# A method for processing perceptual dialectology data

Laura Calaza Díaz[1]   Soraya Suárez Quintas[2]   Rosa M. Crujeiras [1]
Alberto Rodríguez Casal[1]   Xulio Sousa[2]   Jose Ramón Ríos Viqueira[3]

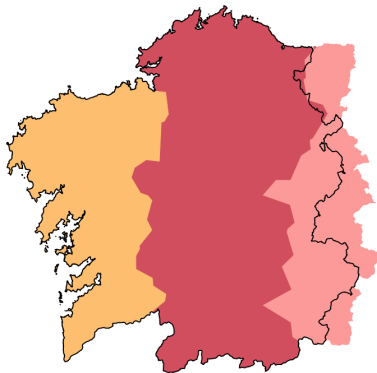[1]Departamento de Estatística e Investigación Operativa. Universidade de Santiago de Compostela.

[2]Instituto da Lingua Galega. Universidade de Santiago de Compostela

[3]COGRADE. CITIUS. Universidade de Santiago de Compostela.

*TecAnDaLi, Tecnoloxías e Análise dos Datos Lingüísticos*

# What is Perceptual Dialectology?

Academics postulate...
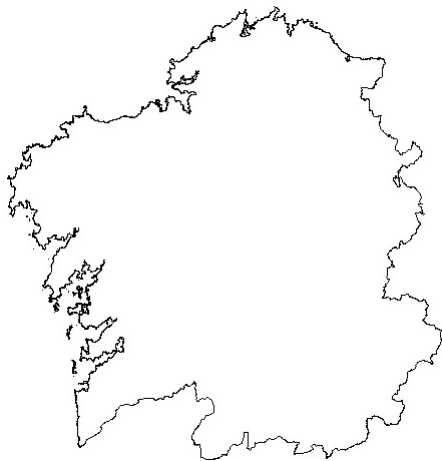


Western    Central    Eastern

Xulio C. Sousa Fernández (2006)
Análise dialectométrica das variedades xeolingüísticas galegas. *Encontro de estudos dialectolóxicos*. Actas, M C. Rola Bernardo / H. Mateus Montenegro, Ponta Delgada: Instituto Cultural de Ponta Delgada, 345-362 - Capítulo de libro
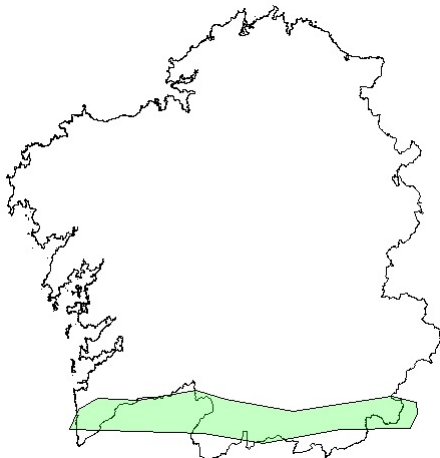
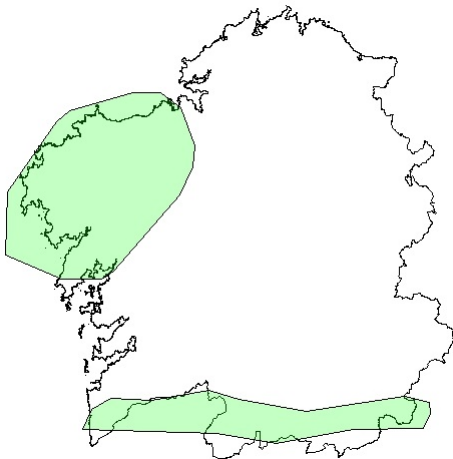# What is Perceptual Dialectology?

What I really perceive...

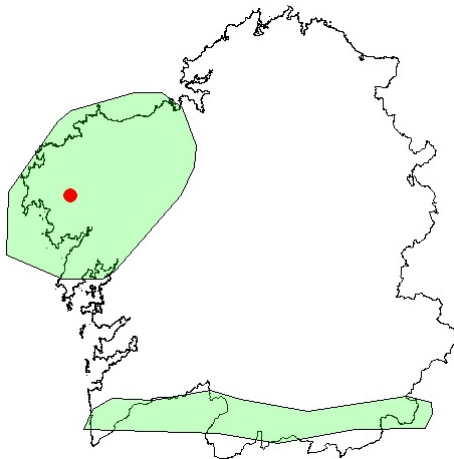# What is Perceptual Dialectology?

What I really perceive...

# What is Perceptual Dialectology?

What I really perceive...

# What is Perceptual Dialectology?

What I really perceive...

# Previous studies

Analyzed issues:

- Data collection (scanning)



- Data visualization (heat maps)

📄 Montgomery, C. and Stoeckle, P. (2013)
Geographic information systems and perceptual dialectology: a method for processing draw-a-map data. *Journal of Linguistic Geography*, 1, 52–85.

Perceptions of Galician dialects
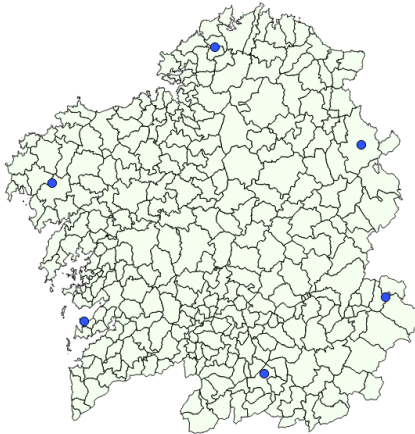
# Our purpose

Perceptions of Galician dialects

Steps...

- Collect and visualize Perceptual Dialectology data
- Discover if people are aware of and recognize regional variations of Galician language.
- Identify factors which influence geographical varieties of Galician language recognition.
- Assess in which way people's perceptions correspond to the geo-linguistic varieties traditionally recognised in Galician studies
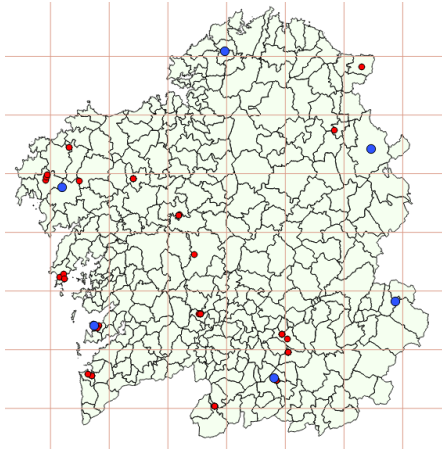
# Survey design

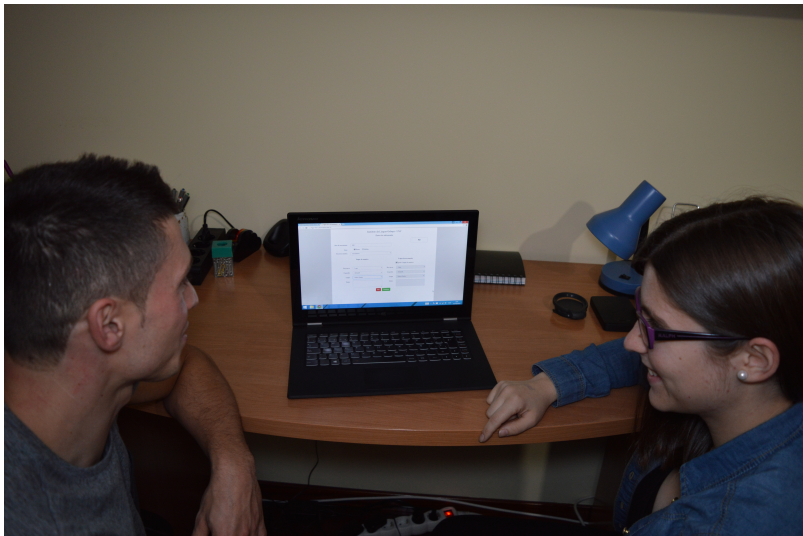Material for survey:

- 7 auditions (from different locations)

# Survey design

Material for survey:

- 50 informants (from different locations)

# Data collection

Instituto da Lingua Galega - USC

Datos do informante

| 19 |

Ano de nacemento | 1990

Sexo ● Home ○ Muller

Nivel de estudos | Universitarios ▾

**Lugar de enquisa**   **Lugar de nacemento**

☑ Igual ó lugar de enquisa

| Provincia | Lugo ▾ | Provincia | Lugo ▾ |
| Concello | A Fonsagrada ▾ | Concello | A Fonsagrada ▾ |
| Lugar | A Barreira ▾ | Lugar | A Barreira ▾ |
| Outro | | Outro | |

Sair Continuar

# Data collection

# Data collection

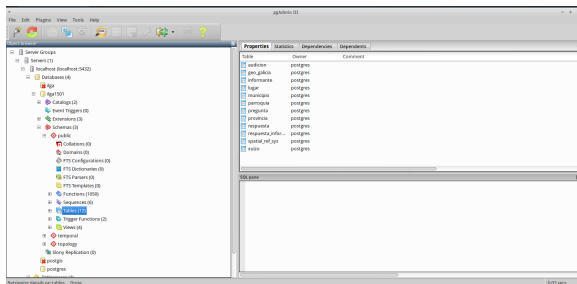# Data collection

Database summary:

- Audition locations
- Respondents information:
    - gender
    - age
    - educational level
    - birthplace
    - 7 dialects geographic perceptions
- 3 divisions of Galician map (Western, Central, Eastern)

# How to process information?

Data acquisition:

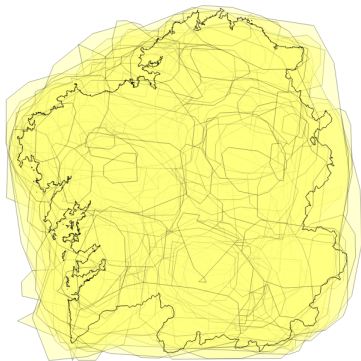PostgreSQL (Structured Query Language)



Possibilities:

- Geographic Information System (GIS) $\Rightarrow$ QGIS
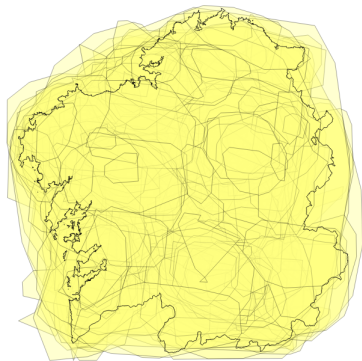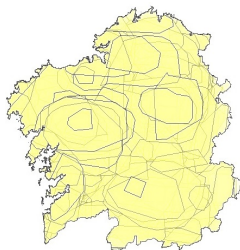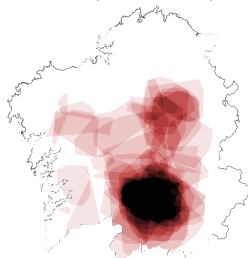- R project $\Rightarrow$ Package `RPostgreSQL`

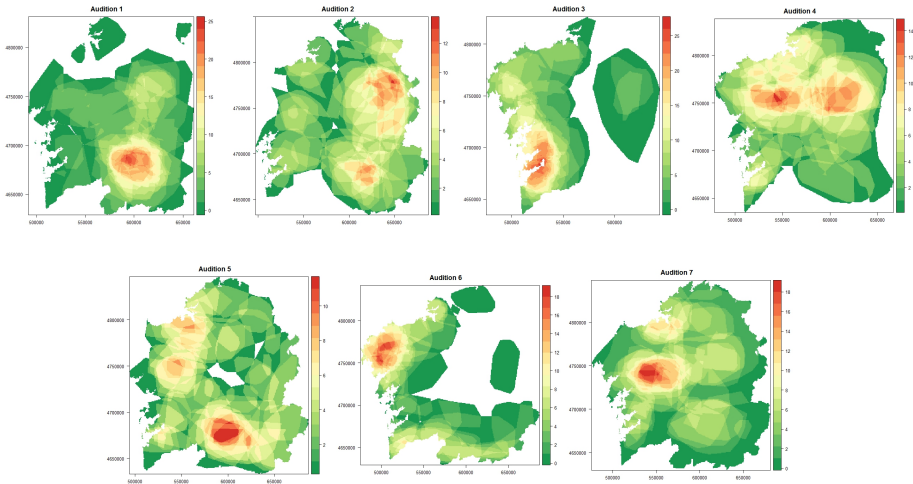Respondents' selections

Respondents' selections

Intersection of perceptions

Heat map (1$^{st}$ audition)

# Data Visualization - R (`rasterVis`)+ QGIS

# How to process information?

Handling geographical data to:

identify which factors influence the recognition of dialectal varieties

comparison with data from traditional dialectology and dialectometric studies

# How to process information?

Handling geographical data to:

identify which factors influence the recognition of dialectal varieties

comparison with data from traditional dialectology and dialectometric studies

$$\Downarrow$$

Assess diferences between sets

## Assess diferences between sets

1<sup>st</sup> case:

Assess diferences between sets

1$^{st}$ case:

Assess diferences between sets

1<sup>st</sup> case:

# How to process information?

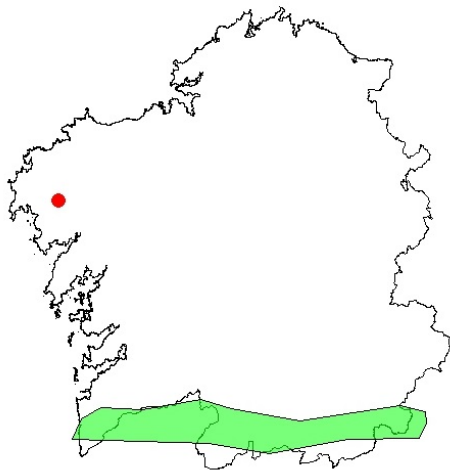## Assess diferences between sets
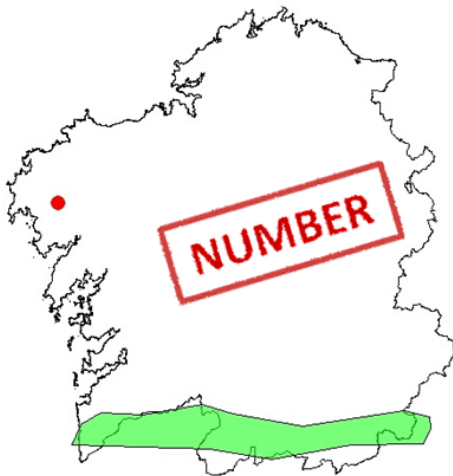
2nd case:



Western

# How to process information?

## Assess diferences between sets

2<sup>nd</sup> case:



Western

# How to process information?

## Assess diferences between sets

2nd case:



Western

# How to process information?

**Alternatives of distances between sets:**

⋆ **Centroid**

   Example (1<sup>st</sup> case)



distance=140297.9

# How to process information?

**Alternatives of distances between sets:**

⋆ **Centroid**

Example (1st case)

Example (2nd case)



distance=140297.9

distance=102477.7

# How to process information?

**Alternatives of distances between sets:**

⋆ **Centroid**

Example (1st case)

Example (2nd case)



distance=140297.9

distance=102477.7

**Alternatives of distances between sets:**

⋆ **Centroid**

⋆ **Hausdorff**

Let be the sets $A, B \subseteq X$, then

$$H(A, B) = \sup_{x \in A} |d(x, A) - d(x, B)|$$

with $d(x, A) = \inf\{\rho(x, a), a \in A\}$, and $\rho(x, y)$ the distance between any two pixels $x, y \in X$

# How to process information?

**Alternatives of distances between sets:**

- ★ **Centroid**
- ★ **Hausdorff**
- ★ **Baddeley**

$$\Delta_b(A, B) = \left[ \frac{1}{n(X)} \sum_{x \in X} |d^*(x, A) - d^*(x, B)|^p \right]^{1/p},$$

here $d^*(x, A) = \min\{d(x, A), c\} = \min\{\inf[d(x, a), a \in A], c\}$

Baddeley, A. J. (1992)
An error metric for binary images. *Robust Computer Vision: Quality of Vision Algorithms*,
W. Förstner and S. Ruwiedel (Eds.), Karlsruhe, Wichmann, pp. 59–78.

**1<sup>st</sup> objective: comparison with data from traditional dialectology and dialectometric studies**



Western

**$1^{st}$ objective: comparison with data from traditional dialectology and dialectometric studies**



Western

# Comparison with traditional dialectology

**1$^{st}$ objective: comparison with data from traditional dialectology and dialectometric studies**



- Normality tests (p.values> 0.05)

  Eastern mean: 55019.33 m
  Western mean: 54227.32 m
  Central mean: 71083.86 m

- Comparison of univariate
  density estimates:
    test equality of distributions
    p.value=0

# Regression model

**2$^{\text{nd}}$ objective: identify which factors influence in the recognition of dialectal varieties**

**For each audition:**
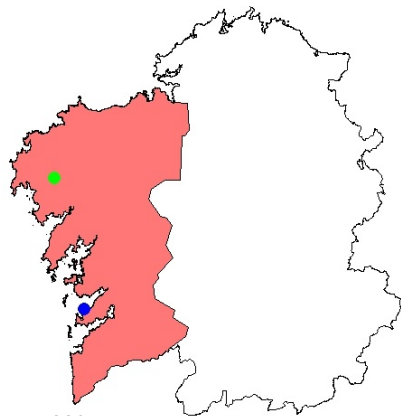
- Response: distance between the audition location and respondents' selections.
- Explanatory variables:
    - ⋆ gender
    - ⋆ age
    - ⋆ educational level
    - ⋆ distance between birthplace and audition location
    - ⋆ distance between birthplace and survey location

all models have been fitted and validated

# Regression model

**$2^{nd}$ objective: identify which factors influence in the recognition of dialectal varieties**

**For each audition:**

- Response: distance between the audition location and respondents' selections.
- Explanatory variables:
  - ⋆ gender
  - ⋆ age
  - ⋆ educational level
  - ⋆ distance between birthplace and audition location
  - ⋆ distance between birthplace and survey location

all models have been fitted and validated

# Regression model

**2<sup>nd</sup> objective: identify which factors influence in the recognition of dialectal varieties**

**To generalize...**

$$Y = X\beta + \epsilon,$$

with,

$Y = \sum_{i=1}^{6} w_i Y_{ik}, w_i = \frac{1}{n}, \forall i = 1, \ldots, 6, k = 1, \ldots, n$

$X$ includes:

- gender
- age
- educational level
- distance between survey place and birthplace

$\epsilon$ = error term

# Results

```
lm(formula = baddeleydista ~ estudos, data = basepromedio)

Residuals:
    Min      1Q  Median      3Q     Max
-21386.6 -5132.4   299.7  5512.6 23109.4

Coefficients:
                     Estimate Std. Error t value Pr(>|t|)
(Intercept)             59691       2693  22.168  < 2e-16 ***
estudosSECUNDARIOS      -5323       3517  -1.514  0.13682
estudosUNIVERSITARIOS  -11372       3376  -3.369  0.00151 **


Residual standard error: 9328 on 47 degrees of freedom
Multiple R-squared:  0.201,     Adjusted R-squared:  0.167
F-statistic: 5.911 on 2 and 47 DF,  p-value: 0.005132
```

■ **Mixed-Effects models**

To introduce random effects for: informants and auditions

■ **(Spatial) correlation between sets**

To avoid the simplification of our database