

# **Analysis and Ranking of Milan's Neighborhoods Operational Guide**

Edoardo Olivieri, Federica Romano, Francesca Verna

January 28, 2025

Version: 1.1

Last updated: 27-01-2024

Operational guide regarding the report:

Olivieri, E., Romano, F., & Verna, F. (2025). *Analysis and Ranking of Milan's  
Neighborhoods*.

## Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>Environment Setup</b>	<b>2</b>
2.1	Software and Libraries . . . . .	2
<b>3</b>	<b>Data Acquisition</b>	<b>2</b>
3.1	Data from OpenStreetMap . . . . .	3
3.2	Data from Yelp API . . . . .	3
3.3	Data from Comune di Milano . . . . .	4
<b>4</b>	<b>Data Processing</b>	<b>5</b>
4.1	Spatial Join . . . . .	5
4.2	Data Cleaning . . . . .	5
4.3	Home Prices . . . . .	5
<b>5</b>	<b>Data Integration and Storage</b>	<b>6</b>
5.1	Connecting to MongoDB . . . . .	6
5.2	Reading the Source Files . . . . .	7
5.3	Creating the Base Neighborhood Documents . . . . .	7
5.4	Appending Points of Interest . . . . .	8
5.5	Appending Home Prices . . . . .	9
5.6	Storing into MongoDB . . . . .	9
<b>6</b>	<b>Queries</b>	<b>9</b>

# 1 Introduction

This operational guide outlines the steps to reproduce the project aimed at analyzing and ranking neighborhoods in Milan. The analysis combines datasets retrieved from the Yelp API and the Dati Comune di Milano portal to rank neighborhoods based on their amenities, such as restaurants, nightlife, universities, and other facilities. The project leverages MongoDB for data storage, Jupyter Notebooks for processing and making interactive visualizations, and PyMongo for querying. The ultimate goal is to provide insights into how neighborhoods are structured and which are most suitable for students, singles/couples and families.

## 2 Environment Setup

To ensure the reproducibility of this project, it is essential to set up the correct computational environment. The project requires Python and a set of libraries for data handling, geospatial processing, and visualization. Additionally, MongoDB is used for data storage and querying.

### 2.1 Software and Libraries

To reproduce this project, install the following:

1. **MongoDB:** Used to store data.
2. **Jupyter Notebook:** Necessary to gather and query the data.
3. **Python libraries:** Install the necessary packages via `pip`:

```
1 pip install pandas geopandas shapely folium requests pymongo
```

- pandas for data manipulation.
- geopandas for geospatial data handling.
- shapely for geometric operations.
- folium for creating interactive maps.
- requests for accessing the Yelp API.
- pymongo for interfacing with MongoDB.

## 3 Data Acquisition

The project relies on three main data sources: OpenStreetMap (in particular Overpass turbo) for the coordinates of the neighborhoods, the Yelp API for amenities and reviews, and the Comune di Milano datasets for home prices and public facilities. These datasets are then combined and stored in MongoDB.

### 3.1 Data from OpenStreetMap

To get the data about the coordinates of neighborhoods in Milan, we need two queries from overpass turbo, which results will be downloaded as geojson files. This first one identifies correctly most of the neighborhoods:

```
1 [out:json][timeout:25];
2 area["name"="Milano"]["boundary"="administrative"]["admin_level"="8"]
3 ->.searchArea;
4 way["place"="quarter"](area.searchArea);
5 out body;
6 >;
7 out skel qt;
```

This second query is necessary to also get coordinates for the neighborhood *Forze Armate*, which can't be retrieved with the first one.

```
1 [out:json][timeout:25];
2 area["name"="Milano"]["boundary"="administrative"]->.searchArea;
3 way["place"="quarter"](area.searchArea);
4 relation["place"="quarter"](area.searchArea);
5 out body;
6 >;
7 out skel qt;
```

Additionally, from overpass turbo we also got a geojson file containing info about supermarkets, using the following query:

```
1 [out:json];
2 area["name"="Milano"]["boundary"="administrative"]->.searchArea;
3 (
4   node["shop"="supermarket"](area.searchArea);
5   way["shop"="supermarket"](area.searchArea);
6   relation["shop"="supermarket"](area.searchArea);
7 );
8 out body;
9 >;
10 out skel qt;
```

### 3.2 Data from Yelp API

To use the Yelp API, you must register for access and obtain an API key by following these steps:

1. Create a Yelp Developer Account:

- Log in with your Yelp account or create a new one if you do not already have one.

### 2. Create a New App:

- Navigate to the *Manage Apps* section in the developer portal.
- Click *Create App* and fill in the required fields:
  - App Name: Give your app a name.
  - App Description: Provide a brief description of your app.
- Once created, you will be provided with an API key.

Once the API key is obtained, using a combination of two custom functions, you can retrieve data about various points of interest:

```
1 HEADERS = {"Authorization": f"Bearer {API_KEY}"}
2 businesses_per_request = 50 # Maximum allowed by Yelp per request
3
4 def make_request(url, params=None):
5     # Makes a request to the Yelp API.
6     response = requests.get(url, headers=HEADERS, params=params)
7     response.raise_for_status()
8     return response
9
10 def search_businesses(location, term="restaurant", limit=businesses_per_request,
11 ↪ offset=0):
12     # Searches for businesses in a given location with pagination.
13     url = "https://api.yelp.com/v3/businesses/search"
14     params = {
15         "location": location,
16         "term": term,
17         "limit": limit,
18         "offset": offset
19     }
20     response = make_request(url, params=params)
21     return response.json().get("businesses", [])
```

## 3.3 Data from Comune di Milano

Download datasets for public facilities, such as schools, parks, and transport, and save them in a dedicated folder.

```
1 geojson_urls = [
2     # Insert here the urls for the geojson files
3 ]
```

```
4 output_directory = "geojson_files"
5 # Download and save each geojson file
6 for url in geojson_urls:
7     try:
8         response = requests.get(url)
9         response.raise_for_status()
10        file_name = os.path.basename(url)
11        output_path = os.path.join(output_directory, file_name)
12        with open(output_path, "wb") as file:
13            file.write(response.content)
14    except requests.exceptions.RequestException as e:
15        print(f"Failed to download {url}: {e}")
```

Regarding the dataset containing informations about home prices, it needs to be downloaded as a csv, as it is the only format available.

## 4 Data Processing

Once the raw data has been gathered from the three primary sources, the next step involves processing and integrating these datasets. The goal of this stage is to clean, standardize, and merge the data into a cohesive structure that allows for storage, analysis and visualization.

### 4.1 Spatial Join

The spatial join is a crucial step to associate point-based data with specific neighborhoods. This step uses the polygon boundaries of neighborhoods obtained from OpenStreetMap. To perform this join we propose the *geopandas.sjoin* function to match points to their corresponding neighborhoods.

```
1 joined = sjoin(point_based_data, polygon_boundaries, how="left", predicate="within")
```

### 4.2 Data Cleaning

For each dataset it is necessary to check for duplicates, potential errors and missing values. Since each dataset is different, there is no universal, one-size-fits-all approach to cleaning the data. For more detailed information about how the data was cleaned consult the original report [Analysis and Ranking of Milan's Neighborhoods](#).

### 4.3 Home Prices

Regarding the dataset containing informations about home prices, it is only available in a csv format, and therefore needs a different treatment compared to the geojson files. Since

it is not a geojson file, and therefore has no coordinates in it, we decided to join it with a geojson file containing the polygons for these Zones. First the csv file needs to be cleaned, and then can be merged:

```
1 merged = zones_polygons_file.merge(cleaned_csv_file, on="Zona", how="left")
```

Once the dataset is enriched with the polygons, using the dataset from OpenStreetMap each entry is matched to their neighborhood:

```
1 joined = gpd.sjoin(polygon_boundaries, merged, how="inner", predicate="intersects")
```

The final step is to group by neighborhood so that one entry per neighborhood is kept. At this stage, some aggregated metrics are saved to serve as general info for each neighborhood:

```
1 final = joined.groupby("Neighborhood").agg({
2     "Compr_min": "min",           # Minimum Compr_min for each Neighborhood
3     "Compr_max": "max",           # Maximum Compr_max for each Neighborhood
4     "Compr_mean": "mean",         # Average Compr_mean for each Neighborhood
5     "geometry": "first"
6 }).reset_index()
```

## 5 Data Integration and Storage

This section describes how the neighborhood and point-of-interest (POI) data is integrated into a unified structure and then stored in MongoDB.

### 5.1 Connecting to MongoDB

We begin by creating a MongoDB client using the appropriate connection string. In our example, we connect to a locally hosted database with a username and password:

```
1 from pymongo import MongoClient
2
3 client = MongoClient("mongodb://admin:DataMan2023!@localhost:27017/")
4 db = client["my_database"]
5 collection = db["neighborhoods"]
```

## 5.2 Reading the Source Files

We have multiple GeoJSON files, one for each category of POI (restaurants, museums, pharmacies, etc.). We use GeoPandas to read these files into GeoDataFrames:

```
1 import geopandas as gpd
2
3 gdf_combined = gpd.read_file("C:/path of the polygons file.geojson")
4 PolyRestaurants = gpd.read_file("C:/path of the restaurants file.geojson")
5 PolyMuseums = gpd.read_file("C:/path of the museums file.geojson")
6 ...
7 # other files
8 ...
```

## 5.3 Creating the Base Neighborhood Documents

We create a dictionary in Python to hold our neighborhood documents. Each key in the dictionary corresponds to one neighborhood. We then iterate through the rows of `gdf_combined`, which contains the neighborhood geometries and names, and create a skeleton document for each neighborhood:

```
1 neighborhood_docs = {}
2
3 for idx, row in gdf_combined.iterrows():
4     nb_name = row["Neighborhood"]
5     geo_json = mapping(row["geometry"])
6
7     neighborhood_docs[nb_name] = {
8         "_id": nb_name,
9         "neighborhood_name": nb_name,
10        "geometry": geo_json,
11        "locations": {
12            "restaurants": [],
13            "museums": [],
14            "nightlife": [],
15            ... # other POIs here ...
16        }
17        "home_prices": {
18            ... # home prices metrics ...
19        }
20    }
```



## 5.4 Appending Points of Interest

We define a function called `append_pois_to_neighborhoods` that takes a `GeoDataFrame`, a target POI key, and an optional dictionary of field mappings. This function: Groups rows by the `Neighborhood` column. Iterates over each group and locates the matching neighborhood in `neighborhood_docs`. Converts the relevant columns in each row into a dictionary (`poi_data`) using `field_mappings` to standardize field names. Appends each POI dictionary to the relevant list inside `neighborhood_docs[nb_name][“locations”][poi_key]`.

```
1 def append_to_neighborhoods(field, gdf, poi_key, field_mappings=None):
2     if field_mappings is None:
3         field_mappings = {
4             col: col
5             for col in gdf.columns
6             if col not in ("Neighborhood", "geometry")
7         }
8     grouped = gdf.groupby("Neighborhood")
9     for nb_name, group_df in grouped:
10        if nb_name not in neighborhood_docs:
11            continue
12        for _, row in group_df.iterrows():
13            poi_data = {}
14            for src_col, dest_col in field_mappings.items():
15                if src_col in row:
16                    poi_data[dest_col] = row[src_col]
17            neighborhood_docs[nb_name][field][poi_key].append(poi_data)
```

We then call the helper function for each `GeoDataFrame`, specifying the correct `poi_key` and providing the relevant `field_mappings`. For instance, for restaurants:

```
1 append_pois_to_neighborhoods(
2     gdf=PolyRestaurants,
3     poi_key="restaurants",
4     field_mappings={
5         "Business Name": "name",
6         "Business Address": "address",
7         "Categories": "category",
8         "Average Star Rating": "avg_star_rating",
9         "Review Count": "tot_ratings",
10        "Price": "price"
11    }
12 )
```

We repeated this process for museums, nightlife venues, pharmacies, etc., ensuring each POI category is correctly mapped to the corresponding dictionary fields.

## 5.5 Appending Home Prices

Once the datasets about locations are added, the informations about the home prices are inserted as follows:

```
1 for idx, row in PolyHomePrices.iterrows():
2     nb_name = row["Neighborhood"]
3     if nb_name in neighborhood_docs:
4         neighborhood_docs[nb_name]["home_prices"]["min_price"] = row["Compr_min"]
5         neighborhood_docs[nb_name]["home_prices"]["max_price"] = row["Compr_max"]
6         neighborhood_docs[nb_name]["home_prices"]["avg_price"] = row["Compr_mean"]
```

## 5.6 Storing into MongoDB

MongoDB requires inserts to be in the form of a list (or single documents). We therefore convert the `neighborhood_docs` dictionary to a list, and used `insert_many` to add all neighborhood documents at once:

```
1 documents_to_insert = list(neighborhood_docs.values())
2 collection.insert_many(documents_to_insert)
```

# 6 Queries

Once the data has been stored in MongoDB, it is possible to perform queries to retrieve and analyze the information. Below are some examples of basic queries using PyMongo to interact with the database. To see more complex queries, for ranking and scoring the Neighborhood, visit the our in depth code in the [MongoDB Integration and Queries](#).

Before querying the database, ensure that you have established a connection as shown earlier:

```
1 from pymongo import MongoClient
2
3 client = MongoClient("mongodb://admin:DataMan2023!@localhost:27017/")
4 db = client["my_database"]
5 collection = db["neighborhoods"]
```

### Example

#### 1. Retrieve All Neighborhood Names

This query fetches the names of all neighborhoods in the database:

```
1 neighborhood_names = collection.find({}, {"_id": 0, "neighborhood_name": 1})
2 for name in neighborhood_names:
3     print(name["neighborhood_name"])
```

### Example

#### 2. Find a specific Neighborhood by name

This query retrieves a document for a specific neighborhood, such as "Duomo":

```
1 neighborhood = collection.find_one({"neighborhood_name": "Duomo"})
2 print(neighborhood)
```

### Example

#### 3. Count the number of POIs in each Neighborhood

This query calculates the total number of restaurants in each neighborhood:

```
1 cursor = collection.aggregate([
2     {"$project": {
3         "neighborhood_name": 1,
4         "restaurant_count": {"$size": "$locations.restaurants"}
5     }}
6 ])
7 for doc in cursor:
8     print(f"{doc['neighborhood_name']}: {doc['restaurant_count']} restaurants")
```

**Example****4. Retrieve Neighborhoods with more than a certain number of POIs**

This query finds neighborhoods with more than 50 restaurants:

```
1 cursor = collection.find(  
2     {"$expr": {"$gte": [{"$size": "$locations.restaurants"}, 50]  
3     }  
4     }, {"_id": 0, "neighborhood_name": 1}  
5 )  
6 for doc in cursor:  
7     print(doc["neighborhood_name"])
```

**Example****5. Find the Most Expensive Neighborhoods**

This query sorts neighborhoods by their average home prices in descending order:

```
1 cursor = collection.find(  
2     {"home_prices.avg_price": {"$exists": True}},  
3     {"_id": 0, "neighborhood_name": 1, "home_prices.avg_price": 1}  
4 ).sort("home_prices.avg_price", -1)  
5 for doc in cursor:  
6     print(f"{doc['neighborhood_name']}: {doc['home_prices']['avg_price']}  
    ↪ average price")
```