Tim auf der Landwehr
t.auf.der.landwehr@student.rug.nl
S2548682
Last revision: December 10, 2013

# Learning from Data
# Assignment 5

## Summary

This is the assignment 5 in the course Learning from Data at the University of Groningen.

## Exercise 1: Cluto-formatter

A script that converts the given file `name_vectors_small.txt` (and later `name_vectors.txt`) to the dense matrix format is provided in `cluto_formatter.py`. This can be run the following way:

```
python3 cluto_formatter.py --in name_vectors.txt --out name_vectors.mat --labels
```

When the parameter `--labels` is used, the files
`[output_filename].rlabel` and
`[output_filename].clabel`
are created in addition.

## Exercise 2: vcluster with cluto

The file that was created in exercise 1 is now processed with vluster.

```
vcluster -showtree -sim=cos -clmethod=agglo -plottree=plot.ps name_vectors_small.mat 7
```

Combining the labels with the classes again, and sorting them by class, the following distribution can be obtained:

```
Class 0:        IsraÃ«l         Rusland        Class 3:
  Donner        Duitsland       Barcelona        Amsterdam      Class 5:
  Zalm          Frankrijk       Vlaanderen       Rotterdam        A
  Bos           Europa          Engeland         Utrecht
  Jan           Ajax            Anderlecht       Den_Haag       Class 6:
  De_Boer       PSV                              Brussel          VVD
  De_Vries      BelgiÃ«                          Parijs           CDA
                ItaliÃ«         Class 2:                          PvdA
                Feyenoord        Balkenende     Class 4:
Class 1:        China            Blair            Ahold
  Bush          Spanje           SchrÃ¶der        ING
  Nederland     Turkije          Kok              KLM
  Irak          Amerika          Berlusconi       Philips
  VS
```

This shows how well this classifier works already for quite small datasets.

## Exercise 3: Large data set

Now the same procedure is used on the big dataset provided in `name_vectors.txt`. The following listing shows the classification in 10 classes, the entities per class are limited to the first 10. From this you can see, that the classification worked quite well.

```
Class 0:                          Lance_Armstrong    Bush
  CD&V             Class 3:      Class 6:            Balkenende
  VLD                Antwerpen     A                 Verhofstadt
  SP.A               West-Vlaanderen B               Blair
  CD&V_/_N-VA        Gent          C                 George_Bush
  VVD                Brussel       Ahold             SchrÃ¶der
  CDA                Leuven        D                 Sharon
  Groen_!            Amsterdam     H                 Berlusconi
  D66                Brugge        ING               Kok
  N-VA               Oostende      Fortis            Abbas
  Vlaams_Belang      Limburg       Delhaize        Class 9:
Class 1:             Kortrijk      Philips           Nederland
  Peeters          Class 4:      Class 7:            BelgiÃ«
  Donner             Bos           MÃ¡xima           Irak
  Zalm               FinanciÃ«n    Maria             Vlaanderen
  Yves_Leterme       Justitie      Astrid            VS
  Verdonk            Sport         Diana             Tom_Boonen
  Janssens           Cultuur       Nathalie          Frankrijk
  Frank_Vandenbroucke Onderwijs    Electrabel        Duitsland
  Bernhard           Kunst         Els               Europa
  Kurt               Economie      Victoria          ItaliÃ«
  Filip              Jeugd         Kim
Class 2:           Class 5:        Kim_Clijsters
  Albert             Filip_Dewinter Class 8:
```

# Exercise 4&5: Label classes

In this task, an approach to define the 10 classes by their containing entities is given. Especially the four tags *MISC*, *ORG*, *LOC* and *PER* are considered.

**Class 0:** shortcuts

**Class 1:** basically names → PER, but *Zalm* is a fish

**Class 2:** name → PER

**Class 3:** cities → LOC

**Class 4:** different abstract things → MISC

**Class 5:** names → PER

**Class 6:** mixed: letters and companies → partly ORG

**Class 7:** mixed: names and companies → partly PER and ORG

**Class 8:** names → PER

**Class 9:** mostly countries → LOC

It is not completely possible to assign perfect class names to the data, but the results are quite sophisticated.