



SAPIENZA
UNIVERSITÀ DI ROMA

DEPARTMENT OF COMPUTER, CONTROL, AND MANAGEMENT
ENGINEERING ANTONIO RUBERTI

**Comparative Study of Classification
Algorithms on Rome Rent Prices Dataset**
MACHINE LEARNING

Students:

Instructor:	Antonio Turco
Federico Fusco	1986183
Fabio Patrizi	Damiano Spadaccini
	1986173

Contents

1	Introduction	3
1.1	Motivation	3
1.2	Objectives	3
1.3	Report Structure	3
2	Dataset Description	4
2.1	Dataset Selection	4
2.2	Data Preprocessing	4
2.2.1	Handling Missing or Noisy Data	5
2.2.2	Complex Feature Extraction	6
2.2.3	Boolean Feature Engineering	7
2.2.4	Adding Geographic Information	7
2.2.5	Derived Metrics	10
2.2.6	Data Validation and Quality Metrics	10
2.2.7	Technical Implementation	10
2.2.8	Limitations and Future Improvements	10
2.2.9	Feature Normalization	12
2.2.10	Data Splitting	12
2.2.11	Feature Engineering (Optional)	12
3	Methodology	12
3.1	Models Implemented	12
3.1.1	Naïve Bayes	12
3.1.2	Logistic Regression	12
3.1.3	Softmax Regression (Optional)	12
3.1.4	Decision Tree	12
3.1.5	Random Forest	12
3.1.6	Support Vector Machine	12
3.2	Hyperparameter Tuning	12
4	Results	12
4.1	Model Performance	12
4.2	Confusion Matrices	12
4.3	ROC Curves and AUC	12
4.4	Training vs. Validation Performance	12
4.5	Computational Cost (Optional)	12
5	Comparative Analysis	12
5.1	Best Performing Models	12
5.2	Model Assumptions and Performance	12

5.3	Overfitting Trade-off	12
5.4	Visualizations	12
5.4.1	Learning Curves	12
5.4.2	Decision Boundaries	12
5.4.3	Feature Importance	12
6	Conclusions	13
6.1	Summary of Findings	13
6.2	Key Takeaways	13
6.3	Limitations	13
6.4	Future Work	13

1 Introduction

1.1 Motivation

1.2 Objectives

The goal of this project is to develop an understanding of foundational supervised learning algorithms by implementing, analyzing, and comparing multiple models on a real-world dataset. Specifically, we aim to:

- Implement and train various classification algorithms
- Evaluate model performance using multiple metrics
- Compare and contrast different approaches
- Analyze the trade-offs between model complexity and performance

1.3 Report Structure

During this assignment, the SEMMA methodology will be followed:

According to (Moine et al., 2011), the SEMMA methodology has five basic phases: Sample, Explore, Modify, Model and Assess. From (Olson and Delen, 2008) SEMMA facilitates the statistical exploration, the visualization techniques and the selection and transforming of the relevant variables in prediction, also can model the variables for prediction processes and later validate the precision of the model.

This report is organized as follows:

- Section 2: Description of the dataset and preprocessing steps
- Section 3: Methodology and model descriptions
- Section 4: Experimental results and evaluation
- Section 5: Comparative analysis and discussion
- Section 6: Conclusions and future work

2 Dataset Description

2.1 Dataset Selection

The dataset is taken from Kaggle and was uploaded 2 years ago by user Tommaso Ramella. It contains data scraped from the website Immobiliare.it about housing announcements in Italy in the year 2023. The author provides the public with two different datasets: the first is about rentals and the second contains purchase listings. Two different versions of each are available, a raw one, consisting of the original data, and a clean one. On the page, the author states the intention to provide a notebook with the methodology applied to parse and clean the data, but we were not able to find the complete version of this notebook (a part of the process can be found in the data publisher's GitHub folder) and as such we assume it was never released. While at the start we were inclined to just use the clean version for rents, in the end we decided against it and tried to gain more insight into the data by doing the work ourselves and trying to get as many features as we could.

The raw dataset contains around 126000 entries. Due to the number of elements in the original set, it was decided to apply machine learning techniques to a smaller subset of interest; that is, the rents in the city of Rome, a sore spot for students and a field with very practical ramifications and consequences. Our objective was to create a way for someone to get a basic understanding of the price a house should have given its characteristics, to avoid being fooled.

The records in the dataset related to Rome (`rome_rents_raw.csv`) are 13276, with 38 columns of information for each of them.

2.2 Data Preprocessing

This section outlines the preprocessing pipeline implemented to transform raw rental listings into a structured dataset suitable for analysis and modeling.

When we first looked at the raw data, we immediately noticed several issues typical of information scraped from real estate websites. The listings weren't standardized—which makes sense, since they were originally posted by different landlords with their own writing styles and conventions.

Numeric fields were particularly messy: prices included currency symbols and "/mese" suffixes, surface areas had "m²" attached to them, and some values had strange encoding artifacts like \x80 instead of the Euro symbol. Categorical information was scattered across Italian and sometimes mixed terminology, making it difficult to group similar entries together.

We also found interesting notation patterns that required careful handling. For instance, landlords often used expressions like "5+" or "3+" to indicate properties with many rooms or bathrooms, rather than listing the exact number. Floor information

proved especially tricky—descriptions ranged from straightforward numbers to Italian terms like ”piano terra” (ground floor) or ”seminterrato” (basement), sometimes combined with additional details about elevators or accessibility features.

Perhaps most challenging were the missing or ambiguous values scattered throughout the dataset. Critical fields like floor information and contract details were often incomplete or expressed in ways that needed interpretation rather than direct extraction.

The first step in the preprocessing pipeline involves standardizing the dataset structure. All column headers in the raw dataset, originally in Italian, are renamed to English to ensure consistency and readability throughout the analysis. For example, `prezzo` is renamed to `price`, `spese condominio` becomes `condo_fees`, `bagni` becomes `bathrooms`, and `quartiere` is mapped to `neighborhood`.

Following this translation, specific cleaning operations are performed on each column to handle missing values and remove noise.

2.2.1 Handling Missing or Noisy Data

Raw data from rental listings often contains inconsistencies, currency symbols, and non-numeric text in numeric fields. The following logic was applied to clean specific columns:

- **Price (price):** The columns are cleaned by stripping non-numeric artifacts such as currency symbols (€) and monthly indicators (/mese). The resulting clean strings are converted into floating-point numbers.
- **Condo Fees (condo_fees):** This column often contains various text indicators for “zero costs” (e.g., “nessuna”, “n.d.”, “nil”, “zero”). These text values are explicitly mapped to 0.0, while valid numeric entries are cleaned and parsed as floats.
- **Surface Area (square_meters):** The unit suffix (m^2) is removed from the values to isolate the numeric surface area, which is then converted to a float.
- **Deposit (deposit):** Similar to the price column, currency symbols and formatting characters are removed to convert the deposit amount into a numerical value.

Some columns contain mixed notations indicating capacity or size limits which require splitting into multiple features:

- **Rooms (rooms):** Listings often use notation like “5+” to denote large properties. This is handled by extracting the base integer (e.g., 5) for the main numeric column and creating a separate boolean flag (e.g., `more_than_5_rooms`) to capture the “greater than” information without corrupting the integer data type.

- **Bathrooms** (`bathrooms`): Similar to rooms, “3+” notations are parsed by extracting the integer 3 and creating a corresponding boolean flag (`more_than_3_bathrooms`).
- **Floor Level** (`floor`): The high variability in floor descriptions is normalized by mapping semantic Italian descriptions to integers. For instance, “piano terra” (ground floor) and “piano rialzato” (raised floor) are mapped to 0, while “seminterrato” (basement) is mapped to -1. Standard numeric floors (e.g., “1° piano”) are parsed directly to their integer equivalents.

2.2.2 Complex Feature Extraction

Beyond the basic cleaning operations, several columns contained rich information that required more sophisticated parsing strategies to extract meaningful features from what initially appeared as unstructured text.

Consider the `floor` column: rather than simply extracting the floor number, we realized that the raw text often contained additional valuable information. For instance, descriptions like “3° piano con ascensore” or “piano terra senza barriere architettoniche” provide not just the floor level, but also accessibility features. We therefore implemented a parser that extracts multiple attributes from this field:

- The numeric floor level (`floor`), with special handling for “seminterrato” (basement, mapped to -1), “piano terra” (ground floor, mapped to 0), and “piano rialzato” (raised ground floor, also mapped to 0)
- Whether an elevator is present (`lift`)
- Accessibility for people with disabilities (`disabled_access`)
- Whether it’s a raised floor (`raised_floor`)
- Whether the property spans multiple floors (`multi_floor`)
- Total number of floors in the building (`total_floors_building`)
- Whether it’s located on the top floor of the building (`last_floor`)

Similarly, the `contract` field proved to be far more complex than anticipated. Italian rental contracts come in various forms with different legal implications, and landlords often specify minimum lease durations and renewal terms directly in their listings. Our parsing approach extracts:

- Whether it’s a rental contract (`rent`)
- Whether it’s a free market rental (`free_rent`)
- Minimum rental period (`rent_min_duration`), often expressed in formats like “3+2” meaning 3 years with a 2-year renewal option

- Renewal terms (`rent_renewal_duration`)
- Whether it's a controlled rent agreement (`controlled_rent`, also known as "canone concordato")
- Whether it's a transitional/short-term contract (`short_term_rent`)
- Whether it's specifically for students (`student_rent`)
- Whether it includes a buyout option (`buyout_rent`)
- Whether it's an income property (`income_property`)

The `rooms_details` column presented another interesting challenge. Rather than a simple room count, listings typically describe the composition of the property—distinguishing between total rooms, actual bedrooms, and other spaces. We broke this down into:

- Total number of rooms (`room_total`)
- Number of bedrooms specifically (`bedrooms`)
- Count of other rooms like living rooms or studies (`other_rooms`)
- Kitchen type (`kitchen_type`), which we further categorized since Italian listings distinguish between different kitchen configurations ("cucina abitabile," "cucinotto," etc.)
- Presence of special amenities like tennis court (`tennis_court`), extracted from this field when mentioned

The `condition` (`stato`) field required splitting into two separate attributes: the overall condition of the property (`stato_condition`) such as "Buono," "Ottimo," "Da ristrutturare," etc., and the renovation status (`stato_renovation`) when applicable. These categorical values are later one-hot encoded to create binary features for each possible state.

The `parking_spaces` field contains detailed parking information that we parsed into multiple numeric features:

- Number of garage/box spots (`garage_box`)
- Number of outdoor parking spaces (`outdoor_parking`)
- Number of common parking spaces (`common_parking`)
- Number of private box spaces (`private_box`)
- Derived boolean flags indicating presence of each parking type (`has_garage_box`, `has_outdoor_parking`, etc.)

To ensure the model treats similar inputs identically, categorical and unstructured text features are transformed into consistent boolean or numeric formats. The `description` column, containing free-text descriptions of rental properties, is converted into a set of boolean features through keyword extraction. The text is scanned for high-value keywords indicating nearby amenities and property characteristics. Specifically, we extract binary features for: subway access (`subway`: “metro” / “metropolitana”), train station proximity (`station`: “stazione” / “treno”), university proximity (`university`: “università”), hospital proximity (`hospital`: “ospedale”), park access (`park`: “parco”), bathroom fixtures (`shower`: “doccia”, `bathtub`: “vasca”), property characteristics (`bright`: “luminoso” / “luminosa”, `quiet`: “silenzioso” / “silenziosa”), storage (`wardrobe`: “guardaroba”), and shopping access (`market`: “mercato” / “supermercato”). Each keyword’s presence in the description generates a corresponding binary indicator column (0 or 1), effectively transforming unstructured textual data into structured features that can be readily utilized by machine learning models. This approach captures semantic information about property location, accessibility, and desirable characteristics that would otherwise be lost in raw text format.

2.2.3 Boolean Feature Engineering

Multiple columns were transformed into comprehensive boolean feature sets:

Property Type: 22 distinct property type flags including `condominium`, various villa types (`single_family_villa`, `semi_detached_villa`, `multi_family_villa`), penthouse, `open_space`, `attic`, `loft`, `building`, `farmhouse`, terraced houses, luxury classifications (`prestigious_class`, `middle_class`, `economical_class`, `luxury_property`), and ownership structures (`full_ownership`, `partial_ownership`, `right_of_surface`, `bare_ownership`, `usufruct`, `timeshare`).

Heating Systems: 16 heating-related features covering system types (`heating_autonomous`, `heating_centralized`), distribution methods (`heating_radiators`, `heating_floor`, `heating_air`, `heating_stove`), and energy sources (`heating_gas`, `heating_methane`, `heating_gpl`, `heating_diesel`, `heating_heat_pump`, `heating_district_heating`, `heating_photovoltaic`, `heating_solar`, `heating_electric`, `heating_pellet`).

Air Conditioning: 6 features describing system configuration (`air_conditioning_autonomous`, `air_conditioning_centralized`, `air_conditioning_prearranged`, `air_conditioning_absent`) and capabilities (`air_conditioning_cold`, `air_conditioning_hot`).

Additional Amenities: 31 binary features extracted from the `other_features` column, including `terrace`, `balconies`, `garden`, `swimming_pool`, `elevator`, `garage_box`, `cellar`, `fireplace`, `furnished`, `furnished_kitchen`, `built_in_wardrobe`, `security_door`, `alarm_system`, `video_intercom`, `fiber_optic`, `tv_system`, `concierge`, `disabled_access`, `high_quality_windows`, various exposure types (`south_facing`, `east_facing`, `double_exposure`, `internal_exposure`, `external_exposure`), and specialized features (`reception`, `jacuzzi`, `tavern`, `tennis_court`, `electric_gate`).

Energy Efficiency: The `energy_efficiency` column is parsed to extract both the energy class letter (A through G) and the numeric energy consumption in kWh/m²·year. The class letter is mapped to a numeric value (`energy_efficiency_class`: A=1, B=2, ..., G=7) while the consumption value is extracted using regex patterns to handle various formatting styles.

Availability: The `availability` field is simplified to a boolean indicator (`is_available`) that checks whether the property is immediately available (“libero”) or not.

Construction Year: The `year_built` column is cleaned by removing any text and converting to integer (`construction_year`).

2.2.4 Adding Geographic Information

Since only the neighborhood is provided, we decided to add some geographic information to help the model better understand the locations. For each neighborhood found in the dataset, we added the following features: a boolean feature indicating if the location falls inside the GRA (Great Ring Road - `InsideGRA`), a zone categorical feature indicating in which position of Rome the property is located (`zone`: center, north, south, east, west), and another categorical feature indicating the administrative district (`municipality`). These geographic features are retrieved from an external mapping file (`quartieri_mapping_with_municipio.csv`) that maps each neighborhood name to its corresponding municipality, zone, and GRA status. This geographic enrichment provides the model with important location context, which is a crucial factor in determining rental prices. After this step, each row has three additional columns: `InsideGRA` (boolean), `zone` (categorical), and `municipality` (categorical).

2.2.5 Data Splitting

2.2.6 Feature Normalization

2.2.7 Feature Engineering

3 Methodology

3.1 Models Implemented

3.1.1 Naïve Bayes

3.1.2 Logistic Regression

3.1.3 Softmax Regression (Optional)

3.1.4 Decision Tree

3.1.5 Random Forest

3.1.6 Support Vector Machine

3.2 Hyperparameter Tuning

4 Results

4.1 Model Performance

4.2 Confusion Matrices

4.3 ROC Curves and AUC

4.4 Training vs. Validation Performance

4.5 Computational Cost (Optional)

5 Comparative Analysis

5.1 Best Performing Models

5.2 Model Assumptions and Performance

5.3 Overfitting Trade-off

5.4 Visualizations

5.4.1 Learning Curves

5.4.2 Decision Boundaries

5.4.3 Feature Importance

6 Conclusions

6.1 Summary of Findings

6.2 Key Takeaways

6.3 Limitations

6.4 Future Work