

Normalisoitu pakkausetäisyys: sovelluksia ja variaatioita

Timo Sand

Kandidaatintutkielma
HELSINGIN YLIOPISTO
Tietojenkäsittelytieteen laitos

Helsinki, 20. lokakuuta 2013

1 Johdanto

Kaikki data on luotu samanveroiseksi, mutta jotkut ovat samankaltaisempia kuin toiset. Esitämme tavan jolla esittää tämä samankaltaisuus, käyttäen uutta samankaltaisuuden metriikkaa (*engl. similarity metric*), joka perustuu tiedoston pakkaamiseen. Metriikka on parametriton, eli se ei käytä datan ominaisuuksia tai taustatietoja, ja sitä voi soveltaa eri aloihin ilman muunnoksia. Metriikka on universaali siten, että se approksimoi parametrin, joka kaikissa pareittain vertailuissa ilmaisee samankaltaisuutta hallitsevassa piirteessä. Se on vakaa siinä mielessä, että sen tulokset ovat riippumattomia käytetystä pakkaajasta [CV05]. Pakkaajalla tarkoitetaan pakkausohjelmaa kuten *gzip*, *ppmz*, *bzip2*.

Pakkaukseen perustuva samankaltaisuus (*engl. Compression-Based Similarity*) on “universaali” metriikka, jonka kehittivät Cilibrasi ja Vitanyi [CV05]. Yksinkertaistettuna tämä tarkoittaa, että kaksi objektia ovat lähellä toisiaan, jos voimme “pakata” yhden objektin huomattavasti tiiviimmin toisen objektin datalla. Abstraktina ideana toimii se, että voimme kuvailla ytimekkäämmin yhden palan toisen avulla, mikäli palat ovat samankaltaisia. Tämän esittelemme luvussa 2 ja samalla käymme läpi mihin teoriaan algoritmi perustuu sekä miten se toimii. Edellä mainitun vakauden esittämiseen voimme käyttää useaa tosielämän pakkausalgoritmiä: tilastollista (PPMZ), Lempel-Ziv -algoritmiin pohjautuvaa hakemistoa (*gzip*), lohkopерusteista (*bzip2*) tai erityistä (Gencompress).

Tarkoituksemme on koota yksittäiseen samankaltaisuuden metriikkaan kaikki todelliset etäisyydet; tehokkaat versiot Hammingin etäisyydestä, Euklidisestä etäisyydestä, Lempel-Ziv etäisyydestä ja niin edelleen. Tämän metriikan pitäisi olla niin yleinen, että se toimii, yhtäläisesti ja samanaikaisesti, kaikille aloille: musiikki, teksti, kirjallisuus, ohjelmat, genomit, luonnollisen kielen määrittelyt. Sen pitäisi pystyä samankaltaisesti havaitsemaan kaikki samankaltaisuudet, joita muut etäisyydet havaitsevat erikseen, palojen välillä.

Kun määrittelemme ryhmän sallittavia etäisyyksiä (*engl. admissible distances*) haluamme sulkea pois epärealistiset, kuten $f(x, y) = \frac{1}{2}$ jokaiselle parille $x \neq y$. Saavutamme tämän rajoittamalla objektien lukumäärän an-

netussa etäisyydessä objektiin. Teemme tämän huomioimalla vain todellisia etäisyyksiä seuraavasti: Määrräämme sopivan ja tietyn ohjelmointikielen, joka toimii tutkielman ajan referenssikielenä. [CV05]

Luvussa 3 esittelemme algoritmin käyttökohteita monelta eri alueelta. Aloitamme siitä, miten yleisesti NCD:n avulla pystymme klusteroimaan tuloksia eri kategorioihin; miten musiikkikappaleet klusteroituvat saman artistin alle, miten kuvantunnistuksessa saamme ryhmitettyä samankaltaiset kuvat ja miten sienten genomeista saamme tarkan lajiryhmityksen.

Syvennymme musiikin, kuvantunnistuksen ja dokumenttien kategorisoinnin tuloksiin luvun lopussa.

Luvussa 4 esitellään NCD:n kestävyyttä ja ongelmia. Ensiksi esitellään NCD:n kohinansietokykyä, eli katsotaan mitä tapahtuu kun lisätään vähitellen kohinaa toiseen tiedostoista, jota pakataan, ja mittaamalla samankaltaisuutta tämän jälkeen [CAO07]. Saamme nähdä miten paljon kohina vaikuttaa NCD:n laskemiin etäisyyksiin ja huonontaako se klusteroinnin tuloksia.

Mikään algoritmi ei ole täydellinen ja niin NCD-algoritmilläkin on ongelmansa. Algoritmissä itsessään ei ole selvää heikkoutta, mutta sen käytössä on otettava pakkaajan valinta huomioon, koska monet suosituista pakkausalgoritmeista ovat optimoituja tietyn kokoisille tiedostoille. Niissä on niin kutsuttu ikkunakoko (*engl. window size*), joka määrittelee mikä tiedostokoko on sopiva [CAO05]. Jos tiedostokoko on pienempi kuin ikkunakoko, niin pakkaus on tehokasta, kun mennään siitä yli, niin pakkauksesta tulee huomattavasti tehottomampaa. Esittelemme tuloksia eri pakkausalgoritmien vertailuista ja mikä näistä algoritmeista on parhaimmaksi havaittu NCD:n kanssa käytettäväksi.

NCD ei ole ainut metriikka, jolla voidaan mitata samankaltaisuutta. Internetiä hyödyntäen on tehty metriikka, joka käyttää hakukoneita samankaltaisuuden tutkimiseen; tämä on nimetty Google samankaltaisuusetäisyydeksi (*engl. Google Similarity Distance*). Tämä toimii myös muilla hakukoneilla kuten Bing. Luvussa 5 esitellemme tämän sekä muita samankaltaisuuden metriikoita.

Lähteet

- [CAO05] Cebrian, Manuel, Alfonseca, Manuel ja Ortega, Alfonso: *Common pitfalls using the normalized compression distance: What to watch out for in a compressor*. Communications in Information & Systems, 5(4):367–384, 2005.
- [CAO07] Cebrian, M., Alfonseca, M. ja Ortega, A.: *The Normalized Compression Distance Is Resistant to Noise*. Information Theory, IEEE Transactions on, 53(5):1895–1900, 2007, ISSN 0018-9448.
- [CV05] Cilibrasi, Rudi ja Vitanyi, Paul M. B.: *Clustering by Compression*. IEEE Transactions on Information Theory, 51(4):1523–1545, Huhtikuu 2005.