**THE USE OF VOLUNTEER GEOGRAPHIC INFORMATION FOR PRODUCING AND MAINTAINING AUTHORITATIVE LAND USE AND LAND COVER DATA**

# Automatic extraction and filtering of OpenStreetMap data to generate training datasets for LULC classification

*Cidália Costa Fonte*

*Joaquim Patriarca*

*Ismael Jesus*

*Diogo Duarte*

*Department of Mathematics – University of Coimbra, Coimbra, Portugal*

*Institute for Systems Engineering and Computers at Coimbra (INESC Coimbra), Portugal*

**24-25 November 2020**

# Summary

- Objective
- The OpenStreetMap (OSM) project and data
- Conversion of OSM data into LULC data
- Study areas
- Methodology to extract and filter OSM to train classifiers
- Results
- Conclusions

# Objective

- Test an automated methodology for generating training data from <u>OpenStreetMap</u> (OSM) to classify Sentinel-2 imagery into Land Use/Land Cover (LULC) classes

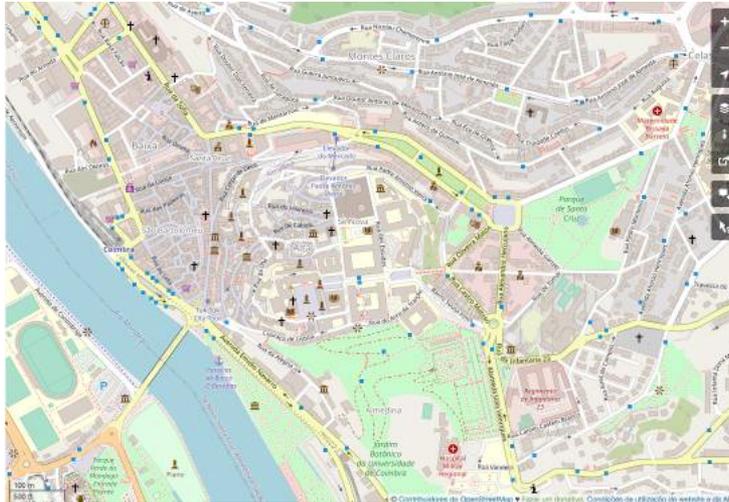- The methodology filters data extracted from OSM to generate high quality training data

# OSM project and OSM data

- OpenStreetMap (OSM) (http://www.openstreetmap.org/)
  - Project created in 2004 in the United Kingdom
  - Objective
    - Create geospatial data with open access
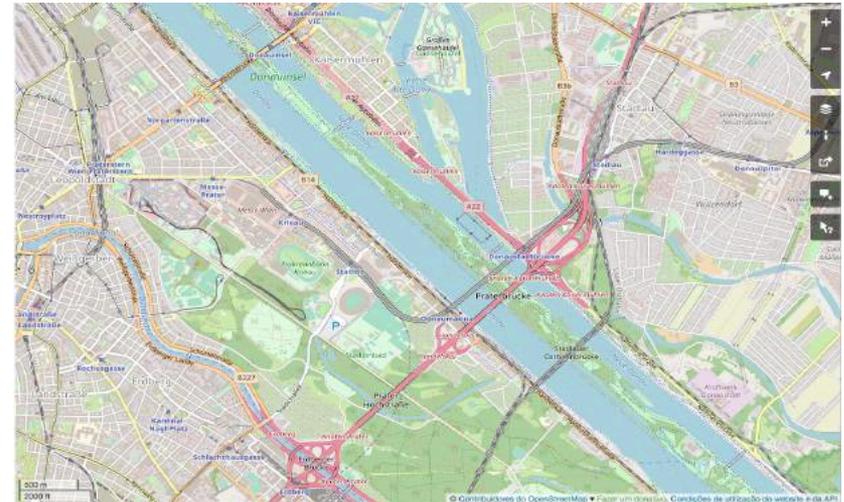


- Geospatial entities available in OSM
  - http://wiki.openstreetmap.org/wiki/Map_Features

# OSM data


OSM – Coimbra (Portugal)


OSM – Vienna (Austria)


OSM – Paris (France)


OSM – Milan (Italy)

# OSM data

- Why use OSM for training?
  - ➢ Detailed data is available
  - ➢ Includes local knowledge
  - ➢ Can be automatically downloaded

# May increase speed + lower costs

# Automate

# OSM data for training

**Main challenge**

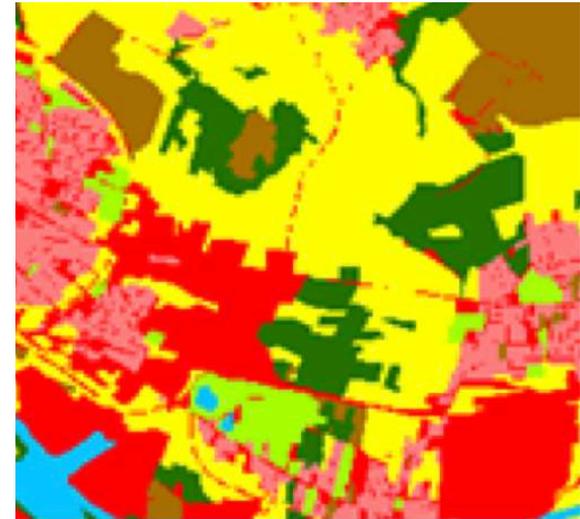Obtain high quality training data

# OSM conversion to LULC maps

- Tool to **<span style="color:red">convert automatically</span>** the data available in OSM into a LULC map

**<span style="color:red">OSM2LULC</span>**

Patriarca, J., Fonte, C.C., Estima, J., Almeida, J.P., Cardoso, A. (2019) Automatic conversion of OSM data into LULC maps: comparing FOSS4G based approaches towards an enhanced performance. Open Geospatial Data, Software and Standards., 4: 11. DOI: 10.1186/s40965-019-0070-2

Fonte, C.C., Patriarca, J., Minghini, M., Antoniou, V., See, L., Brovelli, M.A. (2017). Using OpenStreetMap to create land use and land cover maps: development of an application. In: Campelo, C.E.C., Bertolotto, M., Corcoran, P. (Eds.), *Volunteered Geographic Information and the Future of Geospatial Data*. IGI Global, Hershey, pp. 113 - 137. ISBN: 9781522524465. DOI: 10.4018/978-1-5225-2446-5.ch007

# OSM conversion to LULC maps

- **OSM2LULC**
  - Structured into 6 <u>modules</u>
  - Four <u>versions</u> available (1.2 – 1.4)
    - Using different technologies
      - GRASS GIS
      - PostGIS
      - GDAL
      - Numpy
    - With vector and raster outputs
  - Three output nomenclatures
    - Urban Atlas (UA)
    - Corine Land Cover (CLC)
    - GlobeLand 30 (GL30)
- A modified version of OSM2LULC was used
  - **OSM2LULC-4T**
    - CLC output nomenclature
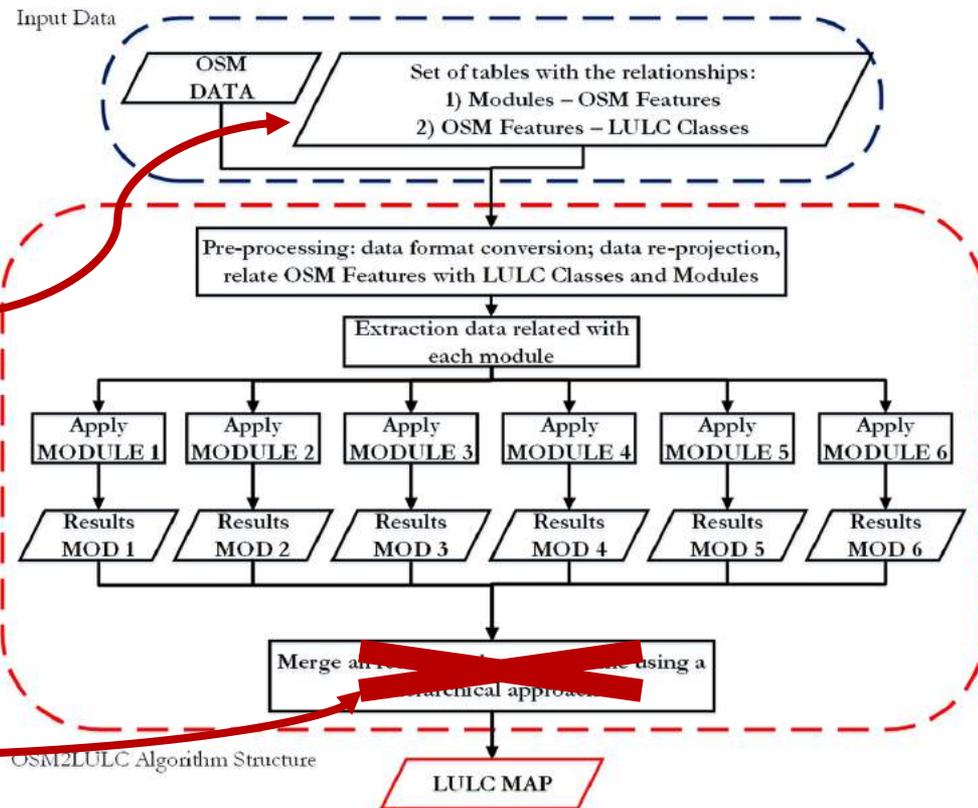
# OSM conversion to LULC maps

- **OSM2LULC-4T**
  - Version 1.2
    - Vector output



Input Data

| OSM DATA | Set of tables with the relationships: 1) Modules – OSM Features 2) OSM Features – LULC Classes |

Pre-processing: data format conversion; data re-projection, relate OSM Features with LULC Classes and Modules

Extraction data related with each module

Apply MODULE 1 — Apply MODULE 2 — Apply MODULE 3 — Apply MODULE 4 — Apply MODULE 5 — Apply MODULE 6

Results MOD 1 — Results MOD 2 — Results MOD 3 — Results MOD 4 — Results MOD 5 — Results MOD 6

Merge an... ...using a ...archical approa...

OSM2LULC Algorithm Structure

LULC MAP

Association OSM features - LULC class

Some modifications

# Methodology to extract and filter OSM to train classifiers

- Concept:
  - The data extracted from OSM is filtered
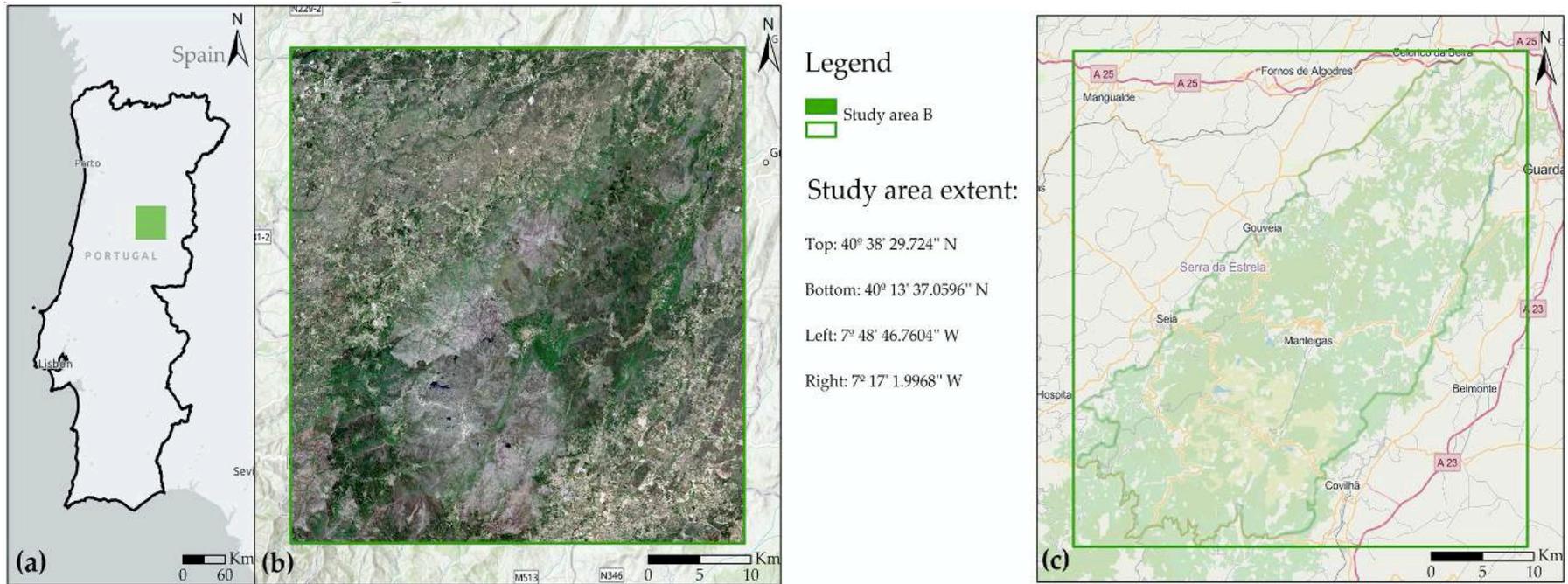    - To remove regions with higher uncertainty or that changed

## Based on:

**Uncertainty due to possible overlaping classes in the pixels**

**Radiometric indices of the images to classify**

# Case studies

- Case study A

# Case studies

■ ## Case study B

# Methodology to extract and filter OSM to train classifiers

- Used data:
  - OSM for training
  - Sentinel-2 multiespectral images
    - Only bands B2, B3, B4 and B8
    - Three images were used for each study area

|  | Satellite | Product Type | Collection date | Sentinel GRID |
|---|---|---|---|---|
| **Study area A** | Sentinel-2A | Level-2A | 2018-03-21 | T29SMC |
|  | Sentinel-2A | Level-2A | 2018-06-19 | T29SMC |
|  | Sentinel-2B | Level-2A | 2018-10-22 | T29SMC |
| **Study area B** | Sentinel-2B | Level-2A | 2018-03-26 | T29TPE |
|  | Sentinel-2A | Level-2A | 2018-06-19 | T29TPE |
|  | Sentinel-2B | Level-2A | 2018-10-22 | T29TPE |

  - Portuguese Land Cover Map (COS 2018)
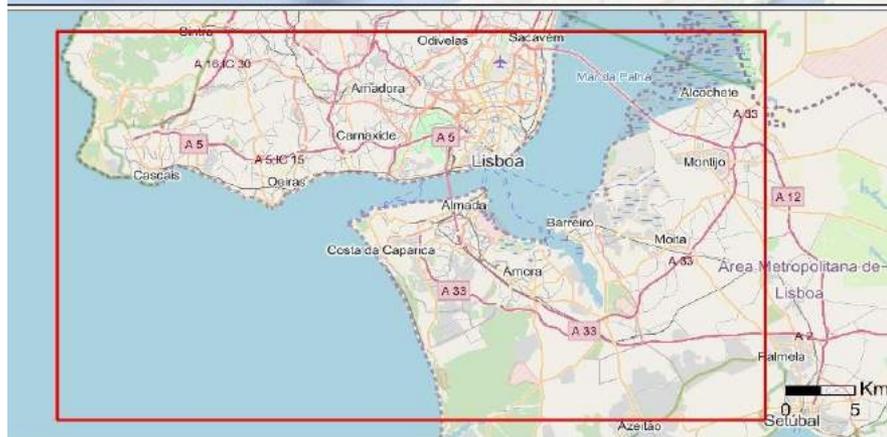    - With nomenclature harmonization
    - As reference data

# Methodology to extract and filter OSM to train classifiers

- Nomenclature harmonization

| Used classes | OSM2LULC | COS 2018 |
|---|---|---|
| 1. Artificial surfaces | 1.1 Urban fabric<br>1.2 Industrial, commercial and transport units<br>1.3 Mine, dump and construction sites<br>1.4.2 Sport and leisure facilities (excluding golf courses) | 1 Artificial surfaces, excluding:<br>- Golf courses (1.6.1.1)<br>- Public gardens and playground (1.7.1.1) |
| 2. Agricultural areas | 2.1 Arable land<br>2.2 Permanent crops<br>2.4 Heterogeneous agricultural areas | 2 Agriculture |
| 3. Herbaceous vegetation | 1.4.1 Green urban areas<br>2.3 Pastures<br>3.2.1 Natural grasslands<br>1.4.2 Sport and leisure facilities (only golf courses) | 3 Herbaceous<br>1.6.1.1 Golf courses<br>1.7.1.1 public gardens and playgrounds |
| 4. Forest areas | 3.1 Forests | 4 Agroforestry<br>5 Forestry |
| 5. Shrublands | 3.2.4 Transitional woodland-shrub | 6 Shrublands |
| 6. Open spaces with little or no vegetation | 3.3 Open spaces with little or no vegetation | 7 Open spaces with little or no vegetation |
| 7. Wetlands | 4 Wetlands | 8 Wetlands |
| 8. Water bodies | 5.1 Inland waters<br>5.2 Marine waters | 9 Water bodies |

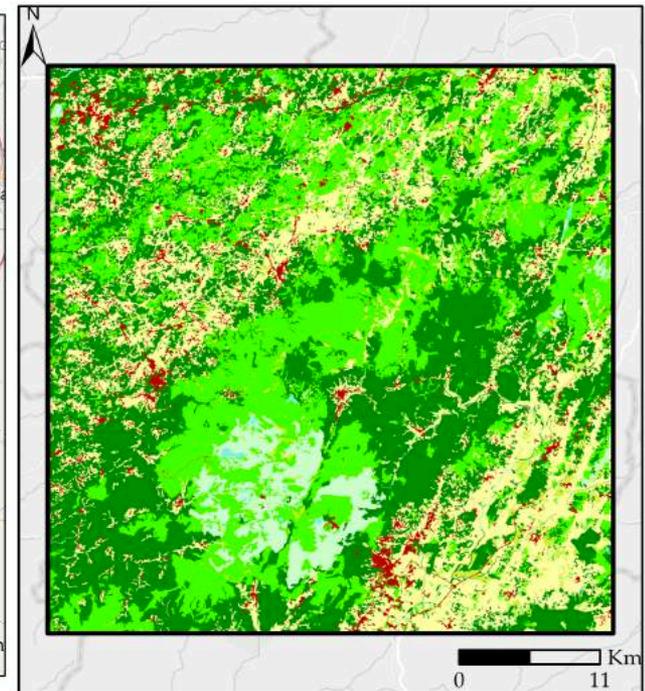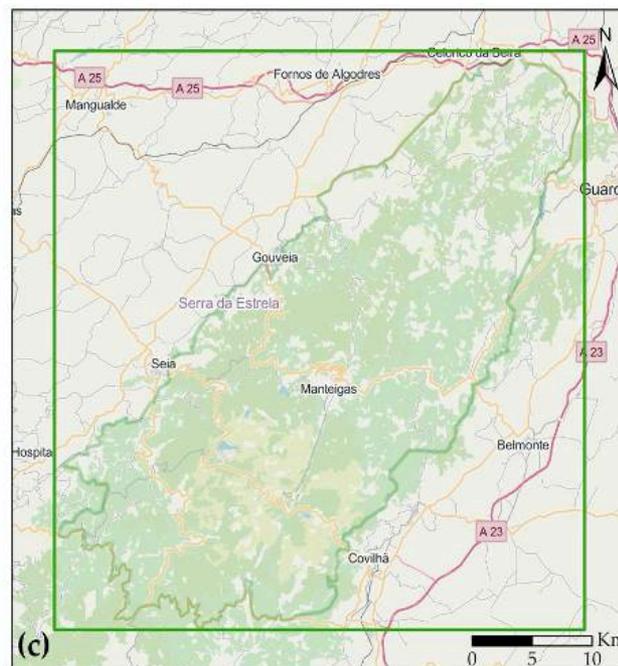# Case studies

- ## Case study A

COS derived reference data



LULC classes

- Artificial surfaces
- Agricultural areas
- Herbaceous vegetation
- Forest areas
- Shrublands

- Study area A
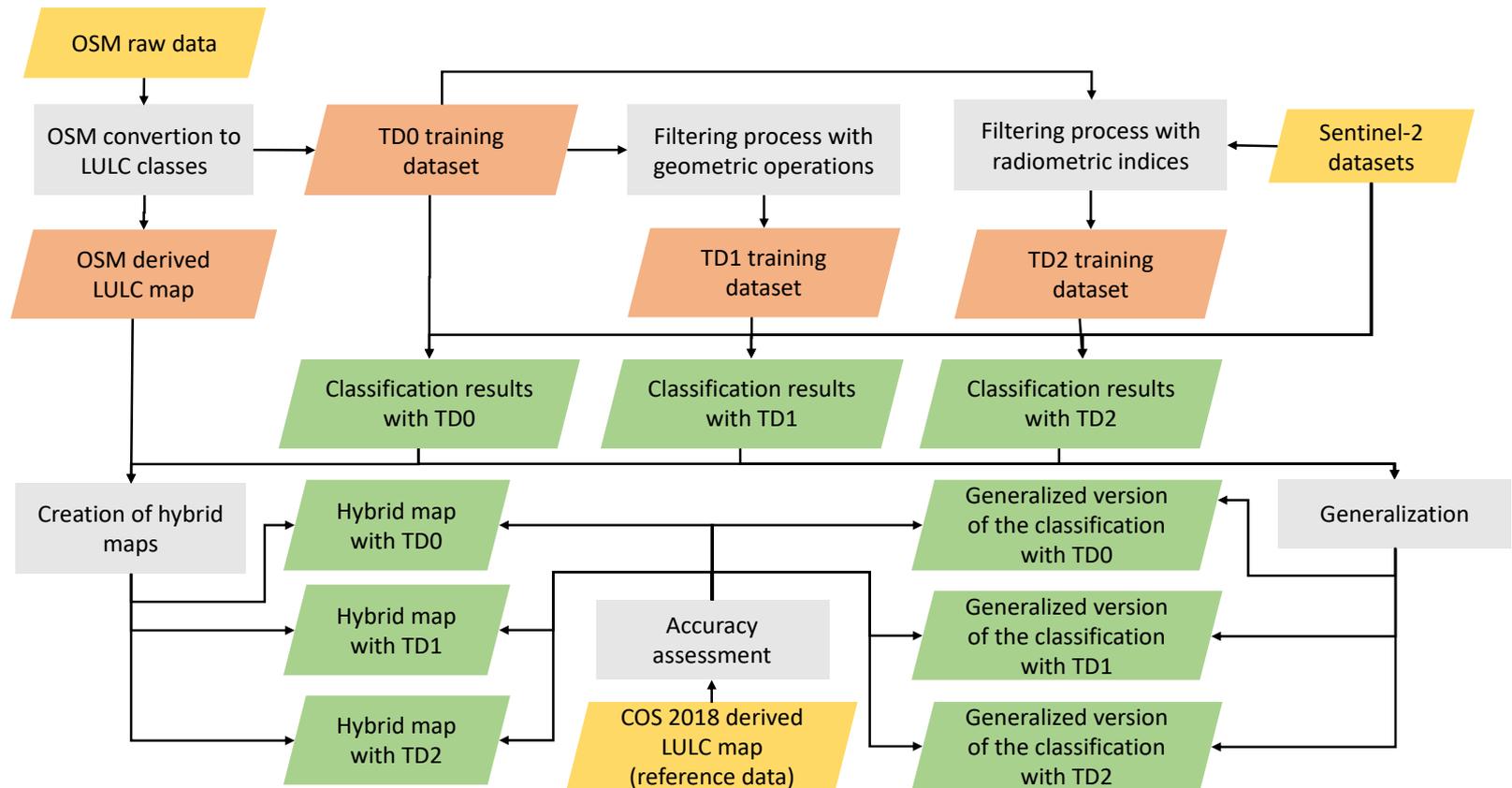- Wetlands
- Water bodies
- Open spaces with little or no vegetation

# Case studies

- ## Case study B

COS derived reference data



LULC classes

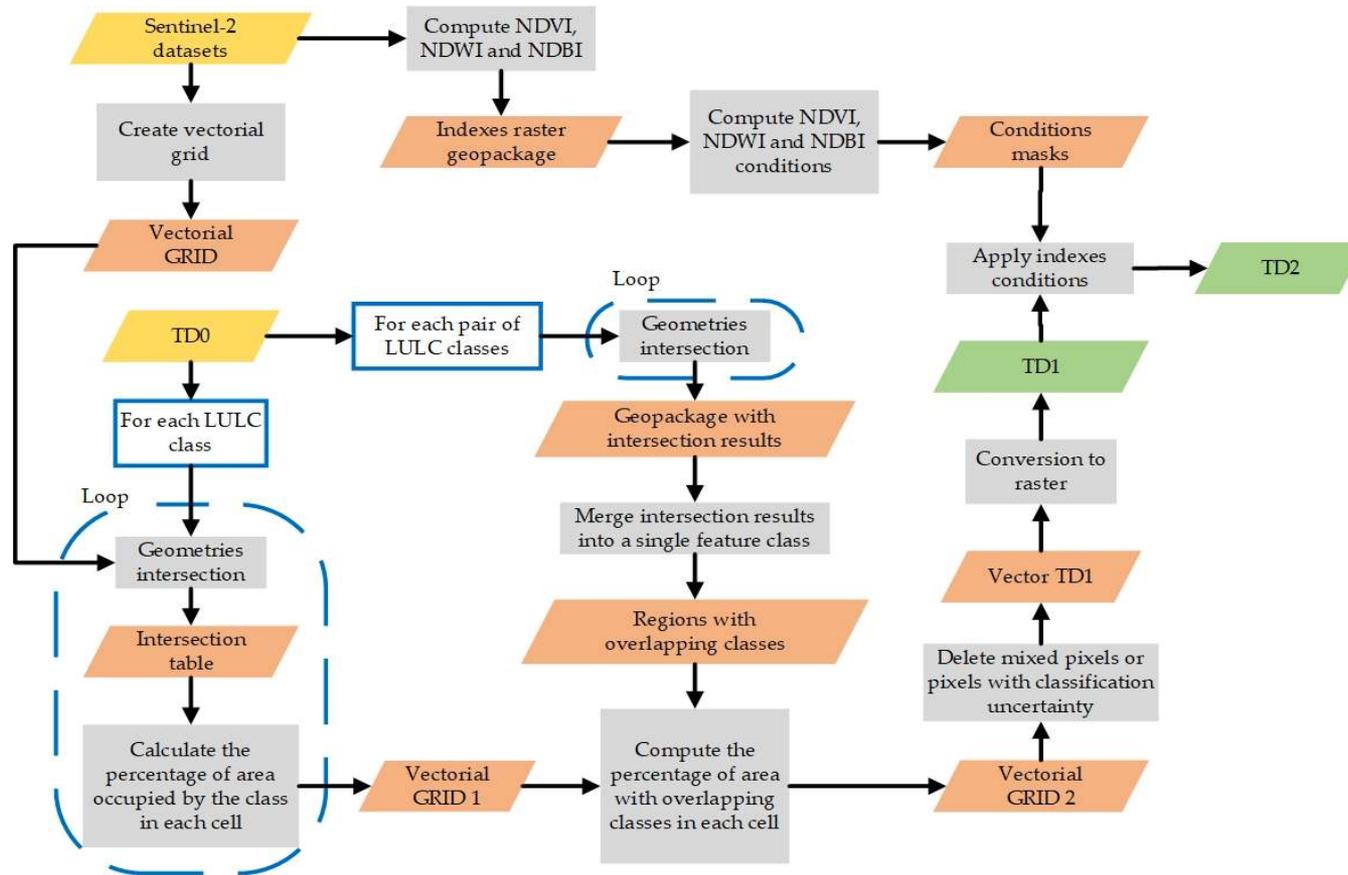| | Study area A |
| :-- | :-- |
| ■ Artificial surfaces | ■ Forest areas | ■ Wetlands |
| ■ Agricultural areas | ■ Shrublands | ■ Water bodies |
| ■ Herbaceous vegetation | ■ Open spaces with little or no vegetation |

# Methodology to extract and filter OSM to train classifiers

■ Workflow

# Methodology to extract and filter OSM to train classifiers

- Filtering procedures

# Methodology to extract and filter OSM to train classifiers

- Training data
  - Classes' separability was computed with the Bhattacharyya distance
- Classification
  - Samples of the training sets were used
    - Due to computational constraints
    - Sample size proportional to class area
  - Random forest classified
- Generalization
  - A majority filter - circular moving window with 5 cells - was applied to remove isolated small regions
- Hybrid maps
  - Data from OSM2LULC + classification results
- Accuracy assessment

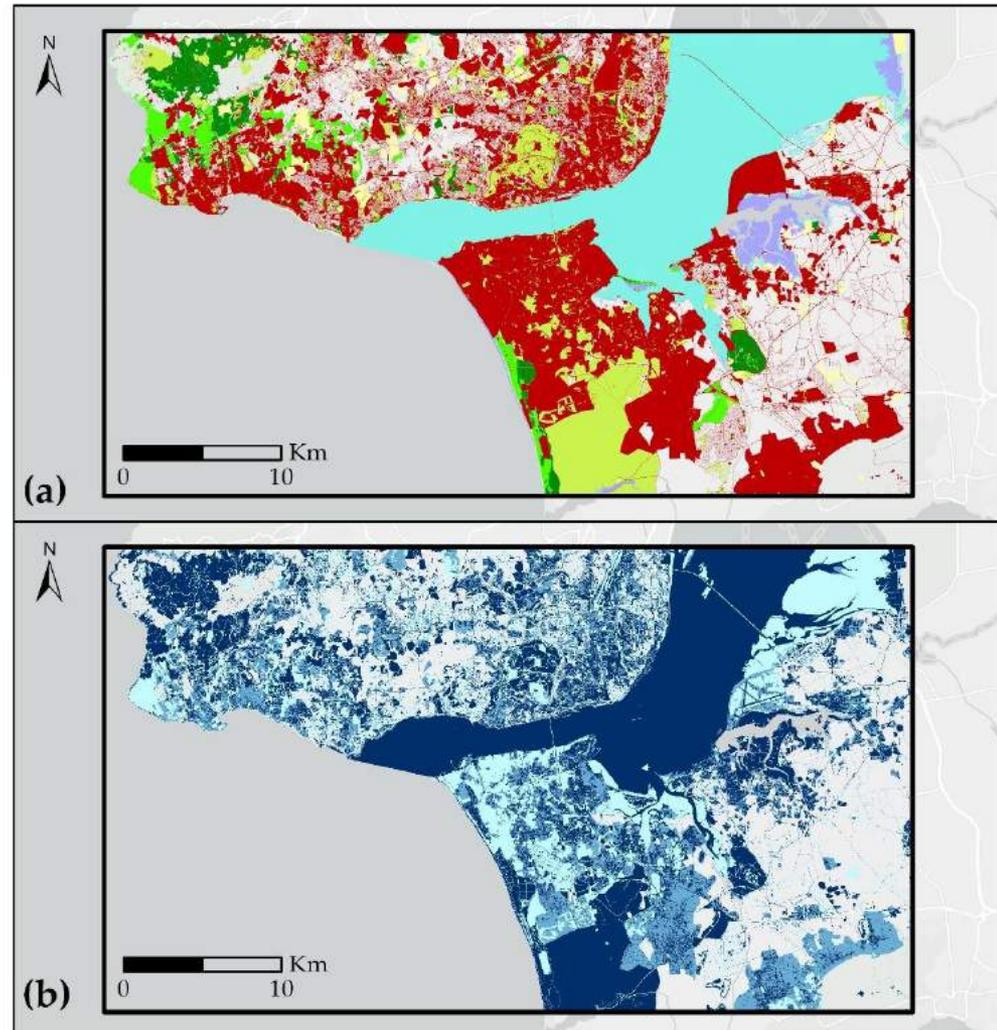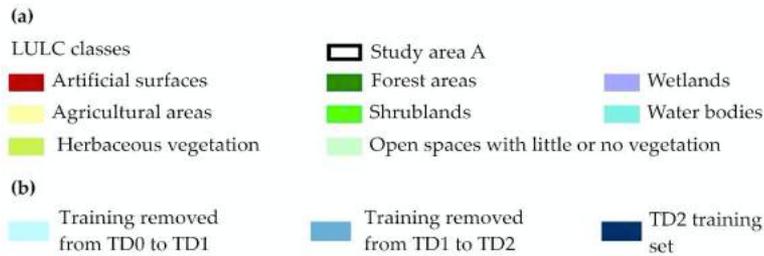# Methodology to extract and filter OSM to train classifiers

- ## Radiometric indices
  - ➤ Conditions used for all the images

| Classes | NDVI / images | NDWI / images | NDBI / images |
|---|---|---|---|
| 1. Artificial surfaces | < 0.3 / all | < 0.0 / all | > 0.0/at least one |
| 2. Agricultural areas | > 0.3 / all | < 0.0 / all | --- |
| 3. Herbaceous vegetation | > 0.3 / all | < 0.0 / all | --- |
| 4. Forest areas | > 0.3 / all | < 0.0 / all | --- |
| 5. Shrublands | > 0.3 / all | < 0.0 / all | --- |
| 6. Open spaces with little or no vegetation | > 0.0/at least one | < 0.0/at least one | --- |
| 7. Wetlands | > 0.0/at least one | < 0.0/at least one | --- |
| 8. Water bodies | < 0.3/at least one | > 0.0 / all | --- |

# Results

- Training data obtained
  - Study area A

(a)

LULC classes

| | | |
|---|---|---|
| ■ Artificial surfaces | ■ Forest areas | ■ Wetlands |
| Agricultural areas | ■ Shrublands | ■ Water bodies |
| Herbaceous vegetation | Open spaces with little or no vegetation | |

☐ Study area A

(b)

| | | |
|---|---|---|
| Training removed from TD0 to TD1 | Training removed from TD1 to TD2 | ■ TD2 training set |

# Results

■ ## Training data obtained

➢ ### Study area B



(a)

| LULC classes | | |
|---|---|---|
| Artificial surfaces (red) | Forest areas (dark green) | Study area A |
| Agricultural areas (yellow) | Shrublands (green) | Wetlands (purple) |
| Herbaceous vegetation (light green) | Open spaces with little or no vegetation | Water bodies (cyan) |

(b)
- Training removed from TD0 to TD1
- Training removed from TD1 to TD2
- TD2 training set

# Results

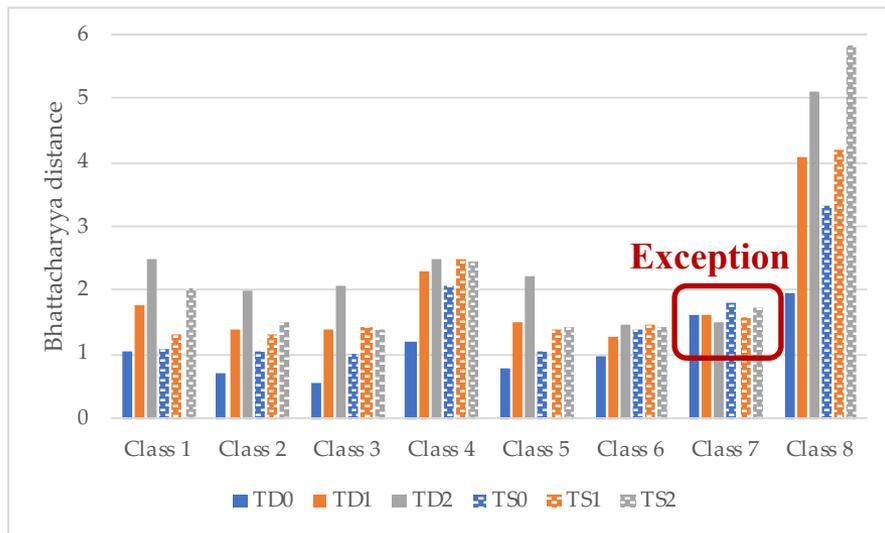■ Percentage of the TD0, TD1 and TD2 datasets belonging to each class for study areas A and B

Most regions excluded from TD0 belong to class 1

| Classes | Study area A | | | Study area B | | |
|---|---|---|---|---|---|---|
| | TD0 | TD1 | TD2 | TD0 | TD1 | TD2 |
| 1. Artificial surfaces | 50.3 | 44.6 | 27.0 | 6.9 | 4.4 | 1.5 |
| 2. Agricultural areas | 2.7 | 2.7 | 3.6 | 12.8 | 11.9 | 13.4 |
| 3. Herbaceous vegetation | 9.6 | 11.5 | 15.4 | 2.3 | 2.0 | 2.0 |
| 4. Forest areas | 4.8 | 5.3 | 7.1 | 36.6 | 37.8 | 42.1 |
| 5. Shrublands | 3.7 | 4.4 | 5.9 | 40.4 | 43.1 | 40.5 |
| 6. Open spaces with little or no vegetation | 0.5 | 0.4 | 0.5 | 0.4 | 0.4 | 0.4 |
| 7. Wetlands | 7.4 | 3.1 | 4.0 | 0.001 | --- | --- |
| 8. Water bodies | 20.9 | 28.0 | 36.4 | 0.8 | 0.4 | 0.2 |

# Results

- **Class separability (the great the better)**
  - ➢ Classes' separability improves with the filtering for most classes

Study area A



Study area B

# Results

# Classification results
## Study area A

TD0

TD2



Artificial surfaces
Agricultural areas
Herbaceous vegetation
Forest areas
Shrublands
Open spaces with little or no vegetation
Wetlands
Water bodies

TD1

Reference data

# Results

## Classification results
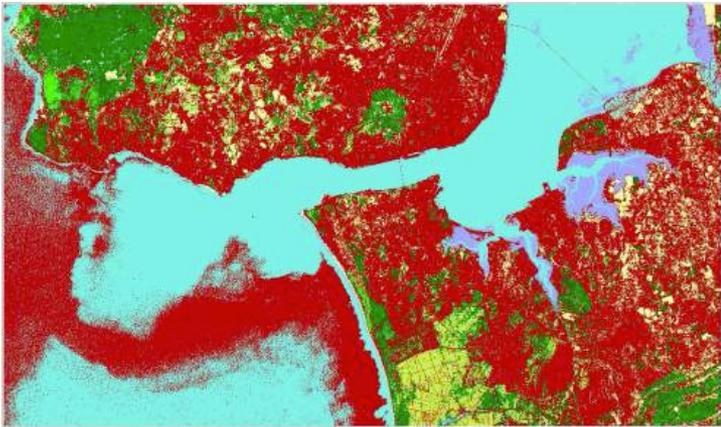### Study area B



TD0

TD2

TD1

Reference data

Legend:
- Artificial surfaces
- Agricultural areas
- Herbaceous vegetation
- Forest areas
- Shrublands
- Open spaces with little or no vegetation
- Wetlands
- Water bodies

# Results

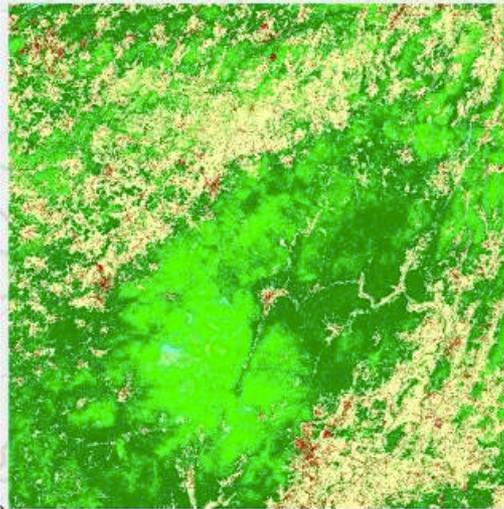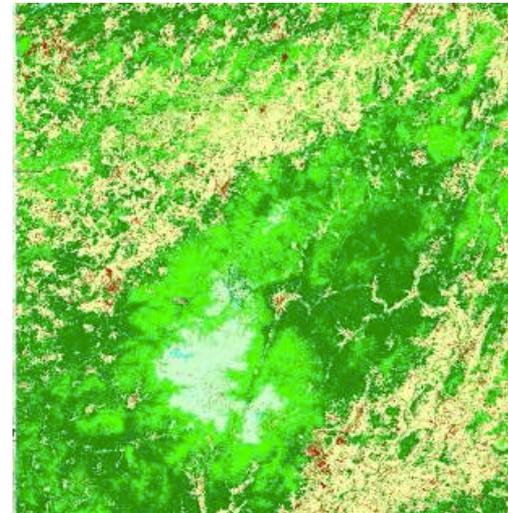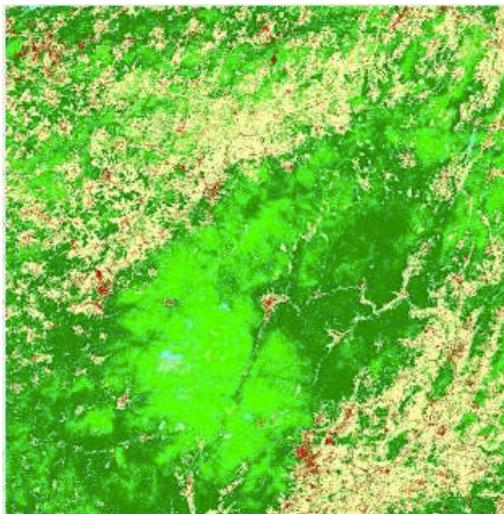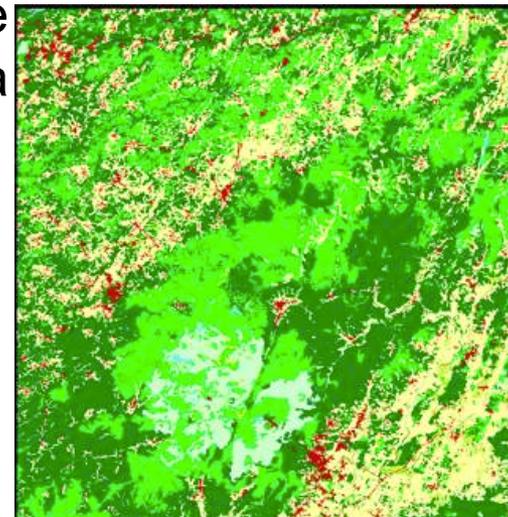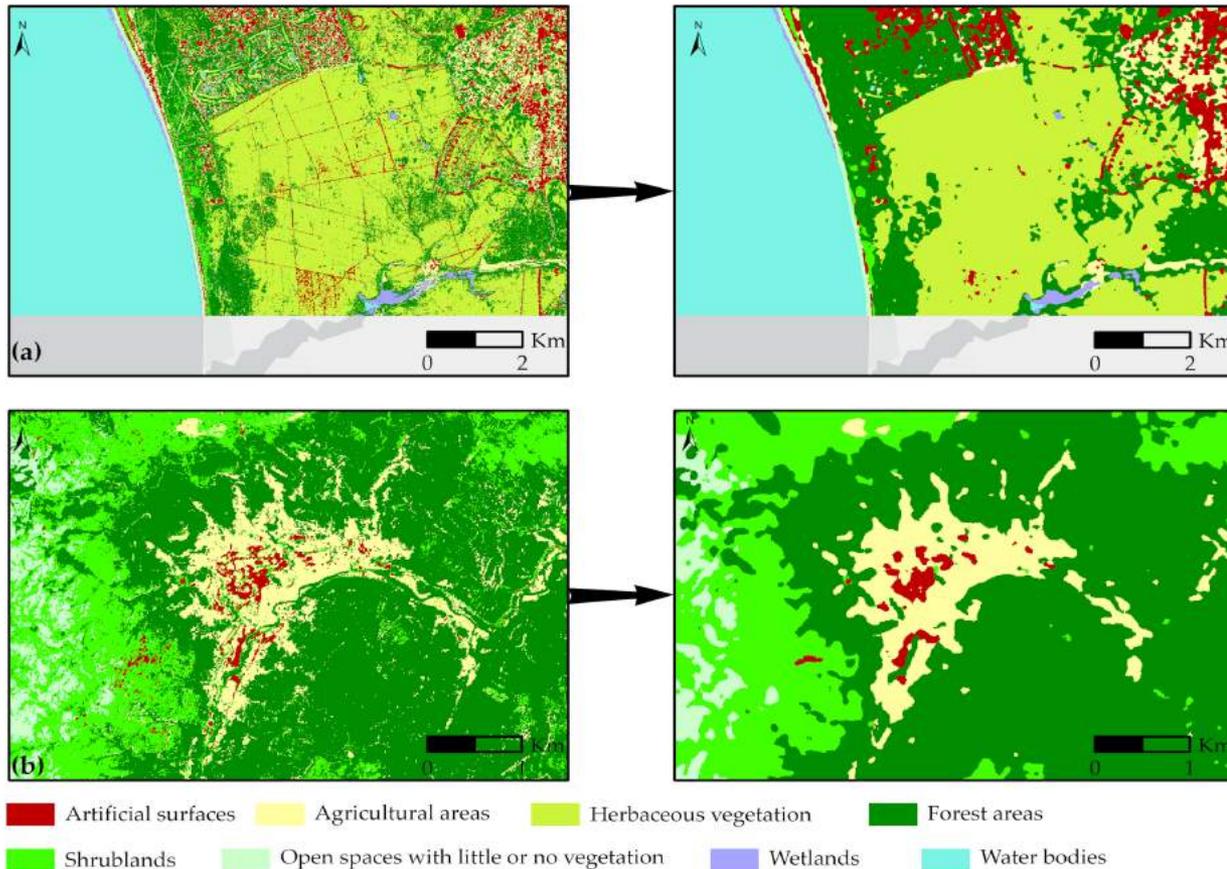- ■ Effect of the generalization step



Artificial surfaces    Agricultural areas    Herbaceous vegetation    Forest areas
Shrublands    Open spaces with little or no vegetation    Wetlands    Water bodies

# Results

- ## Overall Accuracy
  - ➤ Increases with the filtering for
    - the training datasets
    - Classification for study area A

| Dataset | Study area A | | | Study area B | | |
| --- | --- | --- | --- | --- | --- | --- |
| | TD0 | TD1 | TD2 | TD0 | TD1 | TD2 |
| Training datasets | 64 | 74 | 76 | 87 | 89 | 93 |
| Classification results | 55 | 64 | 73 | 65 | 65 | 65 |
| Generalized maps | 55 | 64 | 78 | 69 | 69 | 69 |
| Classification only for regions with OSM data | 69 | 73 | 66 | 66 | 66 | 66 |
| Data obtained with OSM2LULC | 70 | | | 87 | | |

# Results –Study area A

| Classes | USER'S ACCURACY | | | | | PRODUCER'S ACCURACY | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| **Training datasets** | | | | | | | | | | |
| | TD0 | TD1 | TD2 | TD1-TD0 | TD2-TD1 | TD0 | TD1 | TD2 | TD1-TD0 | TD2-TD1 |
| 1. Artificial surfaces | 71 | 81 | 97 | 10 | 16 | 66 | 97 | 95 | 31 | -2 |
| 2. Agricultural areas | 62 | 72 | 72 | 10 | 0 | 73 | 33 | 63 | -40 | 30 |
| 3. Herbaceous vegetation | 12 | 10 | 10 | -3 | 0 | 11 | 42 | 59 | 31 | 17 |
| 4. Forest areas | 75 | 84 | 84 | 9 | 0 | 56 | 24 | 30 | -32 | 6 |
| 5. Shrublands | 38 | 40 | 40 | 3 | 0 | 21 | 41 | 52 | 20 | 11 |
| 6. Open spaces with little or no vegetation | 42 | 47 | 48 | 5 | 0 | 49 | 45 | 48 | -4 | 3 |
| 7. Wetlands | 16 | 34 | 34 | 18 | 0 | 67 | 61 | 78 | -6 | 17 |
| 8. Water bodies | 94 | 97 | 99 | 3 | 1 | 85 | 93 | 93 | 8 | 0 |
| **Classification** | | | | | | | | | | |
| | TS0 | TS1 | TS2 | TS1-TS0 | TS2-TS1 | TS0 | TS1 | TS2 | TS1-TS0 | TS2-TS1 |
| 1. Artificial surfaces | 41 | 48 | 88 | 7 | 40 | 93 | 92 | 62 | -1 | -30 |
| 2. Agricultural areas | 53 | 53 | 42 | 0 | -11 | 40 | 40 | 74 | 0 | 34 |
| 3. Herbaceous vegetation | 4 | 5 | 6 | 1 | 1 | 3 | 5 | 7 | 2 | 2 |
| 4. Forest areas | 72 | 73 | 63 | 1 | -10 | 44 | 47 | 58 | 3 | 11 |
| 5. Shrublands | 41 | 38 | 26 | -3 | -12 | 9 | 14 | 18 | 5 | 4 |
| 6. Open spaces with little or no vegetation | 22 | 25 | 6 | 3 | -19 | 46 | 44 | 46 | -2 | 2 |
| 7. Wetlands | 15 | 28 | 25 | 13 | -3 | 42 | 51 | 60 | 9 | 9 |
| 8. Water bodies | 97 | 99 | 99 | 2 | 0 | 50 | 68 | 95 | 18 | 27 |
| **Generalization** | | | | | | | | | | |
| | TS0 | TS1 | TS2 | TS1-TS0 | TS2-TS1 | TS0 | TS1 | TS2 | TS1-TS0 | TS2-TS1 |
| 1. Artificial surfaces | 41 | 47 | 89 | 6 | 42 | 97 | 97 | 75 | 0 | -22 |
| 2. Agricultural areas | 61 | 60 | 45 | -1 | -15 | 36 | 36 | 85 | 0 | 49 |
| 3. Herbaceous vegetation | 2 | 4 | 6 | 2 | 2 | 2 | 4 | 6 | 2 | 2 |
| 4. Forest areas | 75 | 77 | 69 | 2 | -8 | 44 | 48 | 64 | 4 | 16 |
| 5. Shrublands | 55 | 53 | 42 | -2 | -11 | 6 | 9 | 12 | 3 | 3 |
| 6. Open spaces with little or no vegetation | 36 | 45 | 32 | 9 | -13 | 47 | 44 | 47 | -3 | 3 |
| 7. Wetlands | 16 | 30 | 29 | 14 | -1 | 45 | 53 | 66 | 8 | 13 |
| 8. Water bodies | 97 | 99 | 99 | 2 | 0 | 48 | 67 | 95 | 19 | 28 |

# Results –Study area A

- Examples



Reference　　　Classification　　　Generatization

# Results –Study area B

## Training datasets

| Classes | User's accuracy | | | | | Producer's accuracy | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | TD0 | TD1 | TD2 | TD1-TD0 | TD2-TD1 | TD0 | TD1 | TD2 | TD1-TD0 | TD2-TD1 |
| 1. Artificial surfaces | 48 | 68 | 90 | 20 | 22 | 96 | 98 | 96 | 2 | -2 |
| 2. Agricultural areas | 99 | 99 | 99 | 0 | 0 | 86 | 93 | 99 | 8 | 6 |
| 3. Herbaceous vegetation | 71 | 72 | 69 | 2 | -4 | 89 | 96 | 99 | 7 | 3 |
| 4. Forest areas | 100 | 100 | 100 | 0 | 0 | 94 | 97 | 98 | 4 | 1 |
| 5. Shrublands | 80 | 80 | 86 | 0 | 6 | 98 | 99 | 100 | 2 | 0 |
| 6. Open spaces with little or no vegetation | 97 | 98 | 98 | 1 | 0 | 4 | 4 | 7 | 0 | 3 |
| 7. Wetlands | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| 8. Water bodies | 51 | 93 | 99 | 42 | 6 | 93 | 97 | 97 | 4 | 0 |

## Classification

| Classes | TS0 | TS1 | TS2 | TS1-TS0 | TS2-TS1 | TS0 | TS1 | TS2 | TS1-TS0 | TS2-TS1 |
|---|---|---|---|---|---|---|---|---|---|---|
| 1. Artificial surfaces | 52 | 48 | 58 | -4 | 9 | 47 | 53 | 39 | 6 | -14 |
| 2. Agricultural areas | 63 | 63 | 63 | 1 | 0 | 85 | 84 | 85 | -1 | 1 |
| 3. Herbaceous vegetation | 46 | 47 | 42 | 1 | -6 | 5 | 6 | 4 | 1 | -2 |
| 4. Forest areas | 71 | 72 | 74 | 0 | 2 | 73 | 73 | 70 | 0 | -3 |
| 5. Shrublands | 60 | 60 | 59 | 0 | -1 | 54 | 55 | 54 | 1 | 0 |
| 6. Open spaces with little or no vegetation | 54 | 53 | 41 | -1 | -12 | 8 | 8 | 38 | 0 | 30 |
| 7. Wetlands | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| 8. Water bodies | 93 | 63 | 88 | -29 | 24 | 38 | 64 | 47 | 26 | -17 |

## Generalization

| Classes | TS0 | TS1 | TS2 | TS1-TS0 | TS2-TS1 | TS0 | TS1 | TS2 | TS1-TS0 | TS2-TS1 |
|---|---|---|---|---|---|---|---|---|---|---|
| 1. Artificial surfaces | 77 | 72 | 77 | -5 | 5 | 47 | 55 | 37 | 8 | -18 |
| 2. Agricultural areas | 63 | 64 | 63 | 1 | -1 | 91 | 90 | 91 | -1 | 1 |
| 3. Herbaceous vegetation | 74 | 75 | 56 | 1 | -19 | 2 | 3 | 1 | 1 | -2 |
| 4. Forest areas | 75 | 75 | 78 | 0 | 3 | 78 | 77 | 74 | -1 | -3 |
| 5. Shrublands | 65 | 65 | 65 | 0 | 0 | 56 | 57 | 56 | 1 | -1 |
| 6. Open spaces with little or no vegetation | 78 | 82 | 42 | 4 | -40 | 4 | 4 | 38 | 0 | 34 |
| 7. Wetlands | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| 8. Water bodies | 94 | 78 | 90 | -16 | 12 | 40 | 58 | 47 | 18 | -11 |

# Results

- ## Hybrid maps

  - ### Overall accuracy (%)

| | Class / Gen | | | Hybrid map (HM) | | | HM – Class / HM - Gen | | |
|---|---|---|---|---|---|---|---|---|---|
| | TS0 | TS1 | TS2 | TS0 | TS1 | TS2 | TS0 | TS1 | TS2 |
| **Study area A** | 55 / 55 | 64 / 64 | 73 / 78 | 56 | 62 | 76 | 1 / 1 | -2 / -2 | 3 / -2 |
| **Study area B** | 65 / 69 | 65 / 69 | 65 / 69 | 75 | 75 | 74 | 10 / 10 | 10 / 10 | 9 / 9 |

  - ### The quality of these hybrid products is very much dependent on the characteristics of the region and the data available in OSM

# Conclusions

- In general the filtering processes <u>improved the class separability</u>
  - ➤ Problematic regions were successfully removed from the training datasets
- The accuracy of the classification results and their generalized versions <u>may improve</u> with the filtering
  - ➤ increased for the study area with urban characteristics
  - ➤ remained unchanged for the rural study area

# Conclusions

- Some classes were <span style="color:red">very hard to classify</span>
  - Worse classes:
    - Agricultural areas
    - Herbaceous vegetation
    - Shrublands
    - Open spaces with little or no vegetation

  - The nomenclature also included land use classes

# Conclusions

- An accuracy of up to 78% was achieved with an automated procedure

  ➢ Study area A and training data TD2

- The use of reference data with 1 ha MMU raises problems

Additional tests are under development

Using all Sentinel-2 bands

Without using samples for training

Different approach for accuracy assessment

**THE USE OF VOLUNTEER GEOGRAPHIC INFORMATION FOR PRODUCING AND MAINTAINING AUTHORITATIVE LAND USE AND LAND COVER DATA**

# Quality assessment of Land Use and Land Cover information with VGI

*Cidália Costa Fonte*

*Department of Mathematics – University of Coimbra, Coimbra, Portugal*
*Institute for Systems Engineering and Computers at Coimbra (INESC Coimbra), Portugal*

**24-25 November 2020**

# Quality assessment

Requires the use of reference data

"ground truth"

Difficult to collect

Expensive

Time consuming

Requires human intervention

Hardly automated

# Quality assessment

- **The reference data should be formed by a sample**
  - Probability sample
    - All spatial units have the same probability of being selected

Will VGI be available to provide the reference condition at all locations?

If volunteers can be directed to the required locations

Maybe yes!

If existing VGI will be used

Probably not!

# Types of VGI that may be used for LULC map validation

- **Photographs and descriptions**
  - ➢ Degree Confluence project
  - ➢ Geograph
  - ➢ Panoramio
  - ➢ Flickr
- Volunteer initiatives to **map the world**
  - ➢ Such as OpenStreetMap (OSM)
- **Land cover data** collected by projects
  - ➢ such as Geo-Wiki and VIEW-IT

# Quality assessment

## What about the quality of that VGI?

> **Strategies need to be used so that**
>
> **ONLY HIGHLY RELIABLE VGI**
>
> **is used for LULC map validation!**

**VGI quality assessment**

Let us assume only highly reliable VGI is used

# Can VGI be useful for quality assessment?

- ## Possible strategies:

Using volunteered geographic information (VGI) in design-based statistical inference for area estimation and accuracy assessment of land cover

Stephen V. Stehman[a,*], Cidália C. Fonte[b], Giles M. Foody[c], Linda See[d]

[a] Department of Forest and Natural Resources Management, SUNY College of Environmental Science and Forestry, Syracuse, NY 13210, United States
[b] Department of Mathematics, University of Coimbra, P-3001 501 Coimbra, Portugal/Institute for Systems Engineering and Computers at Coimbra, Portugal
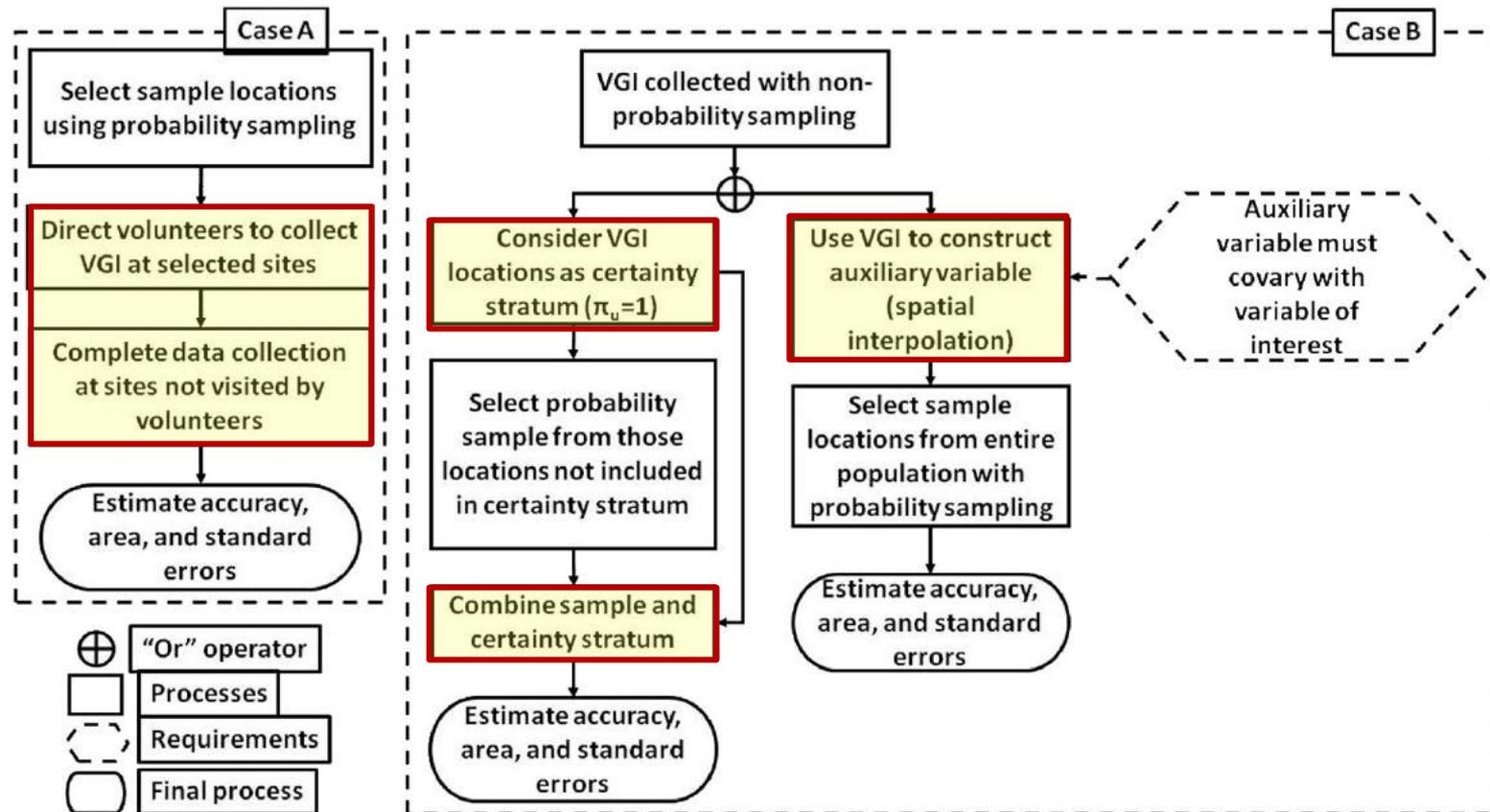[c] School of Geography, University of Nottingham, Sir Clive Granger Building, University Park, Nottingham, NG7 2RD, United Kingdom
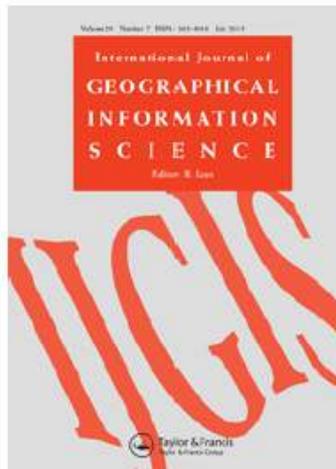[d] International Institute for Applied Systems Analysis (IIASA), Schlossplatz 1, A-2361 Laxenburg, Austria

**Flowchart:**

VGI sample
→ Probabilistic sample?
- Yes → Design-based inference
- No → Use VGI in design-based inference?
  - Yes → Obtain additional data with a probability sample → Combine VGI and probability sample → **VGI contributes to reduced standard error**
  - No → Model-based inference

# Can VGI be useful for quality assessment?

- Possible strategies:

# Use of VGI for the validation of LULC maps

- **Review article**

International Journal of Geographical Information Science

Publication details, including instructions for authors and subscription information:
http://www.tandfonline.com/loi/tgis20

## Usability of VGI for validation of land cover maps

Cidália C. Fonte[ab], Lucy Bastin[c], Linda See[d], Giles Foody[e] & Flavio Lupia[f]

[a] Department of Mathematics, University of Coimbra, Coimbra, Portugal

[b] Institute for Systems Engineering and Computers at Coimbra (INESC Coimbra), Coimbra, Portugal

[c] School of Engineering and Applied Science, Aston University, Birmingham, UK

[d] Ecosystems Services and Management Program, International Institute for Applied Systems Analysis (IIASA), Laxenburg, Austria

[e] School of Geography, University of Nottingham, Nottingham, UK

[f] National Institute of Agricultural Economics (INEA), Rome, Italy

Published online: 17 Mar 2015.

CrossMark

Click for updates

# Use of VGI for the creation and/or validation of LULCM

- Review article

| Types of VGI | Example projects |
| --- | --- |
| Photographs and descriptions | Degree Confluence Project<br>Flickr<br>Instagram<br>Panoramio<br>Geograph |
| Classification of images | Geo-Wiki<br>VIEW-IT |
| Vector maps | OpenStreetMap |

# Examples

- Use of georreferenced photographs
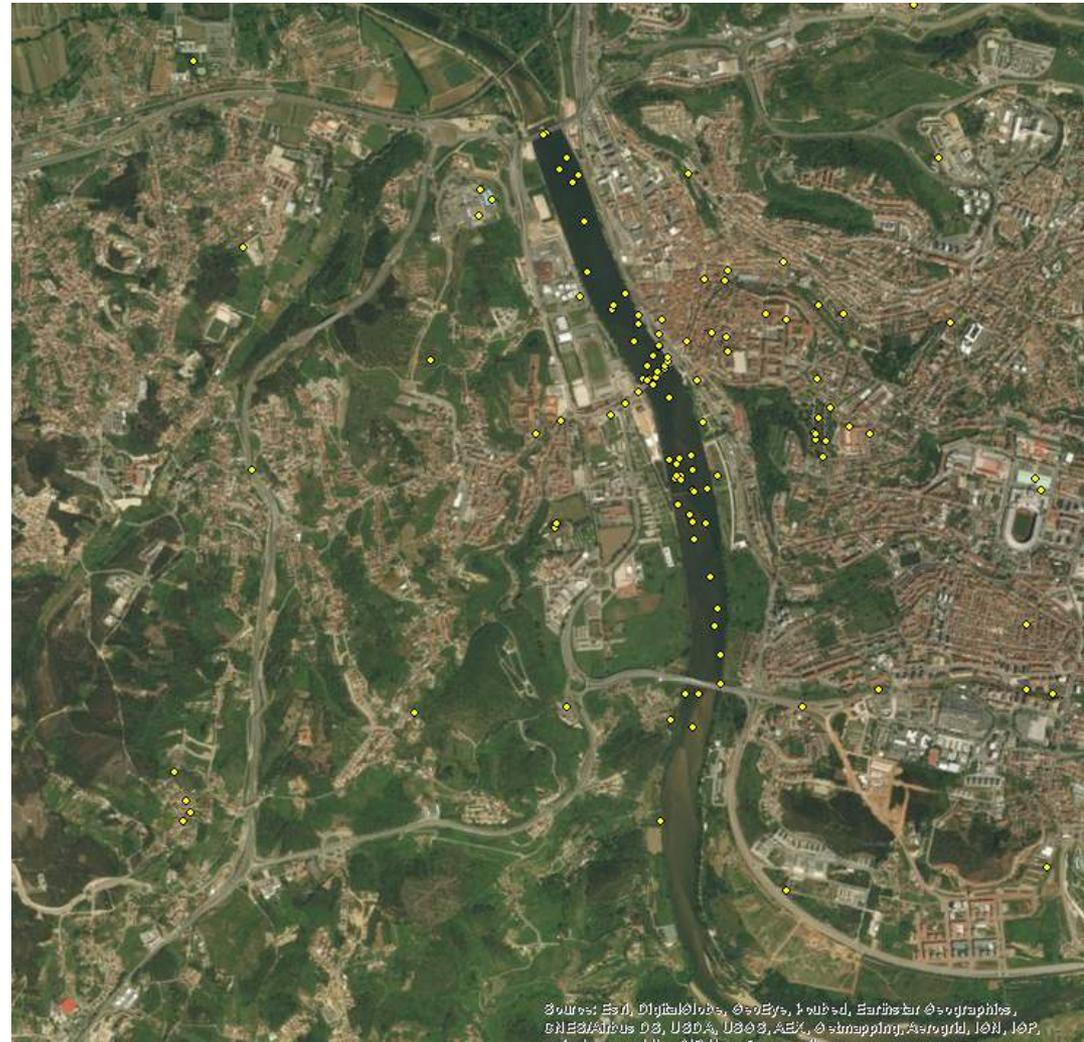- Raises problems:
  - May have positioning errors

# Examples

- Use of georreferenced photographs
- Raises problems:
  - Classify what is visible in the photo

# Examples

- Use of georreferenced photographs

- Raises problems:
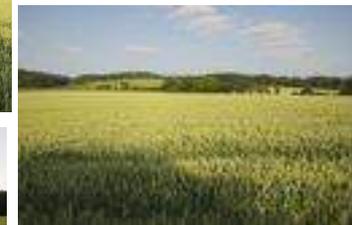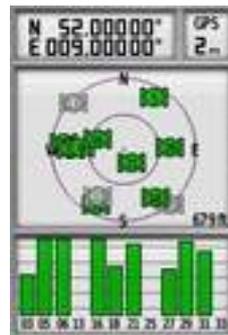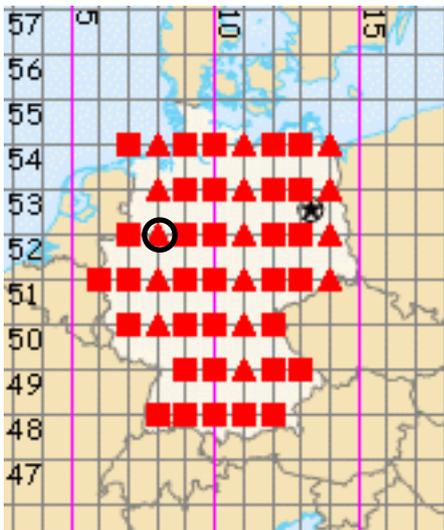  - ➢ Representativeness /geospatial distribution

# Examples



- Use of georreferenced photographs
  - Degree Confluence Project (http://confluence.org/)
    - Project created in1996
    - Collects photos and descriptions at each point with an integer value of latitude and longitude
      - Four photos are collected at each point oriented for the four cardinal directions N, S, E, W



Foody and Boyd, 2013; Iwao et al., 2006, 2011
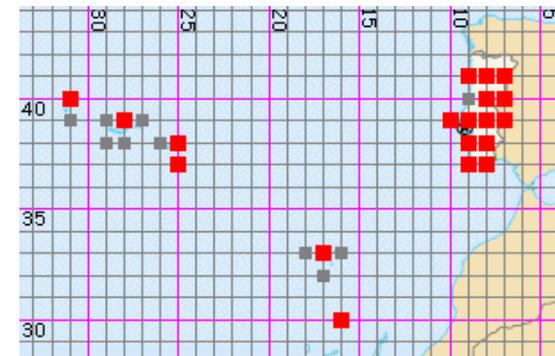
# Examples

- Use of georreferenced photographs
  - Degree Confluence Project (http://confluence.org/)
    - Points regularly spaced
      - Systematic sample of points!
    - Total points: 16 345

  - It is of little use for small regions
    - In Portugal - only 29 points
    - In Portugal main land - 14!

# Examples

- Use OSM as reference data
  - Comparison of reference data obtained from OSM and photointerpretation
    - High correspondence for level 1 classes
    - Problematic for some classes of level 2
      - Photointerpretation may also be problematic
        - Land use classes



International Journal of Geographical Information Science

ISSN: 1365-8816 (Print) 1362-3087 (Online) Journal homepage: http://www.tandfonline.com/loi/tgis20

Assessing the applicability of OpenStreetMap data to assist the validation of land use/land cover maps
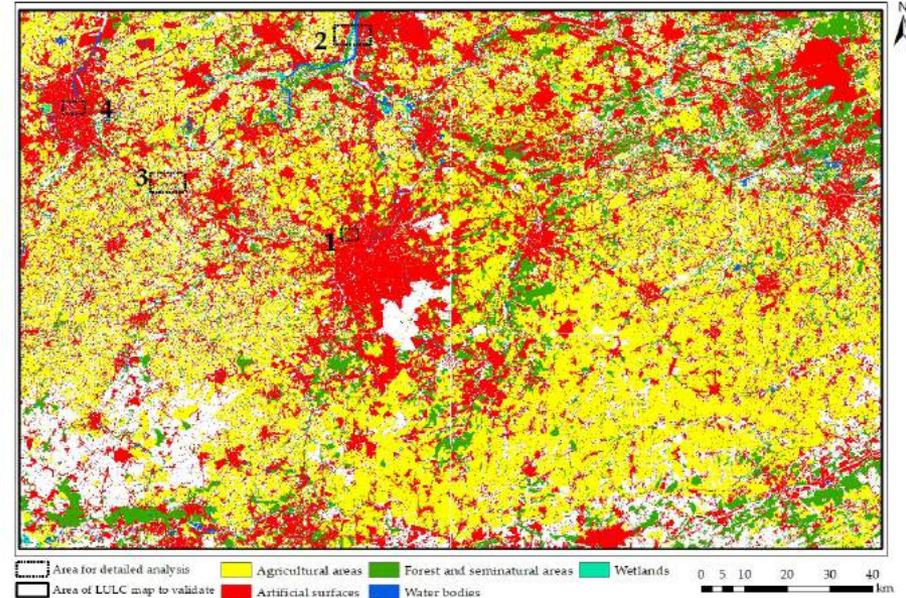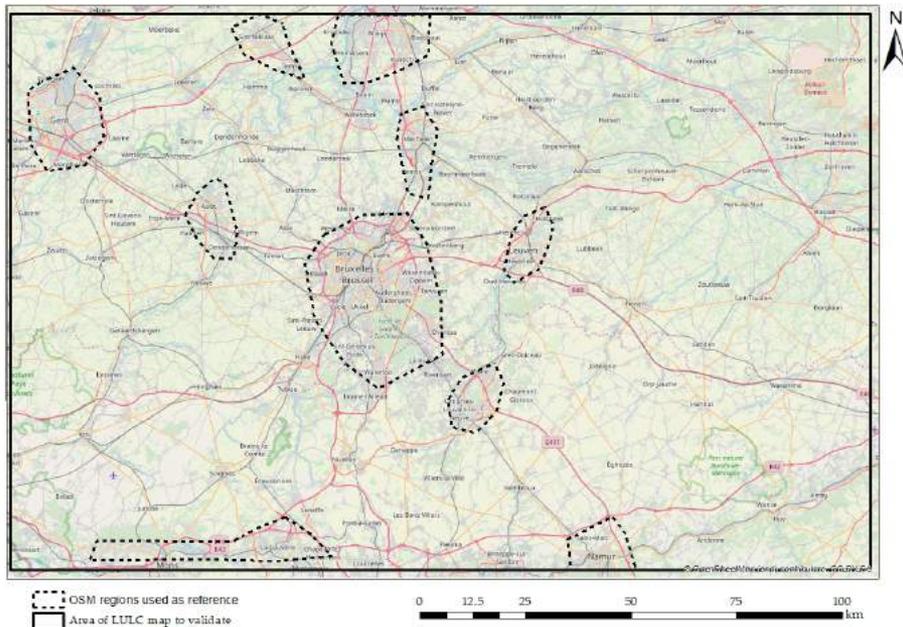
Cidália C. Fonte & Nuno Martinho

To cite this article: Cidália C. Fonte & Nuno Martinho (2017): Assessing the applicability of OpenStreetMap data to assist the validation of land use/land cover maps, International Journal of Geographical Information Science, DOI: 10.1080/13658816.2017.1358814

To link to this article: http://dx.doi.org/10.1080/13658816.2017.1358814

1.1 Urban Fabric
1.2 Industrial, commercial, public, military and private units or transport units
1.3 Mine, dump and construction sites
1.4 Artificial non-agricultural vegetated areas
2.0 Agricultural, semi-natural areas, wetlands
3.0 Forests
5.0 Water

# Examples

- ## Use OSM as reference data

  - ### Stratification



Legend (left map): OSM regions used as reference · Area of LULC map to validate

Legend (right map): Area for detailed analysis · Area of LULC map to validate · Agricultural areas · Artificial surfaces · Forest and seminatural areas · Water bodies · Wetlands

How can we:
- Reduce the validation effort
- Obtain more reliable accuracy results with less effort

# THE USE OF VOLUNTEER GEOGRAPHIC INFORMATION FOR PRODUCING AND MAINTAINING AUTHORITATIVE LAND USE AND LAND COVER DATA

# Thank you !

## Cidália Costa Fonte

Department of Mathematics – University of Coimbra, Coimbra, Portugal

Institute for Systems Engineering and Computers at Coimbra (INESC Coimbra), Coimbra, Portugal

**24-25 November 2020**