

Combining GIS and R for dialectal data analysis and representation

Kristel Uihoaed
Maarja-Liisa Pilvik
Eleri Aedmaa
Siim Antso
University of Tartu

Outline of the presentation

- About the project “Applying Spatial Data in Linguistics”
- The aim of the presentation
- About the Corpus of Estonian Dialects (CED)
- Digitizing dialect atlas data
- Analyzing and visualizing data with R
- Using R to combine data from different sources
- Conclusions

About the project (RuRaKe)

- <http://rurake.keeleressursid.ee/>
- “Ruumiandmete rakendamine keeleteaduses” (Applying Spatial Data in Linguistics)
- Initially:
 - Adding geographical information to the Corpus of Estonian Dialects
 - Digitizing dialect maps of the “Väike Eesti murdeatlas” (“Small atlas of Estonian dialects”, 1955) by Andrus Saareste
- Creating resources for combining and analyzing different data sources with a spatial dimension
- Bringing dialect data closer to people not actively working in the field

The aim of the presentation

- Show, how combining different data sources (atlas data and corpus data) and applications (GIS-systems and statistical computing software R) is beneficial for getting a more diverse picture of linguistic variation and the language phenomena
- 3 alternation phenomena:
 - *sui / suvi* 'summer' (Saareste 1955)
 - *suurem / suuremb* 'bigger' (Saareste 1955)
 - *ei ole / pole* 'is not' (Saareste 1955)

About the Corpus of Estonian Dialects

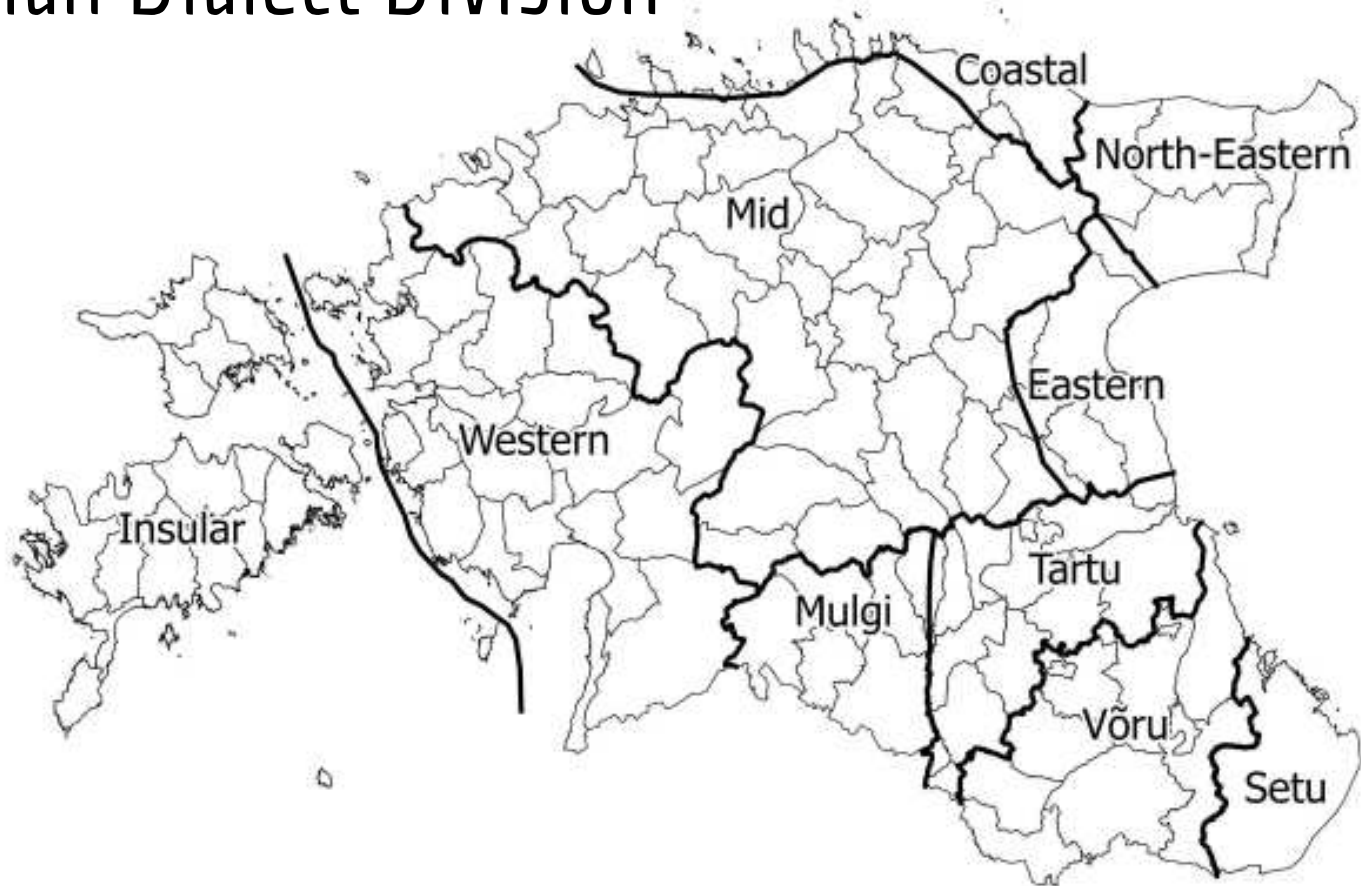
- <http://www.murre.ut.ee/mkweb>
- The electronic corpus contains text from all the Estonian dialects.
- The corpus mainly consists of transcribed and morphologically annotated texts of interview recordings from the 1960s and 1970s, as well as the recordings themselves and the metadata.
- Morphological annotation is done manually. Annotation includes information about lemma, part-of-speech, morphological categories, and meanings.

```

1 <?xml version="1.0" encoding="UTF-8"?>
2 <xml-stylesheet type="text/xsl" href="liivike_tekst.xsl"?>
3 <liivike>
4 <info><tyyp>suuline</tyyp><kihelkond>Hlj</kihelkond><keel>eesti</keel><murre>Banna</murre><kirjaviis>lihtsustatud transkriptsioon</kirjaviis><muu>Bannikumurre, Haljala
murrek, Potju küla, Marie Pihlakas (85-aastane). Lõidistanud (MMH 135, 1959. aastal) ja liitreeerinud (MKI MT 238) Mari Must (MM), Üle kuulunud ja täiendanud M-L. Kalvik
2006. aastal. KJ ? Marie Pihlakas, ? ? neesterahvas, kes juttu eesalt täiendab.</muu><hindistus_lindistaja="Mari Must" aasta="1959">EMH0135</hindistus><liitreeering
liitreeerija="Mari Must" aasta="1959">MKI MT 238</liitreeering><kylä_lõngitõud="26.20667" latitõud="59.46056">Potju</kylä><intervjuerija_id="?" haridus="NA" synniaasta="NA"
synnikoht="NA" vanus="NA" sugu="M">neesterahvas, kes juttu eesalt täiendab</intervjuerija><keelejuht_id="KJ" haridus="NA" synniaasta="1974" synnikoht="Potju" vanus="85"
sugu="M">Marie Pihlakas</keelejuht><intervjuerija_id="MM" haridus="kõrg" synniaasta="1920" synnikoht="Tartu" vanus="39" sugu="N">Mari Must</intervjuerija></info><sisu id
="HLJ_Marie_Pihlakas_EMH0135_synt">
5 <lause id="11" koneleja="KJ"><sona id="11_s1" meta="yahenärk">(-)</sona><sona id="11_s2" lemma="olema" vorm="pers.ind.ipf.sg.3." liik="V">oli</sona></lause><lause id=
"12" koneleja="MM"><sona id="12_s1" meta="intervjuerija">jahh</sona><sona id="12_s2" meta="yahenärk">.</sona></lause><lause id="13" koneleja="KJ"><sona id="13_s1"
lemma="ja" liik="Konj">ja</sona><sona id="13_s2" lemma="giis" liik="ProAdv">giis</sona><sona id="13_s3" meta="yahenärk">=</sona><sona id="13_s4" lemma="e" liik="Par">e
</sona><sona id="13_s5" meta="yahenärk">(<...></sona><sona id="13_s6" lemma="ütlena" vorm="pers.ind.ipf.sg.3." liik="V">*ütles</sona><sona id="13_s7" lemma="yeel" liik=
"Adv">viel</sona><sona id="13_s8" lemma="ema" vorm="ag.all." liik="S">emale</sona><sona id="13_s9" meta="yahenärk">.</sona><sona id="13_s10" lemma="enne" liik="Adv">
*enne</sona><sona id="13_s11" lemma="et" liik="Konj">et</sona><sona id="13_s12" meta="yahenärk">=</sona><sona id="13_s13" lemma="e" liik="Par">e</sona><sona id=
"13_s14" meta="yahenärk">(<...></sona><sona id="13_s15" lemma="ära" vorm="pers.imp.px.sg2." liik="Mn">ära</sona><sona id="13_s16" lemma="giis" liik="Par">giis
</sona><sona id="13_s17" lemma="laps" vorm="pl.part." liik="S">*lapsi</sona><sona id="13_s18" lemma="laskna" vorm="pers.imp.px.neg." liik="V">lase</sona><sona id=
"13_s19" lemma="hant" vorm="ag.sl." liik="S">*hant</sona><sona id="13_s20" lemma="nänena" vorm="inf." liik="V">*nenna</sona><sona id="13_s21" lemma="et" liik="Konj"
>et</sona><sona id="13_s22" lemma="karistana" vorm="pers.imp.px.sg.2." liik="V">karista</sona><sona id="13_s23" lemma="nenad" vorm="pl.part." liik="ProS">neid
</sona><sona id="13_s24" lemma="ka" liik="Par">kaa</sona><sona id="13_s25" meta="yahenärk">.</sona></lause>
6 <lause id="14" koneleja="MM"><sona id="14_s1" meta="intervjuerija">niia</sona></lause>

```

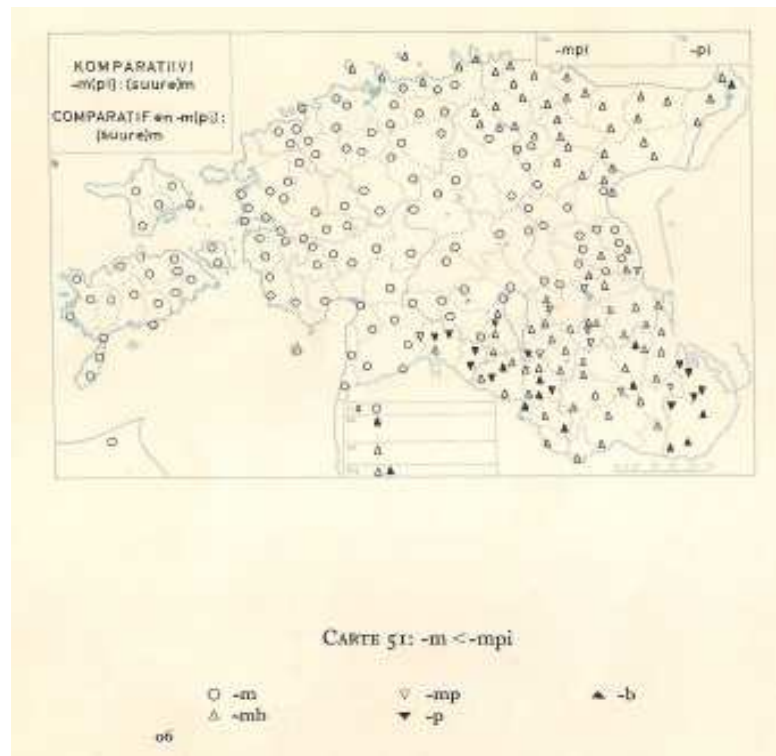
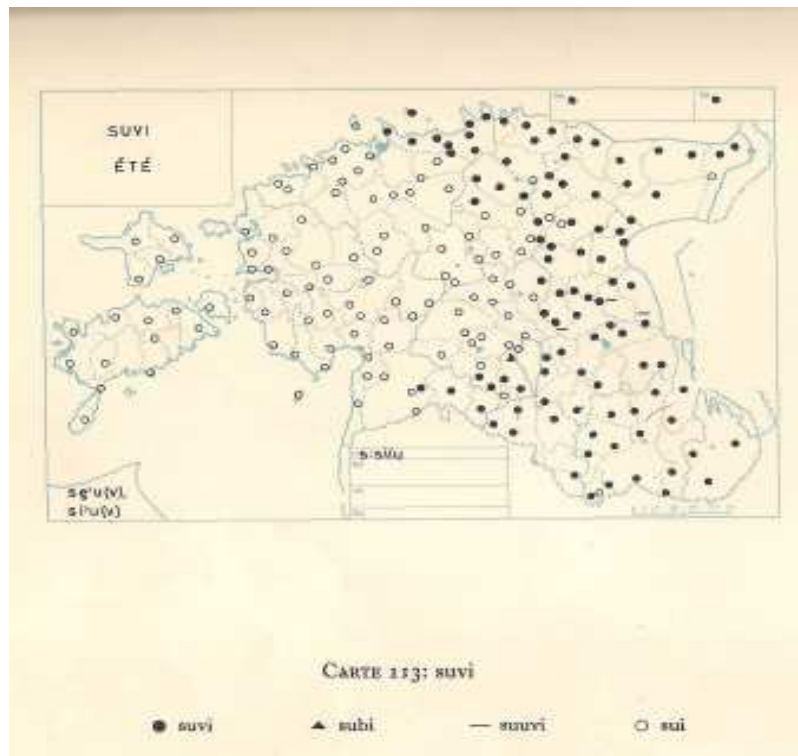
Estonian Dialect Division



Atlas data and using GIS

- A. Saareste's two atlases:
 - 1938/1941 "Eesti murdeatlas" ("Atlas of Estonian Dialects") with 66 maps
 - 1955 "Väike Eesti murdeatlas" ("Small Atlas of Estonian Dialects") with 125 maps
- "Atlas of Estonian Dialects" (AED) contains maps with text (a la Jules Gilliéron)
- In "Small Atlas of Estonian Dialects" (SAED) Saareste also used map symbols to visualize dialect data
- Mostly maps of vocabulary in AED (about 70%)
- SAED contains more morphology and phonology (57%) as well as vocabulary (41%)
- Using QGIS
 - Scanning and georeferencing
 - Assigning phonetically transcribed linguistic data (words, phrases, constructs) to the centroids of villages where the data was collected

Andrus Saareste's dialect maps





CARTE 30: ei ole < e(p)i ole-k

Using R

Perks of using R for dialectal data manipulation:

- Ready made analysis tools support limited number of statistics
- Possibility to manipulate dialect borders
- All the steps of the analysis are in the same environment
- A large number of statistical methods is supported in R

(9144 packages available at

http://cran.r-project.org/web/packages/available_packages_by_date.html)

R packages

Packages for visualization:

Ggplot2 (Wickham 2009)

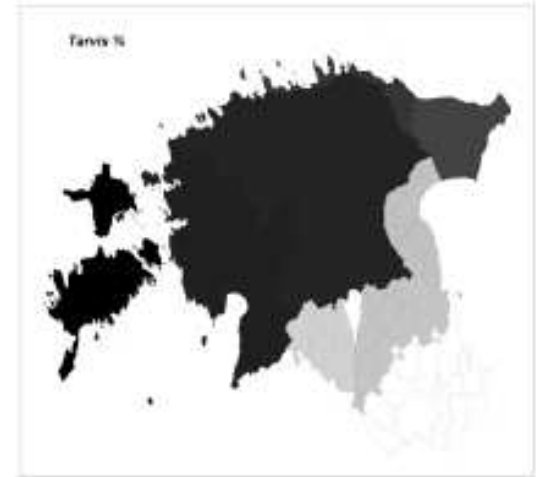
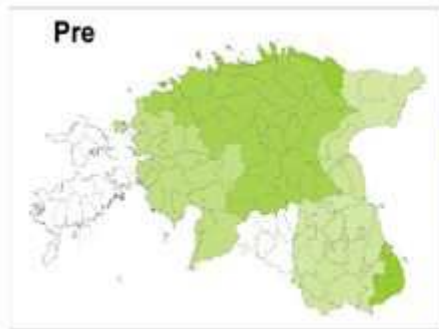
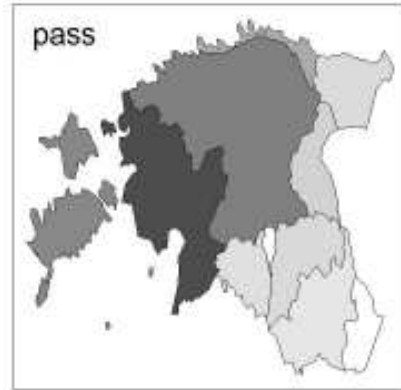
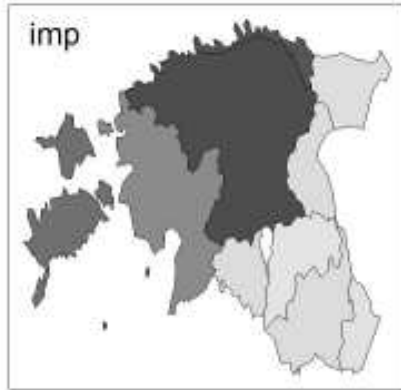
Ggmap (Kahle & Wickham 2013)

Packages for processing geodata:

Rgdal (Bivand et al. 2015)

Rgeos (Bivand & Rundel 2015)

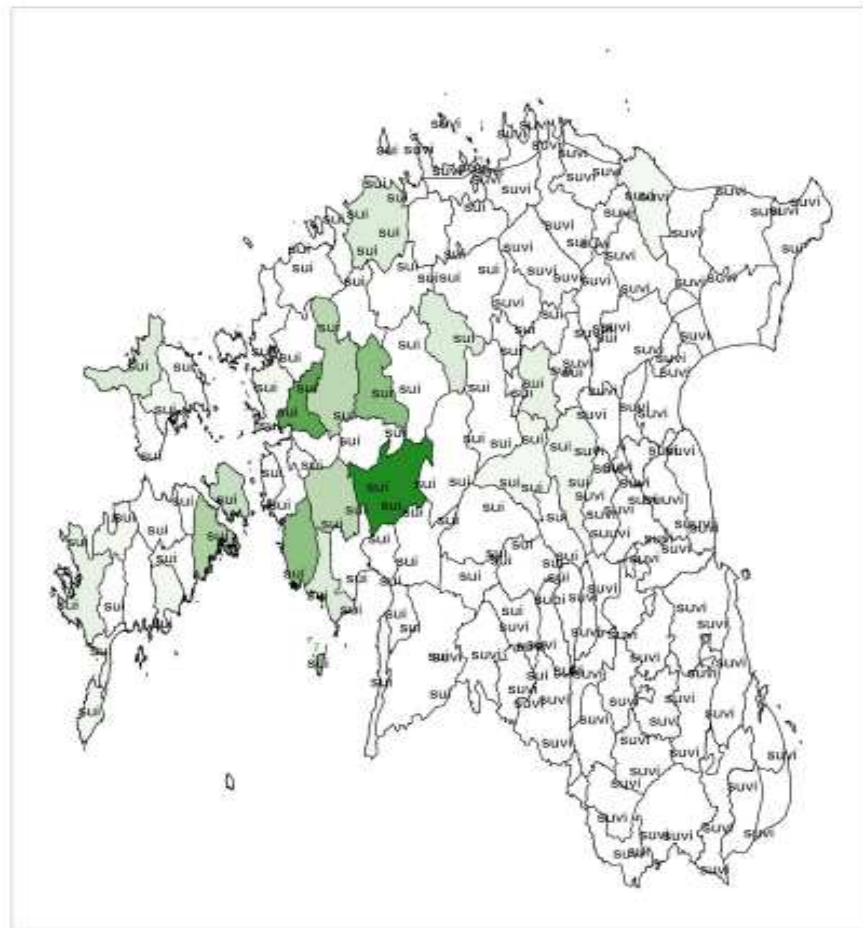
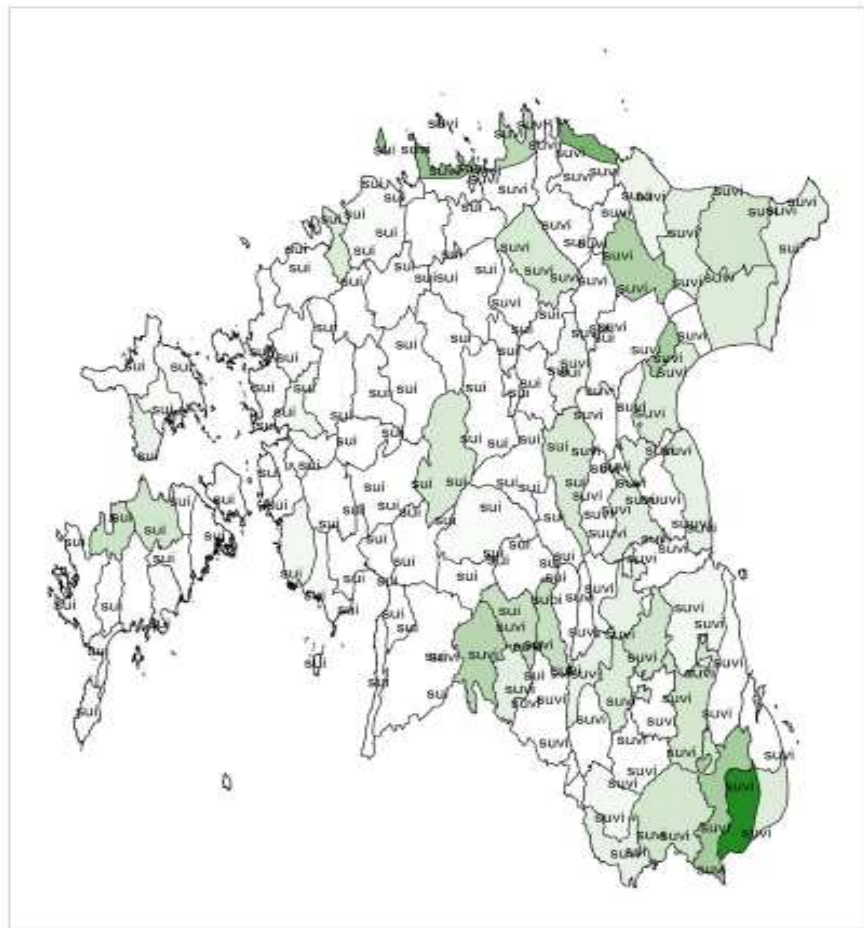
Studies based on CED data

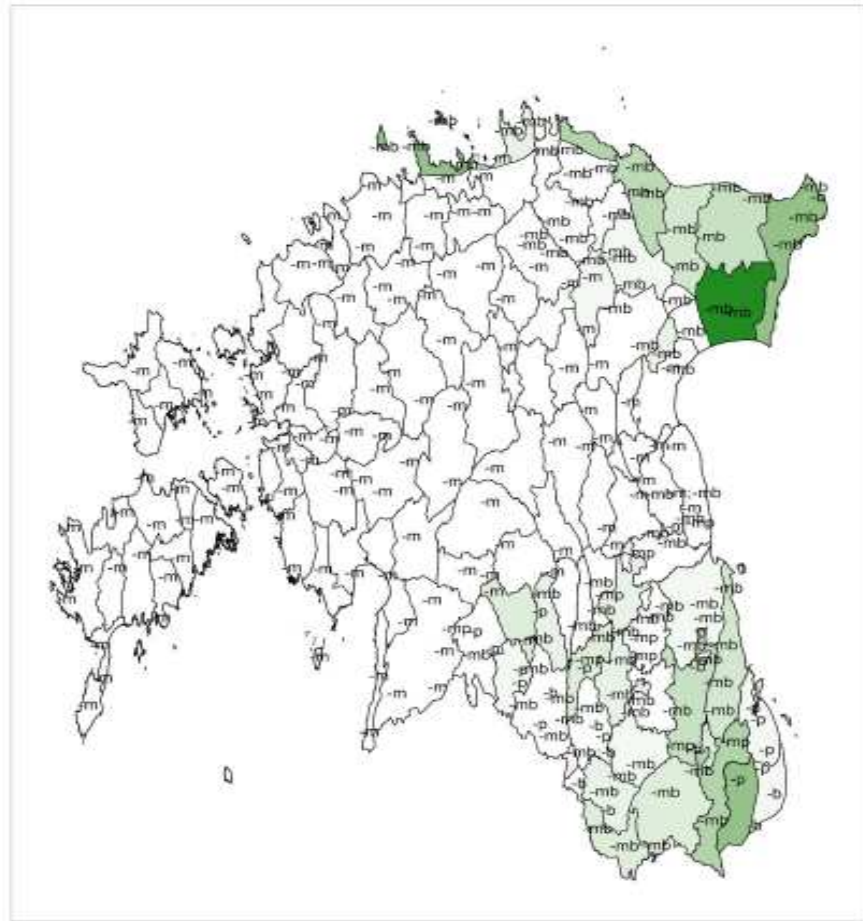
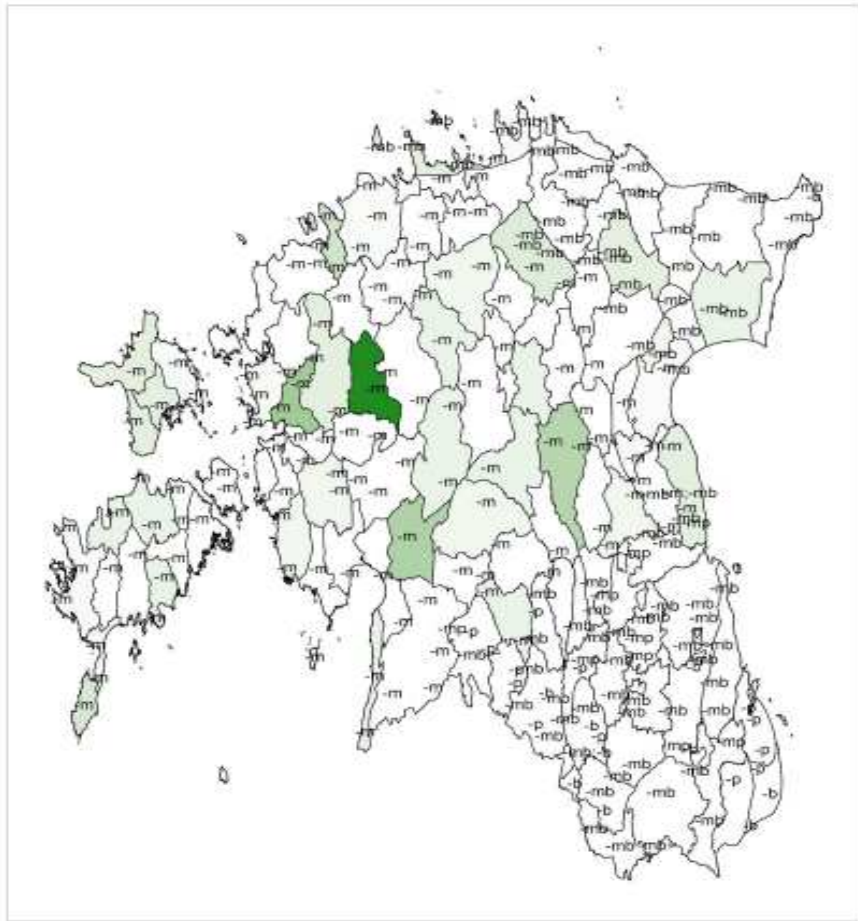


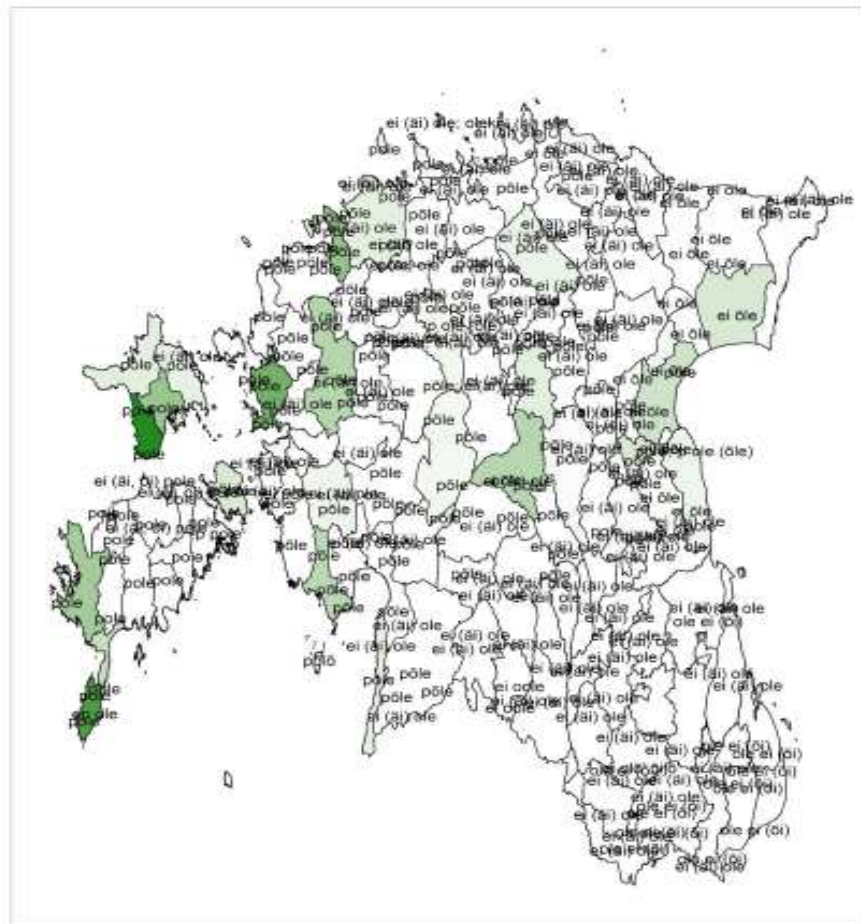
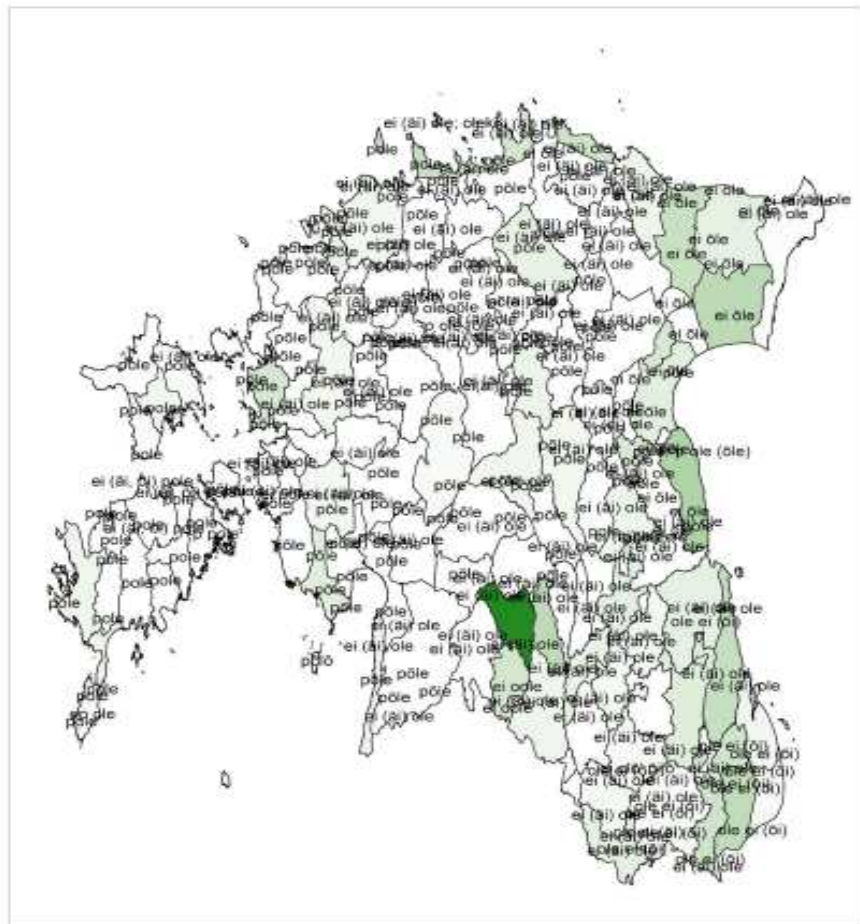
- Uiboaed, Kristel (2013): *Verbiühendid eesti murretes*. Tartu: Tartu Ülikooli Kirjastus. (= Dissertationes Philologiae Estonicae Universitatis Tartuensis).
- Vihman, Virve-Anneli, Liina Lindström, and Kristel Uiboaed (2014): "Varieerumine tarvis-/vaja-konstruksioonides keelekontaktide valguses." *Keel ja Kirjandus* 8-9, pp. 609-630.
- Ruutma, Mirjam. *Kaassõnad eesti murretes*. Diss. Tartu Ülikool, 2016.
- Uiboaed, Kristel 2016. *Spatial Visualization with R*. *spatial-visualization-with-r* 1.0. <http://dx.doi.org/10.5281/zenodo.51473>.

Corpus data

- 973 147 morphologically annotated running words
- *suvi / sui*
 - 375 occurrences from 62 parishes
 - Frequencies normalized based on a mean corpus size
- *suurem / suuremb*
 - 218 occurrences from 62 parishes
 - Frequencies normalized based on a mean corpus size
- *ei ole / pole*
 - 2969 occurrences from 53 parishes
 - Data drawn from a balanced sample of ~40 000 words/dialect (total of 405 287 words)







Shiny

- An RStudio package for creating interactive web applications
- Relatively simple syntax
- No web programming knowledge required
- Enable everything that is possible also in R, among other things, interactive mapping and statistical analysis of dialectal data
- Visualizing dialect data for the “lay person”
- Initial prototypes available at <http://rurake.keeleressursid.ee/index.php/apps/>

In conclusion

- Dialectal data analysis has always had an inherent spatial component
- Corpus data provides frequency information, linguistic atlases information on the spread of certain linguistic phenomena
- Combining the 2 data sources provides us with a more diverse understanding of dialectal variation
- R enables to perform a wide range of operations, starting from data extraction all the way to geographic visualization
- In addition to making dialect data more accessible, Shiny applications would also be a suitable platform for crowdsourcing