# LED-A: a web app for measuring distances in the sound components among local dialects

Wilbert Heeringa & Vincent van Heuven & Hans Van de Velde

FRYSKE 🦋 AKADEMY

Methods XVII

Mainz, August 1, 2022

# Introduction

- We present a new app for measuring dialect variation with Levenshtein distance.
- Proposed by Vladimir Levenshtein in 1965.
- Introduced in dialectology by Brett Kessler in 1995.
- He measured linguistic distances among Irish Gaelic dialects.
- Later widely used for many other language varieties.
- Calculate the cost of changing one string of characters into another.

# Levenshtein distance

- Example: *milk* is pronounced as [mɛlək] in the dialect of Haarlem and as [mɔlkə] in the dialect of Grouw.

- Find cheapest mapping of [mɛlək] → [mɔlkə]:

| cumulative maxim. cost | 1 | 2 | 3 | 3.5 | 4.5 | 5 |
|---|---|---|---|---|---|---|
| | m | ɛ | l | ə | k | |
| | m | ɔ | l | | k | ə |
| actual cost | | 1 | | 0.5 | | 0.5 |

- Raw distance:            Normalized distance:
  $1 + 0.5 + 0.5 = 2$      $2 / 5 = 0.4 = 40\%$

# Aggregated distance

- Example: calculate pronunciation differences between Grouw and Haarlem dialects for 6 words:

|        | Grouw  | Haarlem | cost | max. cost | norm. cost |
|--------|--------|---------|------|-----------|------------|
| work   | ʊɣrk   | ʊɛrək   | 1.5  | 4.5       | 0.33       |
| ship   | skɪp   | sxɪp    | 1    | 4         | 0.25       |
| finger | fɪŋər  | vɪŋər   | 1    | 5         | 0.2        |
| wine   | ʋin    | ʋɛin    | 0.5  | 3.5       | 0.14       |
| house  | huz    | hœys    | 0.5  | 3.5       | 0.14       |
| milk   | mɔlkə  | mɛlək   | 2    | 5         | 0.4        |
|        |        |         | 6.5  |           | 1.46       |

- Raw distance:      Normalized distance:
  $6.5/6 = 1.08$      $1.46/6 = 0.243 = 24.3\%$

# Software

- Visual Dialectometry (VDM).

- DiaTech: related to VDM, but it is a webapp and includes also Levenshtein distance.

- R$u$G/L$^{04}$: set of functions to be entered as command line commands using a keyboard.

- Gabmap: webapp at `gabmap.nl`, Docker version at `https://github.com/pebbe/Gabmap-docker`.

- Python packages: editdistance, LingPy.

- R packages: iL04 (`http://www.let.rug.nl/~kleiweg/L04/R/`), stringdist, alineR, dialectR (at GitHub).

- Goal: make dialectometry as easy as possible:
  - **L**evenshtein **E**dit **D**istance **A**pp,
  - refers to 'Lampje' or Gyro Gearloose's 'Little Helper' in Donald Duck comics; 'Lampje' is Dutch for Little Lamp (LED = Light Emitting Diode = lamp).
- Availability:
  - `https://www.led-a.org`
  - example data sets under 'Examples'.

FRYSKE 🏵 AKADEMY

# Distance measures

- Binary item comparison (Séguy 1973)
- Levenshtein distance plain
  indel is 0.5 or 1
- Levenshtein distance IPA feature-based (Almeida & Braun)
  indel scaled between 0 and 0.5 or joint scaling of subsitions and indels
  between 0 and 1.
- Levenshtein distance PMI-based
  Wieling, Prokić & Nerbonne (2009), Wieling (2012)

# VC-sensitive

The minimum cost is based on an alignment in which:

- a vowel matches with a vowel
- a consonant matches with a consonant

Optionally the user can allow:

- [j], [w], [i] or [u] to match with anything
- [ə] (schwa) or [ɐ] to match with a sonorant

Bolognesi & Heeringa (2002), Heeringa (2004), Wieling et al. (2009)

# Levenshtein

- Both 'raw' and normalized Levenshtein distances can be calculated.
- It is possible to measure distances due to differences in only vowels or only consonants (indels or substitutions).

# Input: transcriptions

- Excel sheet (Microsoft Excel/LibreOffice Calc).
- The transcriptions should be in IPA Unicode (use `https://westonruter.github.io/ipa-chart/keyboard/`)
- Multiple transcriptions per item are possible.

# Format (1)

Rows are locations, columns are items

|  | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|---|---|---|---|---|---|---|---|
| 1 |  | **work** | **ship** | **finger** | **wine** | **house** | **milk** |
| 2 | Delft | wɛrək | sxɪp | vɪŋər | ʋæ˙ĭn | hœ˙z | mɛlək |
| 3 | Grouw | ʋʏrk | skɪp | fɪŋər | ʋin | hu˙z | mɔ˙lkə |
| 4 | Haarlem | ʋɛrək | sxɪp | vɪŋər | ʋɛin | hœỹs | mɛlək |
| 5 | Hattem | ʋɛ̠rək | sxɪp | vɪŋəř | ʋi˙n | ys | mɛlək |
| 6 | Lochem | ʋɑrək | sxɪp | vɪŋəř | ʋin | hys | mɛlək |

# Format (2)

Rows are items, columns are locations

|   | 1 | 2 | 3 | 4 | 5 | 6 |
|---|---|---|---|---|---|---|
| 1 | | **Delft** | **Grouw** | **Haarlem** | **Hattem** | **Lochem** |
| 2 | work | wɛrək | ʋʏrk | ʋɛrək | ʋɛ̝rək | ʋɑrək |
| 3 | ship | sxɪp | skɪp | sxɪp | sxɪp | sxɪp |
| 4 | finger | vɪŋər | fɪŋər | vɪŋər | vɪŋəř | vɪŋəř |
| 5 | wine | ʋæ·ĭn | ʋin | ʋɛin | ʋi·n | ʋin |
| 6 | house | hœ·z | hu·z | hœÿs | ys | hys |
| 7 | milk | mɛlək | mɔ·lkə | mɛlək | mɛlək | mɛlək |

# With LED-A can be processed ...

- vowels, pulmonic consonants and voiced implosives;
- primary stress, secondary stress;
  by processing ' and ˌ as segments
- extra short, normal, half long, long;
  by preprocessing: ă → a, a → aa, aˑ → aaa, aː → aaaa
- aspirated, labialized, palatalized, velarized, pharyngealized, nasalized;
  by averaging with h, w, j, ɣ, ʕ, n respectively.

# Maps

- Easy to create maps, only coordinates are required (no outline required).

# Coordinates

Rows are items, columns are locations

|   | 1 | 2 | 3 |
|---|---|---|---|
| 1 |   | **lat** | **long** |
| 2 | Delft | 52.00667 | 4.35556 |
| 3 | Grouw | 53.09456 | 5.83745 |
| 4 | Haarlem | 52.38084 | 4.63683 |
| 5 | Hattem | 52.475 | 6.06389 |
| 6 | Lochem | 52.15917 | 6.41111 |

Coordinates are taken from GeoNames.

# Output

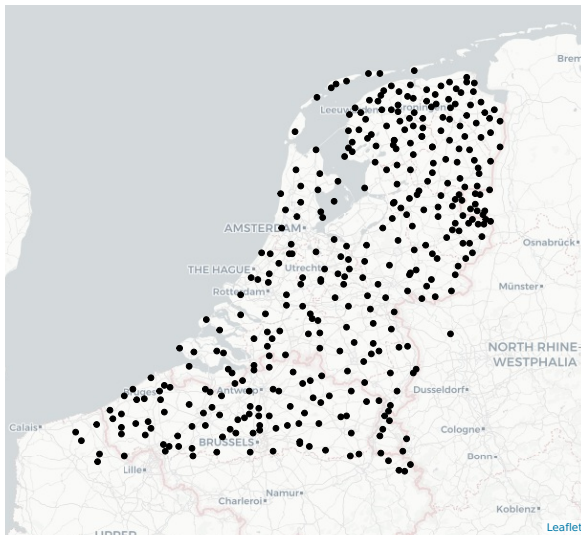- Output: individual word pair distances and/or aggregated distances.

# Visualization

Example data set:

- Reeks Nederlandse Dialectatlassen, compiled by E. Blancquaert and W. Pée.
- Texts from 1922–1975, 1956 local dialects, 139 sentences each.
- We selected 360 dialects, 166 words.
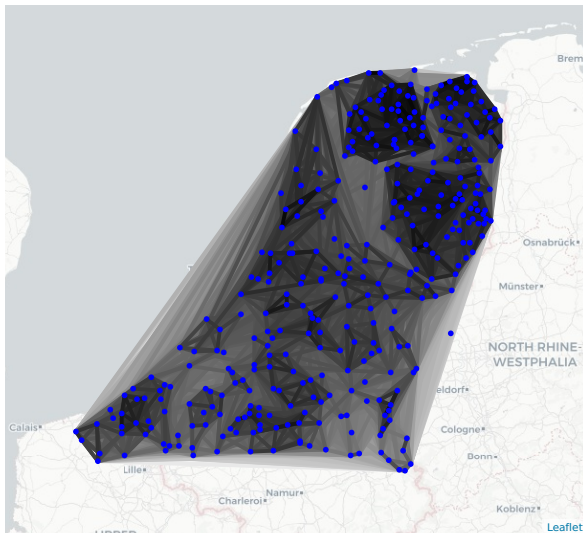- Standard Dutch and Standard German were added.

# Measurements

- We measure distances using PMI-based Levensthein distance.
- Length and diacritics are processed.

Distribution of the 360 varieties in the Dutch dialect area.

# Beam maps

- Introduced by Goebl ($\pm$ 1983).
- The locations of local dialects are connected to each other by straight lines in a map.
- Darker lines represent smaller distances, lighter lines represent larger distances.

Beam map showing Levenshtein distances among local dialects. Darker lines represent larger distances.
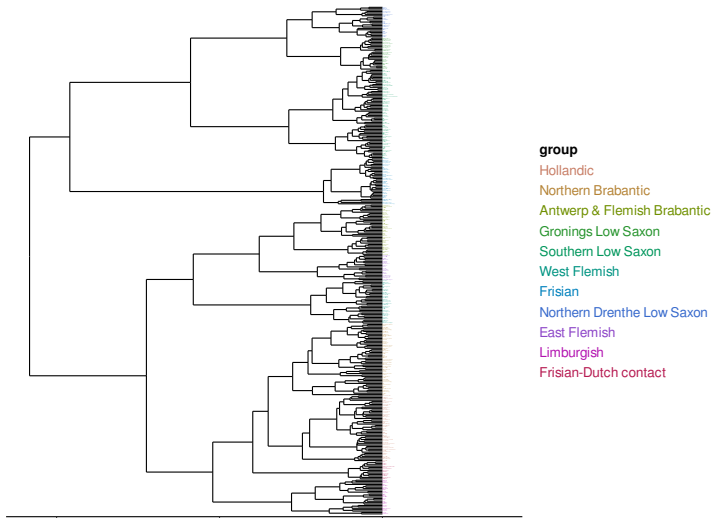
FRYSKE 🎺 AKADEMY

# Cluster analysis

- Introduced by Goebl ($\pm$ 1982) in dialectometry.
- Group objects in such a way that objects in the same group (called a cluster) are more similar (in some sense) to each other than to those in other groups (clusters). (Wikipedia)
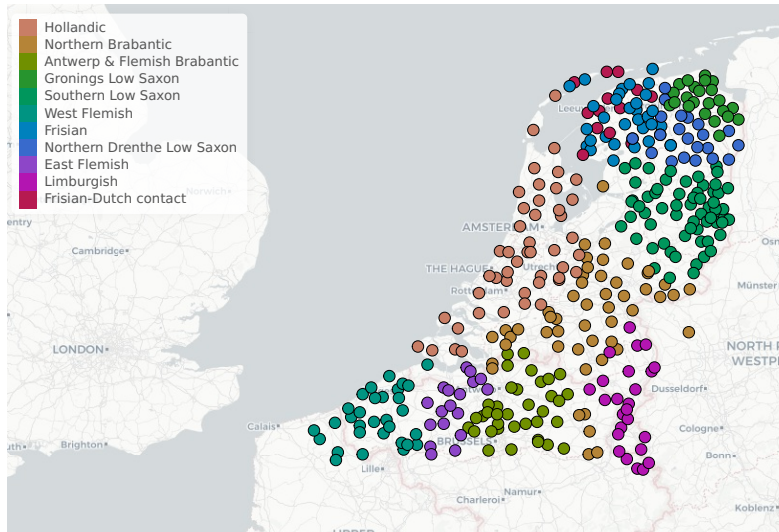
# Cluster methods

In LED-A five cluster methods are available:

- **Single**-**Linkage**: chaining effect, no clear group structure.

- **Complete**-**Linkage**: compact clusters of about equal size.

- **UPGMA**: results reflect the original distances most closely.

- **WPGMA**: in case of a more irregular distribution.

- **Ward's** method: minimizes the variance in the clusters; it usually creates compact, even-sized clusters (Szmrecsanyi, 2012)

**group**
Hollandic
Northern Brabantic
Antwerp & Flemish Brabantic
Gronings Low Saxon
Southern Low Saxon
West Flemish
Frisian
Northern Drenthe Low Saxon
East Flemish
Limburgish
Frisian-Dutch contact

Dendrogram obtained using Ward's method. The tree structure explains 50.8% of the variance in the original distances.

Legend:
- Hollandic
- Northern Brabantic
- Antwerp & Flemish Brabantic
- Gronings Low Saxon
- Southern Low Saxon
- West Flemish
- Frisian
- Northern Drenthe Low Saxon
- East Flemish
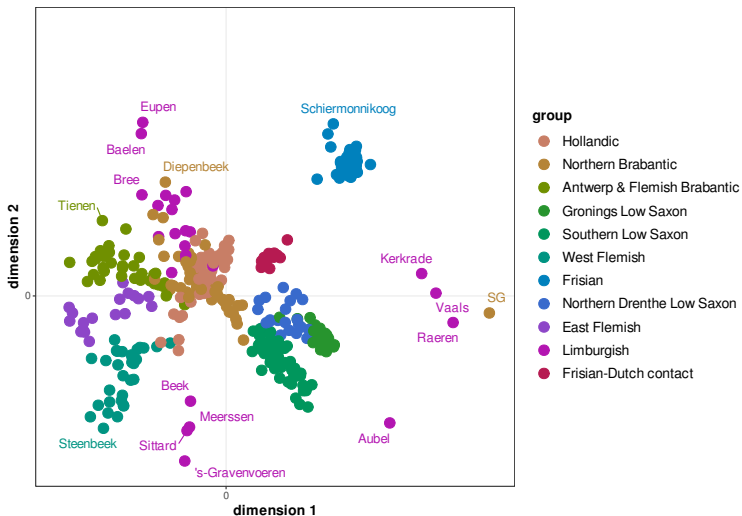- Limburgish
- Frisian-Dutch contact

Eleven groups derived from the dendrogram that was obtained using Ward's method.

# Multidimensional scaling

- Introduced by Embleton (1993) in dialectometry.
- Put the local dialects on a map so that the distances in two-dimensional space reflect the distances in the matrix as closely as possible.
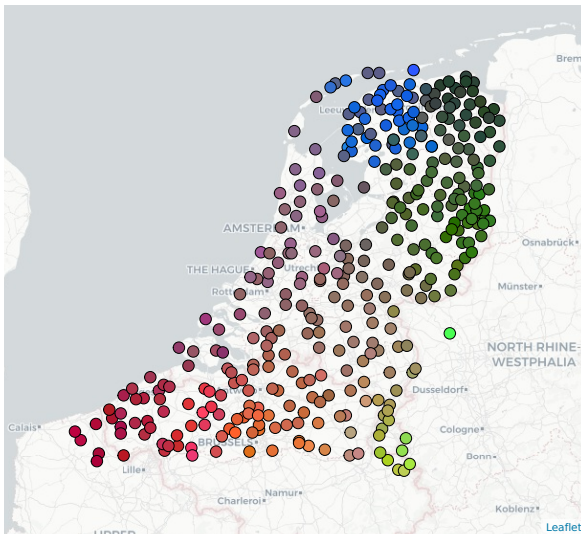
# Multidimensional scaling methods

- **Classical** (metric):
  original algorithm, proposed by Togerson (1952).
- **Kruskal's** non-metric:
  results reflect the original distances most closely.
- **Sammon's** non-linear (metric) mapping:
  points are shown more dispersed.
- **t-SNE** (t-distributed stochastic neighbor embedding):
  reveals (otherwise hidden) patterns, is stochastic.

Using Kruskal's non-metric MDS the 362 dimensions are reduced to 2. They explain 78.5% of the variance in the original distances.
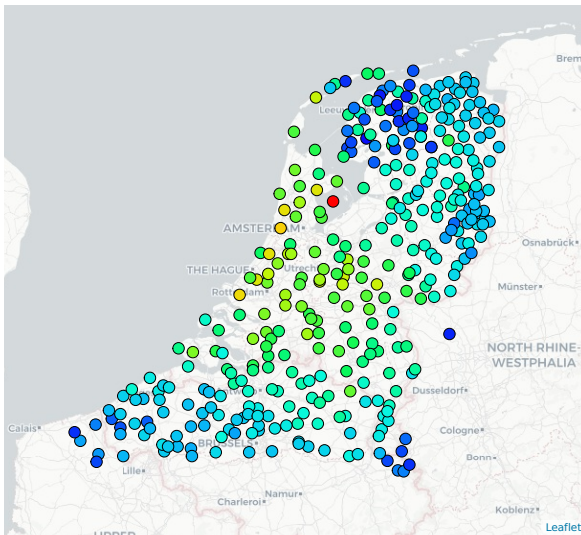
# Multidimensional scaling

- With MDS scaling to three or more dimensions is possible as well.
- Scale to three dimensions so that each local dialect is represented by three values $x$, $y$ and $z$.
- Let $x$ be the intensity of red, $y$ be the intensity of green and $z$ be the intensity of blue.
- Introduced by Nerbonne, Heeringa & Kleiweg (1999).

RGB map where $x$ determines inversely the intensity of red, $y$ determines the intensity of blue and $z$ determines the intensity of green.

# Reference point maps

- Introduced by Goebl ($\pm$ 1982).
- Compare local dialects to a reference point (e.g. standard language, proto-language).
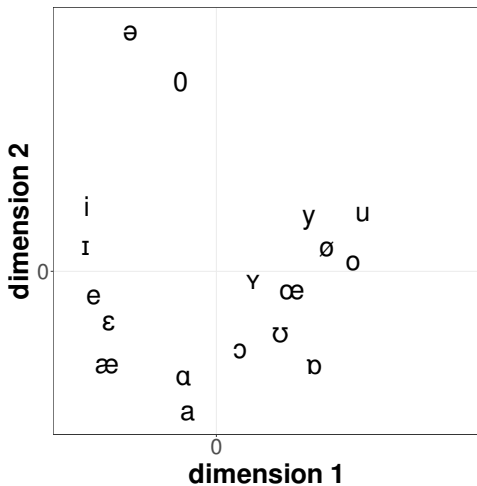- Coloring according to rainbow scheme: red is most similar, blue is most distant.

Dutch dialects compared to Standard Dutch. Red dots represent strongly related dialects, blue dots more remote ones.
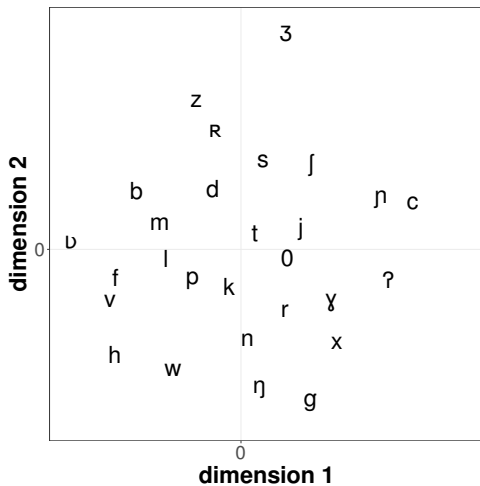
# Segment distances

- Show feature-based or PMI-based segment distances.
- Reduce the distances among the segments to two dimensions with Kruskal's non-metric MDS.

## Vowels

# Consonants

# Thanks!

# Thanks!

The slides are available at:

`led-a.org/slides.pdf`