

Bengio, Ducharme, Vincent: A Neural Probabilistic Language Model

Arto Teräs (ajt@iki.fi)

31.10.2001

Tavoite

- Laskea aiemman materiaalin perusteella seuraavan sanan todennäköisyys

Päätideä

- Käytetään neuroverkkoa, jolle opetetaan yhtä aikaa sekä sanasekvenssit että kunkin sanan ominaisuudet
- ⇒ Pystytään arvioimaan todennäköisyyttä myös silloin, kun sanasekvenssiä ei ole aiemmin esiintynyt

Arto Teräs (ajt@iki.fi) 31.10.2001

2

Yleinen tilastollinen kielimalli

- Ilmaistaan sanan todennäköisyys kaikkien aiemmin havaittujen sanojen perusteella
- $P(w_t^T) = \prod_{t=1}^T P(w_t | w_1^{t-1})$, jossa w_t on t :s sana ja w_t^j sanasekvenssi (w_1, w_2, \dots, w_j)

Ongelmat

- Sanojen yhteisjakauma (joint distribution) kasvaa eksponentiaalisesti huomiottaessa useampia peräkkäisiä sanoja
 - Kompleksisuus $O(V^n)$, jossa V sanaston koko ja n peräkkäisten sanojen määrä
- Esimerkiksi 100000 sanan sanasto ja 10 peräkkäistä sanaa: enimmillään $100000^{10} - 1 = 10^{50} - 1$ vapaita parametria
- Mitä useampia sanoja huomioidaan, sitä harvemmin aiemmin on havaittu täsmälleen sama sekvenssi

Perinteiset ratkaisut

- N-gram -mallit: Tilastoidaan n peräkkäisen sanan yhdistelmät, jossa n on pieni
- Lasketaan todennäköisyys perustuen alemmin kohdatettujen sanojen yhdistelmien lukumäärään
- Kohdattaessa outo sanasekvenssi muodostetaan todennäköisyys lyhyemmän kontekstin perusteella

Arto Teräs (ajf@iki.fi) 31.10.2001

5

Sanojen piirvektorit

- Liitetään jokaiseen sanaan piirvektori $C(w_n)$, joka kertoo sen ominaisuudet
- Ominaisuuksien lukumäärä huomattavasti pienempi kuin sanaston koko, esimerkiksi 30–100 kpl
- Samankaltaisilla sanoilla lähekkäiset piirvektorit
- Vektorit voidaan alustaa satunnaisiksi ja antaa neuroverkon hakea arvot oppimisprosessin aikana

Arto Teräs (ajf@iki.fi) 31.10.2001

7

Samankaltaiset sanat

- Samankaltaiset sanat esiintyvät lauseissa samoissa rooleissa
- ⇒ Lauseen **The cat is walking in the bedroom** esiintyminen opetusmateriaalissa pitäisi lisätä lauseen **A dog was running in a room** todennäköisyyttä, vaikka jälkimmäistä ei olisikaan alemmin havaittu

Arto Teräs (ajf@iki.fi) 31.10.2001

6

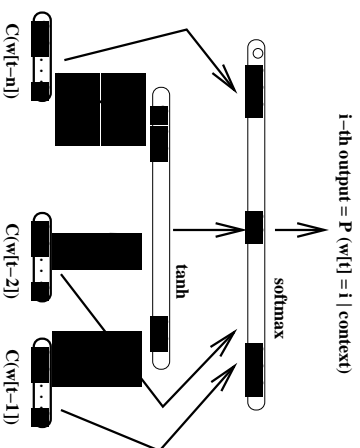
Opetusvaihe

- Tavoitteena oppia hyvä malli $f(w_t, \dots, w_{t-n}) = \hat{P}(w_t | w_1^{t-1})$
- Annetaan syötteeksi kerrallaan n sanan piirvektoriit
- Määritellään sanasekvenssien todennäköisyysfunktio, jonka muuttujina ovat nämä piirvektorit
- Verkko sovitaa samankaltaisesti vektorien arvoja ja todennäköisyysfunktioita
- Kaksi arkkitehtuuria, jotka poikkeavat toisistaan hieman sekvenssin viimeisimmän sanan käsitteilyn osalta

Arto Teräs (ajf@iki.fi) 31.10.2001

8

Mallin kaaviokuva (suora arkkitehtuuri)



Arto Teräs (ajt@iki.fi) 31.10.2001

9

Bengio, Ducharme, Vincent: A Neural Probabilistic Language Model

Tuloksia

- Mallia testattiin noin miljoonan sanan (Brown) ja 34 miljoonan sanan (Hansard) aineistoilla
- Noin 20-30 % parempia tuloksia smoothed trigram-malliin verrattuna perplexity-arvoilla mitattuna
- Parhaat tulokset saavutettiin yhdistämällä malli trigram-mallin kanssa

Arto Teräs (ajt@iki.fi) 31.10.2001

11

Toiminnan pääperiaate

- Kun todennäköisyysfunktioita (log-likelihood) maksimoidaan, samankaltaisissa yhteyksissä esiintyvät sanat saavat lähellä toisiaan olevat piirrevektorit
- Piirteiden ollessa lähellä toisiaan uusi havaittu lause X ei lisää pelkästään saman sanasekvenssin todennäköisyyttä vaan myös kaikkien niiden, jotka ovat lähellä X:ää lauseiden piirreavaruudessa
- Esimerkki: Kissa käveli katolla
Koira käveli kadulla (++)
Ulkona oli jänis (+)
Kauris silli leijona (-)

Arto Teräs (ajt@iki.fi) 31.10.2001

10

Bengio, Ducharme, Vincent: A Neural Probabilistic Language Model

Havaintoja

- Varottava ylisovittusta (over-fitting) pienillä opetusaineistoilla
- Laskentatarckuus tärkeää suurilla aineistoilla
- Laajennan kontekstin huomiointi parantaa tuloksia noin 10 sanaan asti
- Piirrevektorien alustuksella ei suurta merkitystä vertailtaessa täysin satunnaista alustusta ja erään sanojen yleisyyteen perustuvan algoritmin antamia arvoja
- Piirrevektorien arvojen samanaikainen sovitus yhtiä aikaa muiden parametrien kanssa tärkeää

Arto Teräs (ajt@iki.fi) 31.10.2001

12

Nopentuskeinoja

- Listat kaikkein todennäköisimmistä sanoista
- Useimmin esiintyvien sekvenssien esilaskenta
- Puheentunnistuksessa riittää, että tutkitaan akustisesti samankaltaisia sanoja
- Tutkitavan tekstin kappaleiden permutointi voi nopeuttaa sovitusta

Arto Teräs (ajf@iki.fi) 31.10.2001

13

Kotitehtäväkysymykset

1. Esitetyssä mallissa huomioidaan tavallista suurempi määrä edeltäviä sanoja uuden sanan todennäköisyyttä laskettaessa. Mitä teknikoita käytetään, jotta vapaiden parametrien määrä ja laskenta-aika pysyvät kohtuullisina?
2. Miten sanoihin liittyvät piirvektorit (feature vectors) muodostetaan? Miten vektorien valintaa voitaisiin mahdollisesti parantaa?

Arto Teräs (ajf@iki.fi) 31.10.2001

15

Avoimia kysymyksiä

- Ominaisuusvektorien määrittely eri kokoisille yksiköille (parin sanan ryhmät, sanojen osat)
- Eriytyypisten neuroverkkojen ja oppimisalgoritmien vaikutus
- Mitä voidaan päätellä tutkimalla millaiset piirvektorit ovat muodostuneet oppimisprosessin seurauksena
- Muualla tavoin hankitun semanttisen, syntaktisen ja morfologisen informaation hyödyntäminen
- Suorituskyky käytännön sovelluksissa

Arto Teräs (ajf@iki.fi) 31.10.2001

14