinen sana “reality” on

suhteellisen korkealla.

Kokeillaan vielä kollokaatiolaskarista tuttua χ^2 -testiä:

$$\chi^2 = \frac{N(O_{11}O_{22} - O_{12}O_{21})^2}{(O_{11} + O_{12})(O_{11} + O_{21})(O_{12} + O_{22})(O_{21} + O_{22})},$$

missä

$$\begin{aligned} O_{11} &= C(e, f) \\ O_{12} &= C(e, \neg f) = C(e) - C(e, f) \\ O_{21} &= C(\neg e, f) = C(f) - C(e, f) \\ O_{22} &= C(\neg e, \neg f) = N - C(e) - C(f) + C(e, f) \end{aligned}$$

ja N aineiston kaikkien lauseiden määrä. Yli 3.843 arvon saaville sanoille tulokset tarkoittavat, että todennäköisyys sille että yhteisesiintymät olivat sattuman aikaansaamia, on alle 5%.

Suurimmat arvot saaneet sanat ovat oikeanpuoleisessa taulukossa. Testi näyttää toimivan erittäin hyvin: Ainoastaan “fact” ylitti luotettavuusrajan. Toisaalta jos halusimme vaihtoehtoisia käännöksiä, kuten “reality”, joku todennäköisyysarvon antava menetelmä olisi käyttökelpoisempi. Käytännössä käännöstodennäköisyydet etsitään iteratiivisesti EM-algoritmeilla, jolloin rajoitetaan sitä että yksi englannin sana olisi monen suomenkielisen sanan käänнос, mutta esitetyn kaltainen menetelmä voisi toimia todennäköisyyksien alustuksena.

e	$\log\left(\frac{C(e,f)}{C(e)} \cdot \frac{C(e,f)}{C(f)}\right)$	e	χ^2
fact	-3.9758	fact	18.2155
the	-5.6159	the	3.5937
that	-5.8849	that	3.0070
is	-5.9086	is	2.8096
reality	-6.0057	of	2.3485
winkler	-6.2035	reality	2.3166
of	-6.4963	winkler	2.0000
hedgehog	-6.6090	hedgehog	1.3323
a	-6.6577	visarequired	1.0000
and	-6.8386	visaexempt	1.0000
visarequired	-6.8967	veiling	1.0000
visaexempt	-6.8967	valuejudgment	1.0000
veiling	-6.8967	undisputable	1.0000
valuejudgment	-6.8967	stayers	1.0000
undisputable	-6.8967	semipermeable	1.0000
stayers	-6.8967	rulingout	1.0000
semipermeable	-6.8967	roentgen	1.0000
rulingout	-6.8967	residuarity	1.0000
roentgen	-6.8967	regionallevel	1.0000
residuarity	-6.8967	redhaired	1.0000