

Dipartimento di Ingegneria e Scienza dell'Informazione

– KnowDive Group –

Tourism And Waste Management in Trentino

Document Data:

November 13, 2024

Reference Persons:

Maria Amalia Pelle, Gaudenzia Genoni, Yishak Tadele Nigatu

© 2024 University of Trento
Trento, Italy

KnowDive (internal) reports are for internal only use within the KnowDive Group. They describe preliminary or instrumental work which should not be disclosed outside the group. KnowDive reports cannot be mentioned or cited by documents which are not KnowDive reports. KnowDive reports are the result of the collaborative work of members of the KnowDive group. The people whose names are in this page cannot be taken to be the authors of this report, but only the people who can better provide detailed information about its contents. Official, citable material produced by the KnowDive group may take any of the official Academic forms, for instance: Master and PhD theses, DISI technical reports, papers in conferences and journals, or books.



Index:

1	Information Gathering	1
1.1	Dataset "Waste Types"	1
1.1.1	Source Identification	1
1.1.2	Data Collection	1
1.1.3	Data Preprocessing	2
1.2	Dataset "Waste Baskets"	2
1.2.1	Source Identification	2
1.2.2	Data Collection	3
1.2.3	Data Preprocessing	3
1.3	Dataset "Municipality"	4
1.3.1	Source Identification	4
1.3.2	Data Collection	4
1.3.3	Data Cleaning and Merging	4
1.4	Dataset "Waste_production"	5
1.4.1	Source Identification	5
1.4.2	Dataset Collection	5
1.4.3	Dataset Cleaning and Merging	6
1.5	Dataset "Tourist Attraction Places"	6
1.5.1	Source Identification	7
1.5.2	Data Collection	7
1.5.3	Data Preprocessing	8
1.6	Overview of the second iTelos phase	11

Revision History:

Revision	Date	Author	Description of Changes
0.1	October 30, 202	Maria Amalia Pelle	Document created

1 Information Gathering

In this second phase of the project, we searched for and gathered the data resources needed for building the graph. The datasets identified, in particular, are:

1. Waste Types
2. Waste Baskets
3. Municipalities
4. Waste Production
5. Tourist Attraction Places

For each dataset, we dedicated three subsections to source identification, dataset collection, and cleaning/standardization procedures. Furthermore, all datasets and scripts are available on the project's GitHub page.

1.1 Dataset "Waste Types"

This first dataset provides structured information on waste types and disposal methods in the Province of Trento, supporting analysis of local waste management practices.

1.1.1 Source Identification

To obtain data regarding waste production, management, and disposal in the Province of Trento, we referred to Dolomiti Ambiente, a subsidiary of the Dolomiti Energia Group responsible for environmental hygiene services and waste collection. Although Dolomiti Ambiente operates exclusively in the municipalities of Trento and Rovereto, it was selected as a representative case study, under the assumption that waste collection practices in other municipalities within the province are managed in a similar and comparable manner. Future research could expand this study to include an analysis of waste collection regulations in municipalities beyond Trento and Rovereto.

1.1.2 Data Collection

Although Dolomiti Ambiente's website does not provide a dedicated dataset on waste types and their proper disposal methods, the company offers a downloadable PDF document titled "Riciclabolario", which serves as a guide for navigating waste sorting for domestic waste. This document became the primary source for the data extraction process. From this source, an initial CSV file was created, listing all waste types alongside their corresponding disposal categories.

1.1.3 Data Preprocessing

The initial CSV file was subjected to extensive cleaning to extract relevant data into appropriate columns and standardize the structure. The cleaning process began by splitting the *waste_type* column into two: *waste_type* and *disposal_type_1*, separating the main waste category from disposal method information following a hyphen. If *disposal_type_1* contained multiple disposal methods separated by a semicolon, the data was split further into a new column, *disposal_type_2*, to ensure each disposal method was clearly identified. Parenthetical information, such as special handling instructions or waste characteristics, was extracted and moved into separate "notes" columns (*waste_notes*, *notes_1*, *notes_2*), ensuring that relevant details were retained without cluttering the primary columns. Finally, an index column was added to provide a sequential identifier for each row, facilitating easier referencing and further indexing operations as needed.

After performing the above cleaning steps, the final dataset consists of the following columns:

- **index**: A sequential identifier for each row.
- **waste_type**: The main waste category.
- **waste_notes**: Additional information extracted from the *waste_type column*, if any.
- **disposal_type_1**: The primary disposal type.
- **notes_1**: Additional notes extracted from the *disposal_type_1* column, if any.
- **disposal_type_2**: The secondary disposal type.
- **notes_2**: Additional notes extracted from the *disposal_type_2* column, if any.

1.2 Dataset "Waste Baskets"

This dataset provides geospatial information on the distribution of waste baskets, organic bins, and recycling points for various materials across the Province of Trento.

1.2.1 Source Identification

In conducting research on waste basket distribution within the Province of Trento, we found no pre-existing datasets suitable for reuse. Consequently, data acquisition was undertaken using OpenStreetMap (OSM), a widely recognized, community-contributed mapping platform. Given the open-source nature of OSM, it provides freely accessible and editable geographical data, which serves as a valuable resource for mapping facilities and amenities. The community-driven aspect of OSM allows for dynamic updates but also introduces potential limitations in terms of data completeness and accuracy.

1.2.2 Data Collection

The data retrieval was facilitated by Overpass Turbo, a tool that allows for specific and customizable querying of the OSM database. The Overpass Turbo query used to collect data is provided below.

```
[out:json];
area["name"="Provincia di Trento"]->.a;
(
  node["amenity"="waste_basket"](area.a);
  node["amenity"="recycling"](area.a);
  node["amenity"="waste_basket"][["waste:organic"="yes"]](area.a);
  node["amenity"="recycling"][["recycling:glass"="yes"]](area.a);
  node["amenity"="recycling"][["recycling:plastic"="yes"]](area.a);
  node["amenity"="recycling"][["recycling:paper"="yes"]](area.a);
  node["amenity"="waste_disposal"](area.a);
  node["amenity"="waste_collection_point"](area.a);
);
out body;
>;
out skel qt;
```

This query yielded data on a variety of waste management facilities in the Province of Trento, including general-purpose waste baskets, bins for organic waste, and recycling points for specific materials such as glass, plastic, and paper. These data were collected in a GeoJSON file containing 3,883 points of interest and as many as 84 attributes.

Each point's precise location within the Province of Trento can be visualized on an accompanying map (please, allow a few seconds for it to load).

1.2.3 Data Preprocessing

In this case, data cleaning involved the removal of 45 columns containing mostly NaN values or information not relevant to our project, while retaining the 39 most meaningful columns. Furthermore, the GeoJSON file was transformed into a CSV format: Although the dataset is generally sparse, with many rows missing values for most columns, the available data is very useful.

The most important attributes in the dataset are:

- **id**: a unique identifier for each data point.
- **amenity**: it describes the type of facility or service (e.g., waste basket, recycling point, or waste collection point).
- **recycling types**: it indicates the types of materials accepted for recycling, such as glass, plastics, metals, paper, organic waste, and special items.
- **geometry**: it provides spatial coordinates (latitude and longitude) essential for mapping and spatial analysis.

As a final note, given the community-driven nature of the OSM platform, we acknowledge the potential limitations regarding data completeness and accuracy; it is plausible that some areas or types of bins may be underrepresented or inconsistently documented. Nevertheless, in the absence of a more comprehensive source, OSM serves as a viable data source that provides sufficient granularity and geographical coverage for our research objectives.

1.3 Dataset "Municipality"

This dataset contains the geographical boundaries of municipalities in the Province of Trento, along with the corresponding ISTAT code and population data. It can be useful for determining which municipality a user is located in.

1.3.1 Source Identification

The "Municipality" dataset was compiled using two primary sources: population data provided by ISTAT (Italian National Institute of Statistics) and geographic coordinates obtained from Open-StreetMap (OSM). ISTAT is a reliable source for demographic data in Italy, regularly publishing population statistics, including municipal-level data. The Overpass Turbo API was used to retrieve the spatial boundaries of each municipality in the specified region.

1.3.2 Data Collection

Population data for each municipality were collected directly from the ISTAT website in CSV format. These data represent the official population counts recorded on January 1, 2024. The spatial data from OSM, specifically the municipal boundary coordinates, were retrieved using queries to the Overpass Turbo platform, similar to how the waste basket data was collected. The query was designed to capture the administrative boundaries at the municipality level:

```
[out];
area["name"="Provincia di Trento"]
["boundary"="administrative"] ["admin_level"="6"];
rel(area)["boundary"="administrative"] ["admin_level"="8"];
out body;
out skel qt;
```

1.3.3 Data Cleaning and Merging

Once collected, both datasets underwent preprocessing to ensure uniformity and compatibility for merging. The ISTAT dataset was formatted with consistent column headers, and most columns were dropped, leaving only the ISTAT code and total population. Similarly, the geographic data from Overpass Turbo were structured into a GeoJSON format to facilitate easy data

merging and visualization. To create a comprehensive table that includes both population and geographic data, the ISTAT code was used as a key to merge the two datasets. This approach was preferred over using city names, as it ensures consistency—especially given that some municipalities have undergone name changes in the last five years. The GeoJSON file was enhanced by appending a "Total Population" tag to each municipality entry, allowing a single dataset to reflect both spatial and demographic information for each municipality.

The final dataset includes the following attributes for each municipality:

- **Name:** The name of the municipality
- **postal_code:** The postal code (CAP) of the municipality
- **ref:** The ISTAT code of the municipality
- **Totale:** The population count as of January 1, 2024
- **geometry:** The coordinates of the polygon defining the municipality's borders

As mentioned in the previous section, we acknowledge potential limitations regarding the completeness and accuracy of the data from the OSM platform.

1.4 Dataset "Waste_production"

This dataset contains a table with the annual waste production in tons for all the cities in the Province of Trento, covering the years from 2014 to 2022. Some data are provided in an aggregated form.

1.4.1 Source Identification

The data was sourced from the Italian Institute for Environmental Protection and Research (ISPRA), which provides comprehensive information on waste production across Italy. The datasets are publicly available on the ISPRA website.

1.4.2 Dataset Collection

These data were collected and compiled by ISPRA, offering insights into waste management trends over the specified period. The dataset were collected from their website.

1.4.3 Dataset Cleaning and Merging

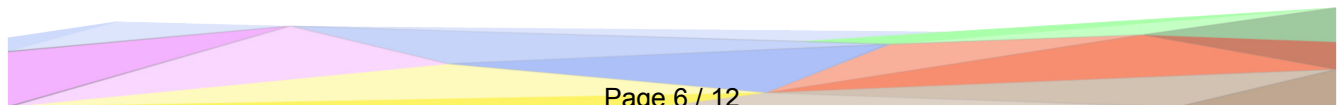
To ensure consistency, the data from all years were combined into a single dataset, with each record labeled by the corresponding year. During this process, unnecessary columns with missing or irrelevant data were removed to improve the quality and relevance of the dataset. The cleaned dataset was then saved as a single CSV file presenting the following attributes:

- **Comune:** The name of the municipality
- **Istat:** The ISTAT code of the municipality
- **Frazione organica (t):** Organic waste produced in tons
- **Ing. misti a recupero(t):** Organic waste produced in tons
- **Carta e cartone (t):** Paper and cardboard waste produced in tons
- **Altro RD (t):** General waste produced in tons
- **Legno (t):** Wooden waste produced in tons
- **Metallo (t):** Metal waste produced in tons
- **Plastica (t):** Plastic waste produced in tons
- **RAEE (t):** Electronic waste produced in tons
- **Selettiva (t):** Selective demolition waste produced in tons
- **Tessili (t):** Fabric waste produced in tons
- **Vetro (t):** Glass waste produced in tons
- **Rifiuto da costruzione e demolizione (t):** Construction and demolition waste produced in tons
- **Anno:** Year in which data was collected

Some of the rows present missing values especially for those waste that are less common and require specific processing like clothes, construction material and electronic waste.

1.5 Dataset "Tourist Attraction Places"

Tourist attraction places encompass a diverse range of categories, including natural locations, cultural landmarks, and facilities tailored to enhance visitor experiences. In this project, we considered a wide array of tourist attractions, including:



-
- **Natural Attractions:** Such as protected areas, lakes, rivers, beaches, peaks, viewpoints, caves, waterfalls, and springs.
 - **Accommodation Facilities:** Including hotels, holiday apartments, and houses.
 - **Dining and Hospitality:** Encompassing food and drink establishments such as Restaurants, Cafes and Bars.
 - **Cultural Sites:** Including museums, historic sites, and other cultural attractions.
 - **Entertainment and Recreation:** Such as amusement parks, recreational facilities, skiing and winter sports, and other adventure locations.

This comprehensive categorization reflects the broad appeal and variety of attractions available to tourists, from immersive natural environments to comfortable accommodations and lively entertainment options.

1.5.1 Source Identification

The data for this project was sourced from OpenStreetMap via the Overpass Turbo API. We executed multiple queries tailored to each attraction type to gather comprehensive information. The resulting data was retrieved in GeoJSON format and subsequently processed to align with our specific project requirements.

1.5.2 Data Collection

The following is a list of tourist attractions places which were taken into consideration in this project:

- Protected Areas
- Lakes and Rivers
- Beaches
- Peaks and Viewpoints
- Caves
- Hotel and Accommodations
- Holiday Apartments and Houses
- Food and Drink Establishments
- Cultural Attractions

- Amusement and Recreational Facilities
- Skiing and Winter Sport Facilities
- Waterfall and Springs

The following query is executed in Overpass API but. to make the data manageable, the queries are run separately and the data is downloaded separately. For instance for *Hotels and Accommodation* the following query is executed:

```
[out:json];
area["name"="Provincia di Trento"]->.searchArea;
(
  // Hotels and Accommodations
  node["tourism"="hotel"](area.searchArea);
  node["tourism"="guest_house"](area.searchArea);
  node["tourism"="hostel"](area.searchArea);
  node["tourism"="camp_site"](area.searchArea);
  node["tourism"="caravan_site"](area.searchArea);
  node["tourism"="chalet"](area.searchArea);
  node["tourism"="alpine_hut"](area.searchArea);
  node["building"="hotel"](area.searchArea);
);
out center;
```

While for *Food and Drink Establishment* the following query is executed and the result is saved in a separate geoJSON file. The total Overpass query is found in the project's GitHub repository.

```
[out:json];
area["name"="Provincia di Trento"]->.searchArea;
(
  // Food and Drink Establishments
  node["amenity"="restaurant"](area.searchArea);
  node["amenity"="cafe"](area.searchArea);
  node["amenity"="bar"](area.searchArea);
  node["amenity"="pub"](area.searchArea);
  node["amenity"="fast_food"](area.searchArea);
);
out center;
```

1.5.3 Data Preprocessing

The GeoJSON data provides extensive details about each type of attraction, resulting in a sparse dataset with information that varies significantly across attraction types. For instance, food and drink establishments are typically privately owned and include contact details such as phone numbers and emails. In contrast, natural attractions like lakes and rivers lack ownership information, as they are not associated with individual proprietors. Therefore, we carefully selected

specific features for each type of attraction to best represent them, ensuring the dataset effectively serves the purpose of our project. Table 1 provides comprehensive information on each processed attraction type.

Table 1: Summary of Generated Data Files and Their Properties

Generated File Name	Columns	Unique Values	NaN Values	Table Size
<i>caves.csv</i>	ID Name Latitude Longitude	489 242 489 489	0	489 x 4
<i>artworks.csv</i>	ID Name Latitude Longitude Artist Name Artwork Type Description Website	419 210 419 417 98 11 14 2	0 208 0 0 292 0 90 -	419 x 8
<i>memorials.csv</i>	ID Name Latitude Longitude Memorial Type Historic Type Description	504 222 504 504 14 3 80	0 244 0 0 258 0 425	504 x 7
<i>gallery_and_museum.csv</i>	ID Name Latitude Longitude Website Type Street City Postcode House Number	71 69 71 71 11 2 26 18 16 21	0 2 0 0 61 0 46 47 47 47	71 x 10

Generated File Name	Columns	Unique Values	NaN Values	Table Size
<i>food_and_drink_establishments.csv</i>	ID Latitude Longitude Name Cuisine Operator Street City Postcode House Number Website Phone Email	2293 2293 2292 1942 139 280 645 163 66 199 334 619 204	0 0 0 192 1598 2007 1185 1314 1298 1217 0 0 0	2293 x 13
<i>holiday_apartments_and_houses.csv</i>	ID Latitude Longitude Name Tourism Type Street City Postcode House Number Website Phone Email	244 244 244 236 1 128 33 21 80 79 89 82	0 0 0 2 45 37 40 51 66 22 49 44	244 x 12
<i>hotels_and_accommodations.csv</i>	ID Latitude Longitude Name Tourism Type Street City Postcode House Number Cuisine Description Operator Website Phone Email	759 759 759 688 7 303 114 56 111 8 6 37 365 326 264	0 0 0 28 1 356 388 379 365 739 753 721 0 0 0	759 x 15
<i>lakes.csv</i>	ID Central_Latitude Central_Longitude Name	38 31 38 38	0	38 x 4

Generated File Name	Columns	Unique Values	NaN Values	Table Size
<i>rivers.csv</i>	ID Central_Latitude Central_Longitude Name	138 137 138 25	0	138 x 4
<i>peaks_and_viewpoints.csv</i>	ID Latitude Longitude Name Description Historic Amenity Height Website	2542 2540 2539 1676 22 4 3 2 3	0 0 0 799 2516 2533 2524 2539 2052	2542 x 9
<i>protected_areas.csv</i>	ID Latitude Longitude Name Source Website Protection Title Leisure	3 3 3 3 3 3 2 1	0	3 x 8
<i>skiing_and_winter_sports.csv</i>	ID Latitude Longitude Name Sport Description	1292 1292 1291 461 1 1	0 0 0 644 1191 1291	1292 x 6
<i>waterfall_and_spring.csv</i>	ID Latitude Longitude Name Amenity Description	295 295 295 98 3 5	0 0 0 196 264 290	295 x 6

1.6 Overview of the second iTelos phase

In the second phase of the project, we began by discussing together the data we needed and the resources we could rely on. The decision regarding which datasets to use was not straightforward and required careful reflection on the purpose outlined in the first phase: for instance, we had to determine whether to create two distinct datasets for tourist facilities and tourist attractions or to combine both domains into a single dataset. We later decided to merge them under tourist attraction places based on our purpose definition. Once this decision was made collec-

tively, each of us took responsibility for individual tasks related to data collection and cleaning. Specifically, Amalia worked on the "Municipalities" and "Waste_Production" datasets, Gaudenzia on the "Waste Type" and "Waste Basket" datasets, and Yishak on the "Tourist Attraction Places" and "Location" datasets.

This division of labour guided the structure of the report, with each of us writing the sections corresponding to our respective datasets. We focused on detailing the data source, collection methods, and data cleaning process for each dataset.

One of the main challenges we faced during this phase was finding data related to waste management in the Province of Trento. Waste management regulations vary across municipalities, although the differences are relatively minor. To address this issue, we decided to refer to the Dolomiti Ambiente website, which provides guidelines for the municipalities of Trento and Rovereto, assuming that other municipalities follow similar regulations to those of the provincial capital.

For spatial data, we primarily relied on OpenStreetMap, utilizing Overpass Turbo to execute custom queries. The platform, being voluntary and collaborative, has both advantages and limitations. On the positive side, the available data is extensive, accurate, and reliable, making it a valuable resource for our study. However, due to its collaborative nature, the data is necessarily partial, with potential gaps or inconsistencies in some areas.

If, as we continue with the project, we find that additional data is required, we may consider reusing datasets from previous course projects, particularly those related to tourism in the Province of Trento.

Finally, this phase required us to revise the ER model, as the initial version needed adjustments. It is possible that further refinements will be necessary as the project progresses.

All raw data, pre-processing and querying scripts, and processed datasets are hosted in our GitHub repository and can be found under the 'Phase Two' section.