

Hva kan Wikipedia fortelle oss om kaffe?

Lars G. Johnsen

Magnus Breder Birkenes

Nasjonalbiblioteket

Wikipedia-akademiet 2015



- Wikipedia emneord og gruppering av innhold
- Strukturering av innhold
- Gruppering av artikler med innholdsord

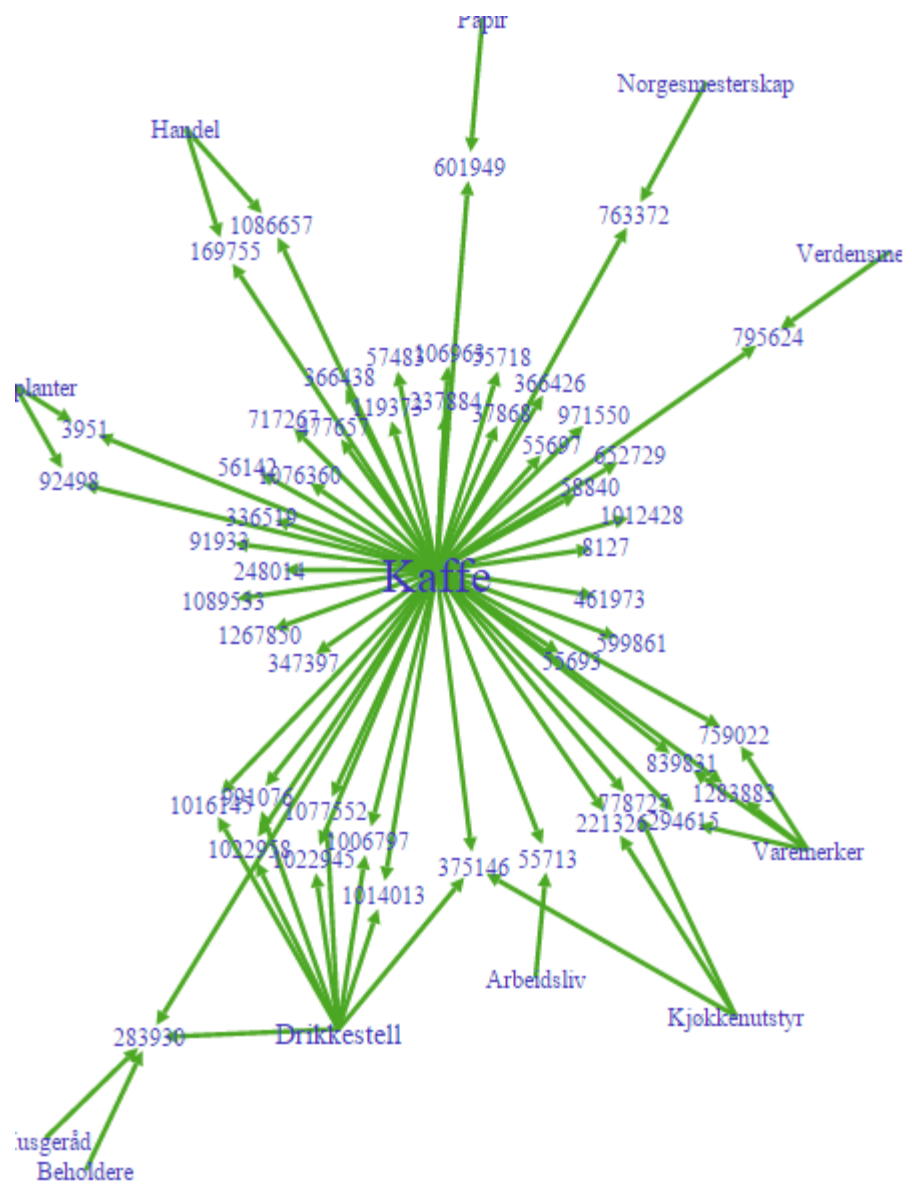
Relasjoner mellom innholds- og emneord

- Emner og tema
- Tekstklassifisering

Fra Wikipedia til data

- MySQL (wikipediadata)
- SQLite (arbeidsbase)
- Python (skripting)
- Gephi (visualisering)

F	Art.ID	Emne	Innhold
1	894284	Kortvinger	til
1	31563	Geologi	sies
1	26945	Rasisme	pasifister
1	189391	restkategori	to
1	407914	Bær	Vaccinium
1	478710	Tinget	etterkommes
1	526951	Kvinner	pensjonerte
1	187815	Menn	resulterte
1	1257701	Barokkomponister	Karl
1	947705	restkategori	på
1	410947	Menn	Trapezunt
1	134951	Göteborg	innbyggerantall



4	USA	iranske
1	Trekirker	bakover
1	restkategori	SkyCargo
1	Tanagarer	lokalt
1	Litauen	rådet
1	Sprengstoff	jobb
1	Geofysikk	rombølger
1	Jesuittordenen	Juizhou
1	Portal:Færøyene/artikler	átti
2	Menn	verdenssensasjon
1	Tyskere	gjeldende
4	Dramakomedier	Frances

PMI (Pointwise Mutual Information)

Dividerer

Innholdsord i Emne/Emnetotalen

med

Innholdsordtotalen/Totalen

Og vokter med

Innholdsord i Emne

Kaffe

kaffen starbucks iskaffe
kaffekanner arabica kaffe kaffepulveret
barista robusta coffea
nespresso kokende friele
café kaffehuset espresso kokekaffe
spesialkaffe löfbergs
bønnene liberica latte beholderen
lait baristakunst coffee
kaffebønner kaffetraktere kjeldsberg
gevalia kaffekjeler kapslene

Ekteskap

brudgommen ektefeller
gifte
bryllup **ekteskap**
likekjønnet bruden vielse
bryllupet gamos
vigsel parforhold ektevigsel
kvinne ektefellene inngått
forloveren inngå ekteskapsloven
gifter skilsmisse
hieros vigsler trolovelse
særeie felleseie inngåelse
polygami ektepakt
vigsleren ekteskapslov
ekteskapet likekjønnede

Bilmerker

maserati sportsbiler bilen
lastebiler

selskapet motors ford
bil modell

company produksjon fiat
rover daimler

fabrikken buick merket bilene

bilmerke motorer opel

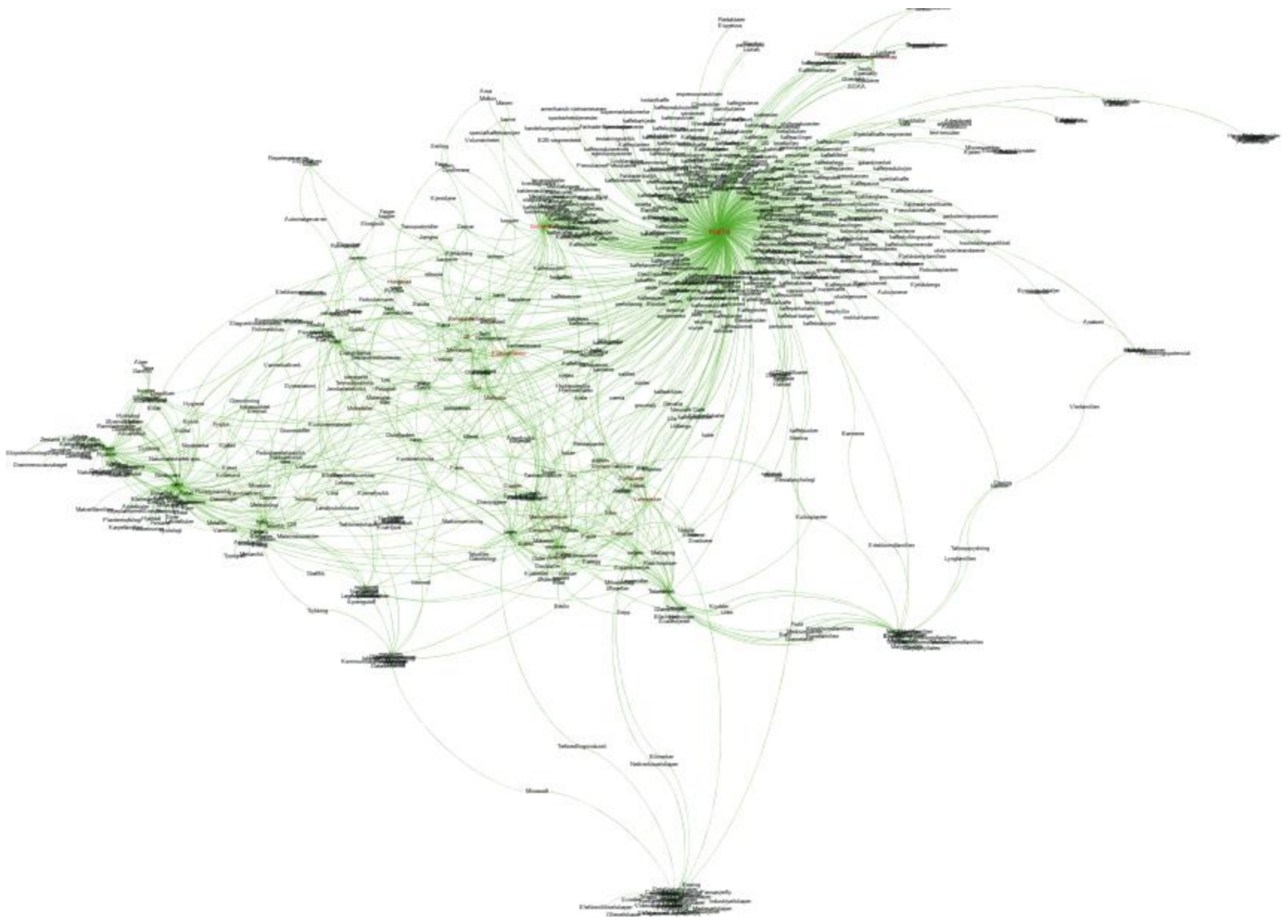
motor biler modeller

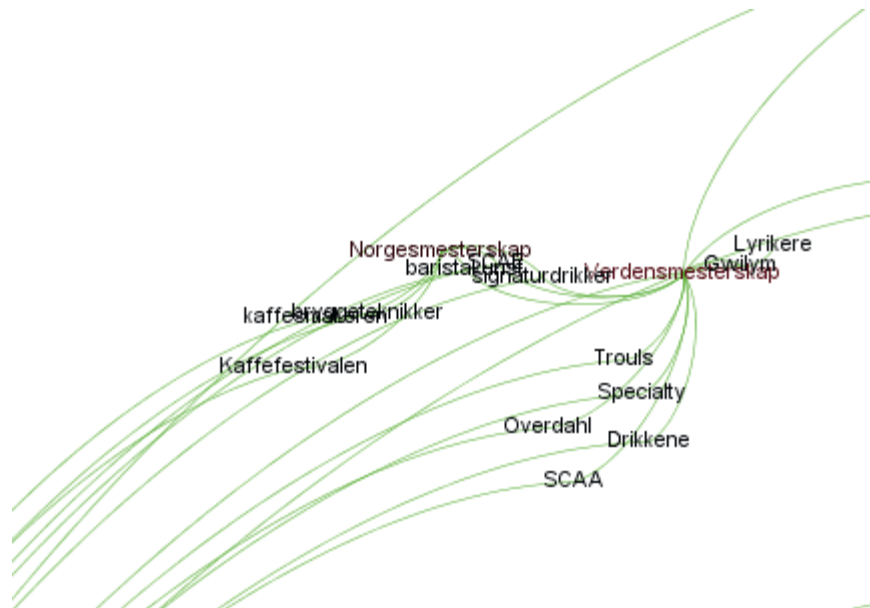
cadillac produserte chrysler

modellen modellene bilprodusent

Kaffe og undergrupper

- Innholdsordene står også i tekster med forskjellige emneord





Hva med gruppering på innholdsord?

- La «kaffe» være alle artikler som inneholder ordet «kaffe».

Assosiasjon

- For to ord «kaffe» og «kopp» i en artikkel
- Tell opp forekomster av «kaffe»
- Tell opp forekomster av «kopp»
- Ta produktet og del på størrelsen av artikkelen.
- Jo mindre artikkel, desto høyere score
- Jo fler «kaffe» eller «kopp», desto høyre score

Leksem vs. ordformer

Leksem vs. ordformer

