

Evaluation of machine-learned semantic clusters: Emergent concept structures vs. expert taxonomies

Seppo Nyrkkö

In this work, I study the correspondence between linguistic distributional syntactic analysis of term definitions and hand-made semantic term taxonomies. Automated syntactic dependency analysis is used for Finnish and English term definitions. Hand-made taxonomies written by human experts are compared to the machine learnt concept structures.

From Wikipedia texts to syntactic data, and further into high-dimensional vector space

For each term analyzed, the syntactic context is projected into a high-dimensional input vector space. Unsupervised machine learning is applied to the numerical syntactic evidence.

Data analysis methods

- SOM - Self-Organizing Map - 2D / 3D lattice representation of the high-dimension input space, model based on brain study
- GTM - Generative Topographic Map, a statistical analogy to SOM
- PCA, ICA - separation of *feature* and *noise* dimensions
- t-SNE - dimension reduction retaining local structures

Emergent conceptual structures can be seen in the grouping of frequent occurrences of highly similar features in an observed dataset. Emergent taxonomies are extracted from the unsupervised clustering. Human-designed semantic web ontologies and taxonomy structures are compared against the emergent clusters. The model can be evaluated even with short dictionary definitions.

Questions answered

Semantic Reasoning:

Which semantic features are seen in syntactic distributions?

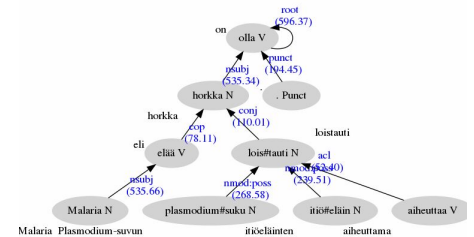
Hyperonymy		Hyponymy
Partonymy		Meronymy
Synonymy		Antonymy
Causes		Effects

Conceptual space model:

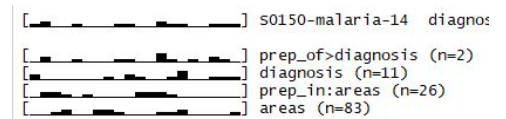
Can geometric approach be used in concept representation? Which theoretical feature dimensions are seen in the syntactic evidence?

- Colors, Sizes, Physical dimensions
- Shapes, Tastes, Symbols

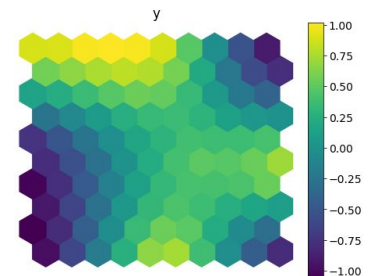
```
0: Malaria Malaria N ['Nom', 'Sg'] -> 1 nsobj 535.66
1: eli elia V ['Pst', 'Act', 'Sg'] -> 2 cop 78.11
2: horkka horkka N ['Nom', 'Sg'] -> 3 nsbj 535.34
3: on olla V ['Pres', 'Act', 'Sg'] -> 3 root 266.37
4: Plasmodium-suvun plasmodium#suku N ['Gen', 'Sg'] -> 7 nmod:poss 268.58
5: itioeläinten itioeläin N ['Gen', 'Pl'] -> 7 nmod:poss 239.51
6: aiheuttaa aiheuttaa V ['Pass', 'Nom', 'Sg', 'AgPrc'] -> 7 acl 52.48
7: loistauti loistauti N ['Nom', 'Sg'] -> 2 conj 118.01
8: . Punct [] -> 3 punct 194.45
```



dependency syntax analysis



high-dimensional vector space



emergent concept structures in output space

