

# **Normalisoitu pakkausetäisyys**

Timo Sand

Kandidaatintutkielma  
HELSINGIN YLIOPISTO  
Tietojenkäsittelytieteen laitos

Helsinki, 13. lokakuuta 2013

## Sisältö

<b>1</b>	<b>Johdanto</b>	<b>1</b>
<b>2</b>	<b>Normalisoitu Pakkausetäisyys</b>	<b>1</b>
<b>3</b>	<b>Klusterointi</b>	<b>3</b>
<b>4</b>	<b>Käyttökohteet</b>	<b>3</b>
4.1	Klusteroinnin tulokset . . . . .	3
4.2	Kuvantunnistus . . . . .	3
<b>5</b>	<b>Kohinansietokyky</b>	<b>3</b>
<b>6</b>	<b>Huomioitavaa kompressorin valitsemisessa</b>	<b>4</b>
<b>7</b>	<b>Muita samankaltaisuuden metriikoita</b>	<b>4</b>
7.1	Google Similarity Distance . . . . .	4
	<b>Lähteet</b>	<b>4</b>

# 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. Se on parametrin siten että se ei käytä ominaisuuksia tai taustatieto datasta, ja sitä voi soveltaa eri alueisiin ja aluerajojen yli ilman muunnoksia. Se on yleinen siten että se approksimoi parametrin joka esittää samankaltaisuutta hallitsevassa piirteessä kaikissa pareittain vertailuissa. Se on tukeva siinä mielessä, että sen menestys esiintyy riippumatta siitä mitä kompressoria käytetään.

Seuraavaksi esittelemme algoritmin toimintaperiaatteen ja sen koostumuksen.

Luvussa 3 esittelemme algoritmin käyttökohteita monelta eri alueelta

Luvussa 4 esittelemme algoritmin ominaisuuksia ja ongelmia (Kohinansieto, kompressorin valinta)

Luvussa 5 kosketamme muita samankaltaisuuden metriikoita kuten Google Similarity Distance

# 2 Normalisoitu Pakkausetäisyys

**Kolmogorov kompleksisuus** Lyhimmän binääriohjelman pituus, joka palauttaa  $x$  syötteellä  $y$ , on *Kolmogorov kompleksisuus*  $x$ :stä syötteellä  $y$ ; tämä merkitään  $K(x|y)$ . Pohjimmillaan Kolmogorov kompleksisuus tiedostosta on sen äärimmäisesti pakatun version pituus.

**Normalisoitu Informaatioetäisyys** Artikkelissa [CV05] on esitelty *informaatioetäisyys*  $E(x, y)$ , joka on määritelty lyhimpänä binääriohjelmana, joka syötteellä  $x$  laskee  $y$ :n ja syötteellä  $y$  laskee  $x$ :n. Tämä lasketaan seuraavasti [CVdW04]

$$E(x, y) = \max\{K(x|y), K(y|x)\}. \quad (1)$$

Normalisoitu versio informaatioetäisyydestä ( $E(x, y)$ ), jota kutsutaan *normalisoiduksi informaatioetäisyydeksi*, on määritelty seuraavasti

$$NID(x, y) = \frac{\max\{K(x|y), K(y|x)\}}{\max\{K(x), K(y)\}}. \quad (2)$$

Tätä kutsutaan *samankaltaisuuden metriikaksi*, koska tämän on osoitettu [CV05] täyttävän vaatimukset etäisyyden metriikaksi. *NID* ei kuitenkaan ole laskettavissa tai edes semi-laskettavissa, koska Turingin määritelmän mukaan Kolmogorov kompleksisuus ei ole laskettavissa. Nimittäjän approksimointi annetulla kompressorilla  $C$  on  $\max\{C(x), C(y)\}$ . Osoittajan paras approksimaatio on  $\max\{C(xy), C(yx)\} - \min\{C(x), C(y)\}$  [CV05]. Kun *NID* approksimoidaan oikealla kompressorilla, saadaan tulos jota kutsutaan *normalisoiduksi pakkausetäisyydeksi*. Tämä esitellään formaalisti myöhemmin.

**Normaali Kompressor** Seuraavaksi esitämme aksioomia, jotka määrittelevät laajan joukon kompressoreita ja samalla varmistavat *normalisoidussa pakkausetäisyydessä* halutut ominaisuudet. Näihin kompressoreihin kuuluvat monet tosielämän kompressorit.

Kompressor  $C$  on *normaali* jos se täyttää seuraavat aksioomat,  $O(\log n)$  termiin saakka:

1. *Idempotenssi*:  $C(xx) == C(x)$  ja  $C(\lambda) = 0$ , jossa  $\lambda$  on tyhjä merkkijono,
2. *Monotonisuus*:  $C(xy) \geq C(x)$ ,
3. *Symmetrisuus*:  $C(xy) == C(yx)$  ja
4. *Distributiivisuus*:  $C(xy) + C(z) \leq C(xz) + C(yz)$ .

**Normalisoitu Pakkausetäisyys** Normalisoitua versiota *hyväksyttävästä etäisyydestä*  $E_c(x, y)$ , joka on kompressorin  $C$  pohjautuva approksimaatio *normalisoidusta informaatioetäisyydestä*, kutsutaan nimellä *Normalisoitu Pakkausetäisyys (NCD)* [CV05]. Tämä lasketaan seuraavasti

$$NCD(x, y) = \frac{C(xy) - \min\{C(x), C(y)\}}{\max\{C(x), C(y)\}}. \quad (3)$$

$NCD$  on funktioden joukko, joka ottaa argumenteiksi kaksi objektia (esim. tiedostoja tai Googlen hakusanoja) ja tiivistää nämä, erillisinä ja yhdistettyinä. Tämä funktioden joukko on parametrisoitu käytetyn kompressorin  $C$  mukaan.

Käytännössä  $NCD$ :n tulos on välillä  $0 \leq r \leq 1 + \epsilon$ , joka vastaa kahden tiedoston eroa toisistaan; mitä pienempi luku, sitä enemmän tiedostot ovat samankaltaisia. Tosielämässä pakkausalgoritmit eivät ole yhtä tehokkaita kuin teoreettiset mallit, joten virhemarginaali  $\epsilon$  on lisätty ylärajaan. Suurimmalle osalle näistä algoritmeista on epätodennäköistä että  $\epsilon > 0.1$ .

Luonnollinen tulkinta  $NCD$ :stä, jos oletetaan  $C(y) \geq C(x)$ , on

$$NCD(x, y) = \frac{C(xy) - C(x)}{C(y)}. \quad (4)$$

Eli etäisyys  $x$ :n ja  $y$ :n välillä on suhde  $y$ :n parannuksesta kun  $y$  pakataan käyttäen  $x$ :ää, ja  $y$ :n pakkauksesta yksinään; suhde ilmaistaan etäisyytenä bittien lukumääränä kummankin pakatun version välillä.

Kun kompressor on normaali niin  $NCD$  on normalisoitu hyväksyttävä etäisyys, joka täyttää metriikan yhtälöt, eli se on samankaltaisuuden metriikka.

### 3 Klusterointi

### 4 Käyttökohteet

#### 4.1 Klusteroinnin tulokset

#### 4.2 Kvantunnistus

### 5 Kohinansietokyky

Kun  $NCD$ :tä käytetään kahteen eri tiedostoon toista näistä voi pitää kohinalisena versiona ensimmäisestä. Progressiivisen kohinan lisääminen tiedostoon voi tuottaa tietoa mittarista(measure) itsestään. Tämän vastaavuuden perusteella voimme tehdä teoreettisen päätelmän odotetusta kohinan lisäämisen vaikutuksesta algoritmiin, mikä selittää miksi  $NCD$  voi saada suurempia arvoja kuin 1 joissain tapauksissa. [CAO07]

## 6 Huomioitavaa kompressorin valitsemisessa

## 7 Muita samankaltaisuuden metriikoita

### 7.1 Google Similarity Distance

#### Lähteet

- [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.
- [CVdW04] Cilibrasi, Rudi, Vitanyi, Paul ja Wolf, Ronald de: *Algorithmic Clustering of Music*. Web Delivering of Music, International Conference on, 0:110–117, 2004.