

T-61.281 Luonnollisten kielten tilastollinen käsittely

Harjoitus 11, ti 6.4.2004, 8:30-10:00 Klusterointi, Konekääntäminen. Versio 1.0

- 2-ulotteiseen avaruuteen on projisoitu sanoja. Sanat, jotka esiintyvät samankaltaisissa lauseissa ovat lähekkäin tässä avaruudessa. Piirrä dendrogrammit taulukossa 1 olevien sanojen klusteroinnin muodostumiselle käyttäen
 - yksinkertaisen linkin (single link) klusterointia
 - kokonaisen linkin (complete link) klusterointia

| koordinaatit | sana |
|--------------|----------|
| (-4,2) | puukko |
| (-2,2) | tuppi |
| (-2,-1) | kaivuri |
| (-3,-2) | zetor |
| (1,-3) | kenraali |
| (2,2) | maija |
| (2.5,1) | matti |
| (4,2) | jens |

Taulukko 1: *klusteroitavat sanat*

- Etsi k-means -klusteroinnilla edellisen tehtävän sanoista 3 klusteria. Oletetaan, että klusterien lähtöarvoiksi on arvottu $(2,3)$, $(2,0.5)$ ja $(4,1)$.
- Taulukossa 2 on kaksiulotteiseen avaruuteen projisoituja sanoja.
 - Käytetään knn-luokitinta 3:lla naapurilla. Arvioi sanoille, joille ei ole merkitty sanaluokkaa, paras sanaluokka.
 - Korvaa samaan sanaluokkaan kuuluvien sanojen vektorit niiden keskiarvolla. Vertaa uusia sanoja kahteen muodostuneeseen prototyyppi-vektoriin. Minkälainen luokittelutulos nyt tulee ?
 - Mitä etua ja haittoja b)-kohdan menetelmällä on verrattuna a)-kohdan menetelmään ?
- Olet klusteroinut samankaltaiset sanat ryhmiin. Perinteisessä trigrammimallissa arvioidaan seuraavan sanan todennäköisyys perustuen edellisiin sanoihin ($P(w_n|w_{n-1}, w_{n-2})$). Mallin kokoa rajoittaaksesi haluat nyt kuitenkin arvioida seuraavan sanan todennäköisyyden niin, että käytät historiana vain edellisten sanojen klustereita, etkä sanoja sinällään. Johda tällaisen kieliopin matemaattinen muoto.

| sana | luokka | koordinaatit |
|--------------|--------------|--------------|
| vetää | verbi | (2,5,3) |
| työntää | verbi | (1,2) |
| nostaa | verbi | (3,2) |
| moukari | substantiivi | (1,1) |
| naama | substantiivi | (-1,2) |
| hius | substantiivi | (-5,1) |
| heittää | ? | (2,7,2,7) |
| kihartuminen | ? | (-3,2) |
| kuula | ? | (0,5,2) |

Taulukko 2: kn , data

| w | $P(w)$ | w_1 | w_2 | $P(w_1 \rightarrow w_2)$ |
|----------|--------|----------|--------|--------------------------|
| it | 0.18 | it | den | 1.0 |
| becomes | 0.05 | becomes | blir | 0.7 |
| clean | 0.01 | becomes | klär | 0.3 |
| eats | 0.1 | turns | blir | 0.7 |
| the | 0.12 | turns | vänder | 0.3 |
| seats | 0.02 | into | ∅ | 1.0 |
| turns | 0.07 | clean | ren | 0.9 |
| into | 0.11 | clean | städa | 0.1 |
| a | 0.21 | a | ∅ | 1.0 |
| reindeer | 0.01 | reindeer | ren | 1.0 |
| and | 0.13 | and | och | 1.0 |
| lichen | 0.01 | eats | äter | 1.0 |
| | | the | ∅ | 1.0 |
| | | seats | laven | 0.1 |
| | | seats | stolar | 0.9 |
| | | lichen | laven | 1.0 |

Taulukko 3: Vasemmalla unigrammikielimalli, oikealla käännöstodennäköisyydet.

5. Etsit ratkaisua hevostmiehiä pitkään pohdituttaneeseen ongelmaan, “Varför får hästen inte gå i bastun?”. Ratkaisun ongelmaan tuntevat vain ruotsalaiset (“Den blir ren och äter laven”). Osaat englantia ja käytössäsi on sekä taulukon 3 kielimalli ja käännöstiedot. Sinulla on kaksi vahvaa ehdokasta vastauksen käännökseksi:

- It becomes clean and eats the seats
- It turns into a reindeer and eats lichen

Kumpi on todennäköisempi ?