



# UNIVERSITÀ DI PISA

Artificial Intelligence and Data Engineering

Data Mining and Machine Learning

## *AirBnb Price Estimation using Machine Learning*

Project Documentation

---

*TEAM MEMBERS:*

Edoardo Fazzari

Mirco Ramo

Academic Year: 2020/2021

# Contents

|          |   |           |
|----------|---|-----------|
| <b>1</b> | <b>Introduction</b>   | <b>2</b>  |
| 1.1      | Goals . . . . .   | 2         |
| 1.2      | Initial Dataset . . . . .   | 2         |
| <b>2</b> | <b>Preprocessing</b>  | <b>4</b>  |
| 2.1      | What can be preprocessed now? . . . . .                             | 4         |
| 2.2      | Data Cleaning and Reduction . . . . .                               | 4         |
| 2.2.1    | Removing Noisy and Irrelevant attributes . . . . .                  | 4         |
| 2.2.2    | Removing Redundant Attributes . . . . .                             | 4         |
| 2.2.3    | Additional Removed Features . . . . .                               | 5         |
| 2.3      | Dealing with Missing Fields . . . . .                               | 5         |
| 2.4      | Data Transformation . . . . .                                       | 5         |
| 2.4.1    | Attribute Formats Transformation . . . . .                          | 5         |
| 2.4.2    | One Hot for the Amenities . . . . .                                 | 6         |
| 2.4.3    | Amenities' One Hot refinements . . . . .                            | 6         |
| 2.5      | Preprocessing Implementation . . . . .                              | 9         |
| <b>3</b> | <b>Classification</b>   | <b>10</b> |
| 3.1      | Strategies . . . . .  | 10        |
| 3.1.1    | Train and Test Splitting . . . . .                                  | 10        |
| 3.1.2    | Chosen Classifiers . . . . .  | 10        |
| 3.2      | Building Classification Models . . . . .                            | 11        |
| 3.2.1    | Data Loading . . . . .  | 11        |
| 3.2.2    | Classifier Definition . . . . .                                     | 11        |
| 3.2.3    | Classifier Implementation . . . . .                                 | 12        |
| 3.3      | Performance Evaluation and Effects of Attribute Selection . . . . . | 13        |
| 3.4      | Evaluation of Significance of the Classifiers' Results . . . . .    | 14        |
| <b>4</b> | <b>AirBnb Price Estimator App</b>                                   | <b>15</b> |
| 4.1      | Introduction . . . . .  | 15        |
| 4.2      | Use Cases List . . . . .  | 15        |
| 4.3      | Application Guide . . . . .   | 15        |
| 4.4      | Implementation Details . . . . .                                    | 15        |
| 4.4.1    | Technologies . . . . .  | 15        |
| 4.4.2    | Input Validation . . . . .  | 15        |
| 4.4.3    | Model Loading and Prediction . . . . .                              | 16        |
| <b>5</b> | <b>References</b>   | <b>17</b> |

# 1 — Introduction

Housing prices are an important reflection of the economy, and pricing a rental property on Airbnb, therefore, can be a challenging task for the owner as it determines the number of customers for the place. On the other hand, customers have to evaluate an offered price with minimal knowledge of an optimal value for the property. <sup>1</sup>

## 1.1 Goals

The aim of this paper is to explain the choices and the strategies we adopted on the project and development of **AirBnb Price Estimator**, whose aim is to help owners to decide the most correct price for their *BnB*. In order to accomplish it, we started from web-scraped data, we performed all the preprocessing needed for having a suitable dataset and then we built several classifiers, using different strategies, in order to determine the one that predicts best the class attribute. All these classifiers have been tested using more than one method and the analysis of the results guided us in the choice of the best classifier. Since the class attribute is numeric, we had two possible choices:

1. Discretize the attribute, choosing the most appropriate algorithm
2. Keep it numeric, using regression algorithms for the classification purposes

The first approach is surely easier, but it would not be as helpful as the second one for our application purposes: suggesting a precise value to an owner will give him/her a more accurate advice rather than a range.

The regression model that generalizes best the class feature has then been chosen as the “heart” of AirBnB Price Estimator: the application asks users to input the required fields that correspond to the attributes needed by the classificatory. On these fields bases, it simply outputs to the user the suggested price for night.

## 1.2 Initial Dataset

The fee of a real estate is strictly related to the city where it is located, thus would force us to recompute for every different city each step we make over and over, wasting time. Thus we decide to consider only one location (i.e, city), big enough to have a huge number of records. We chose *New York City*. The dataset, related to all the **BnB** situated in *NYC*, is taken from <http://insideairbnb.com/get-the-data.html> <sup>2</sup>. The scraping has been performed the 10th November 2020 (the scraped data has been made available by third parties).

The *initial dataset* is composed by 35821 instances and 74 columns. In order not to confuse the reader with useless information, the attribute list is not here reported, however on the following chapter regarding preprocessing we justify in detail the actions performed on the data. For now, the most important things to know is that:

1. The dataset is composed by various mixed features regarding the estate, the host and he geographical position
2. The source file has been published without respecting exactly the *csv* format, thus some frameworks and *csv* file handlers are not able to parse it

---

<sup>1</sup>GitHub repository for the project: <https://github.com/edofazza/AirBnBPriceEstimator>

<sup>2</sup>In the case of the dataset form the website inside the *listing.csv* file, related to NYC, is updated, we stored the dataset used on Google Drive. It can be downloaded at: <https://drive.google.com/file/d/1KQ2yB6eJ0rbSZoyL5fxWJjq72i40NyQn/view?usp=sharing>

3. The initial data is very dirty: there are missing values, redundant attributes, lists of strings embedded in a single column, pointless features and so on. For all of these problems, a suitable solution has been provided and it is fully reported on the next chapter.

## 2 — Preprocessing

### 2.1 What can be preprocessed now?

Since we have to deal with classification problems, before preprocessing data we must be sure not to apply supervised filters on the whole dataset. If we need supervised filters, we must split the dataset in training set and test set before applying them. However, in our case there is no need for a supervised filters but attribute selection, so we performed all the unsupervised preprocessing operations at the beginning, postponing the attribute selection after the training/test split.

### 2.2 Data Cleaning and Reduction

Since the Weka framework was not able to parse correctly the source file, for the following operations we used *Microsoft Excel* and *Apple Numbers*.

#### 2.2.1 Removing Noisy and Irrelevant attributes

The first preprocessing operation is the attribute reduction, indeed we decided to remove all the features that are not domain-specific, nor useful for classification purposes. Among them we deleted:

- **IDs:** *Bnb\_ID*, *scrape\_id*, *host\_id*
- **URLs:** *listing\_url*, *picture\_url*, *host\_url*, *host\_thumbnail\_url*, *host\_picture\_url*
- Scraping informations like *last\_scraped\_date*, *calendar\_last\_scraped*
- Useless information on the rent for classification purposes like *name*, *description*, *first\_review\_date*
- Useless information about the host like *host\_name*, *host\_location* (it is not the place in which the rent is, but where the host lives), *host\_about*, *host\_has\_profile\_pic*

#### 2.2.2 Removing Redundant Attributes

The next step was the reduction of attributes explicitly redundant. For all of these we did not compute the  $\chi^2$  test because of the explicit correlation between the features

- *host\_neighbourhood* and *neighbourhood* **w.r.t.** *neighbourhood\_cleansed* are explicit redundancies; *neighbourhood\_group\_cleansed* has been demonstrated to be very highly correlated with *neighbourhood\_cleansed* ( $P\{\text{independence}\} < 0.05$ )
- *host\_total\_listings\_count* completely equal to *host\_listing\_count*
- *host\_verification* **w.r.t.** *host\_identity\_verified*
- *minimum\_minimum\_nights*, *minimum\_maximum\_nights*, *maximum\_minimum\_night*, *maximum\_maximum\_nights*, *minimum\_nights\_avg* and *maximum\_nights\_avg* are redundancies of the attributes *minimum\_nights* and *maximum\_nights*
- *review\_scores\_accuracy*, *review\_scores\_cleanliness*, *review\_scores\_checkin*, *review\_scores\_communications*, *review\_scores\_location*, *review\_scores\_value* are rounded values that have been more precisely combined in another already-existing attribute that is *review\_scores\_rating*
- *calculated\_host\_listings\_count* (+ category) are redundant attributes of *listings\_count* valid only for their respective category

### 2.2.3 Additional Removed Features

Eventually, the other columns that have been discarded are:

1. Attributes that strongly depend from the instant in which the scraping has been performed like *has\_availability\_now*, *availability\_30*, *availability\_60*, *availability\_90*, *number\_of\_reviews\_ltm*, *number\_of\_reviews\_l30*
2. Empty attributes like *license*, *bathrooms*
3. Others, like *latitude* and *longitude* that are not more useful than the easier data on the neighborhood

## 2.3 Dealing with Missing Fields

Once the previous steps were performed, we achieved a **Weka-convertible CSV file**. Thus, the following steps have been implemented using the *Java Weka API*. The original dataset contains several missing values sparse in more than one attribute. The number of missing values is enough relevant to discourage the instance deletion, although they are pretty easy to handle, indeed:

- The majority of them are on numeric attributes characterized by low variance, like the *response\_rate* ( $\sigma = 26.17$  on a 0 – 100 interval) or the *review\_score\_rating* ( $\sigma = 9.52$  on a 0 – 100 interval). In these cases, replacing the missing value with the mean does not introduces big error rates.
- Some missing values can be easily inferred with the simple analysis by the domain expert: for example, the attribute *bedrooms* contains missing values, but only when the corresponding *beds* value is 1: it's reasonable that the missing value for *bedrooms* is 1 as well.
- Supervised approaches for missing values could have guaranteed more precise results, but since we are exploiting a regression problem and we didn't split training and test sets yet, using a supervised filter here was an error.

## 2.4 Data Transformation

### 2.4.1 Attribute Formats Transformation

Some of the attributes persisted after the cleaning and reduction phases of the precedent sections were in a not suitable state for the knowledge discovery operation that will be done in chapter 3. Thus, we decide to transform those attribute to make them adequate for classification. The features we transformed are: *amenities* (we will discuss them in chapter 2.4.2), *bathrooms*, and *price*.

The *bathrooms* in the scraped file were two columns, one empty and removed and the other one containing string values as "Shared bath", "1 shared bath", "1.5 bath". We manage to format it in two columns, one referring to the number of bathrooms inside the building and the other telling if the bathrooms are shared or not.

| bathrooms | bathroomsShared |
|-----------|-----------------|
| 1         | 0               |
| 1         | 0               |
| 1         | 0               |
| 1         | 0               |
| 1         | 0               |
| 1         | 0               |
| 1         | 0               |
| 1         | 0               |
| 1         | 0               |
| 1         | 0               |
| 1         | 1               |

The *price* attribute in the original file was a string made of \$ plus the amount (e.g., \$170). We transformed it simply by removing the dollar sign.

#### 2.4.2 One Hot for the Amenities

*Amenities*, in the original file, were list of string inputed by the used that contain all the facility the specific BnB offers to the guests.

["Wifi", "Air conditioning", "Kitchen", "Cable TV", "TV", "Elevator", "Heating"]

["Cable TV", "Essentials", "Washer", "Heating", "Air conditioning", "Kitchen", "TV", "Dryer", "Wifi"]

Structured in this way they were difficult to process. Because many machine learning models need their input variables to be numeric, we need to transform these categorical variables using the *one-hot encoding*. One-hot encoding is a frequently used method to deal with categorical data, since many machine learning models need their input variables to be numeric, categorical variables need to be transformed in the pre-processing part.

"How did we proceed to create the one-hot table?" First, we collect all the distinct values from the entries in the csv amenities' column and we created new columns using them, through a first scan of the column. Then we made a second scan to compute the values of each of the newly created attribute: if the attribute name is the same of one is contained in the amenities' list we are considering, then we put a value equal to 1, otherwise we put 0. Obtaining something like this:

| Cleaning_before_checkout | Coffee_machine | Concierge | Crib |
|--------------------------|----------------|-----------|------|
| 1                        | 1              | 0         | 0    |
| 0                        | 1              | 0         | 1    |
| 0                        | 0              | 0         | 0    |
| 0                        | 0              | 0         | 0    |

#### 2.4.3 Amenities' One Hot refinements

The number of columns obtained was very huge, more than three hundreds, so we needed to reduced them in some ways. We notice that few of them were irrelevant and the most of

them were redundant, for all of these we did not compute the  $\chi^2$  test because of the explicit correlation between the features (e.g., Keurig coffee machine, Coffee maker, Nespresso machine) but we simply merge them.

The irrelevant columns we removed through code, they were:

- *Limited housekeeping u2014 on request*: put just by one user and outdated nowadays;
- *Safe*: put just by one user and without a consistent meaning. Without knowing if safe stands that the building is safe or that the BnB is located in a safe location, we decide to remove it.

The columns that were similar to each other were merged forming a new columns made of the values of the grouped columns. Merging was done through code, summing they values in each position:

$$new\_column = \sum_i^{similar\ columns} columns_i$$

In the case of values inside the *new\_column* greater than 1, we convert back to 1. The merged columns were (*new name for the column*  $\rightarrow$  *columns merged*):

- **Toiletries**  $\rightarrow$  *1802 Beekman toiletries, Diptyque toiletries, Gilchrist & Soames toiletries, Malin+Goetz toiletries, Natura toiletries, Toiletries, Neil George toiletries, C.O. Bigelow toiletries, comfort zone toiletries, Appelles toiletries, Cote Bastide Argan toiletries, Bio Beauty toiletries, Le Labo toiletries, Elemis toiletries, MOR toiletries*
- **Stove**  $\rightarrow$  *Stainless steel gas stove, Wolf stainless steel gas stove, Viking stainless steel gas stove, Frigidaire stainless steel gas stove, Ge stove, We provide a portable gas stove in our kitchenette. gas stove, GAS COOK TOP ONLY NO OVEN gas stove, GE stove, Frigedare 30 inches stainless steel gas stove, Magic Chef gas stove, Samsung stainless steel gas stove, Stove, Electric stove, LG Stove stainless steel electric stove, Frigedare stainless steel gas stove, Stainless steel electric stove, Stainless steel stove, Gas stove, Single burner countertop range electric stove, Induction stove, 2 burner hot plate electric stove, Small portable induction stove electric stove, Two Burner Electric Cook-Top electric stove, GE electric stove, Stovetop works - Oven does not gas stove*
- **Refrigerator**  $\rightarrow$  *Magic Chef refrigerator, LG refrigerator, Samsung refrigerator, small refrigerator, Stainless Steel Frigedare refrigerator, Undercounter refrigerator, bloomberg refrigerator, Gaggenau refrigerator, Kenmore refrigerator, Undercounter Refrigerator refrigerator, Bosch refrigerator, Subzero refrigerator, Beko refrigerator, Whirlpool refrigerator, Americana refrigerator, Magic Chef refrigerator, Refrigerator, Frigedare Stainless Steel refrigerator, Sub Zero refrigerator, LG smart Tech refrigerator, Inc refrigerator, Frigidaire refrigerator, GE refrigerator, Ge refrigerator*
- **Sound system**  $\rightarrow$  *Built-in sound system in the apartment. sound system, Tivoli Audio Bluetooth sound system, Bluetooth sound system, Unknown - you can plug right into phone sound system with aux, Sonos over WiFi with built-in speakers throughout the house and backyard sound system, BOSE sound system with Bluetooth and aux, Marshall sound system with Bluetooth and aux, Marshall Bluetooth sound system, Bose Surround Speaker System in All Rooms sound system with Bluetooth and aux, Roku Bluetooth sound system, roku tv Bluetooth sound system, Echo Dot Bluetooth sound system, bose speaker Bluetooth sound system, Sound system, Bose sound system, Sound system with Bluetooth and aux, Yamaha Bluetooth sound system, Sonos sound system, Sonos Bluetooth sound system, Marshall sound system with Bluetooth and aux, Harman Kardon Bluetooth sound system, Cambridge Audio Bluetooth sound system, Sound system with aux, Samsung Bluetooth sound system, Bose sound system with Bluetooth and aux, Bose Bluetooth sound system, Yamaha sound system with Bluetooth and aux*



- **Linens** → *Supmia linens, Sferra linens, Bed linens, Frette linens, Sferra linens, linens*
- **Breakfast** → *Cooked-to-order breakfast available u2014 \$30 per person per day, Breakfast buffet available u2014 \$25 per person per day, Complimentary breakfast, Cooked-to-order breakfast available u2014 \$25 per person per day, Complimentary continental breakfast, Hot breakfast available u2014 \$20 per person per day, Complimentary hot breakfast, Cooked-to-order breakfast available for a fee, Breakfast, Cooked-to-order breakfast available u2014 \$15 per person per day, Continental breakfast available u2014 \$29 per person per day*
- **Air conditioning** → *ICE Air conditioner, Central air conditioning, Air conditioning, Portable air conditioning*
- **Dryer** → *Dryer, Hair dryer, Dryer u2013u00a0In unit, Dryer u2013 In building*
- **Extra pillows** → *Extra pillows and blankets, Bed sheets and pillows*
- **Parking** → *Free street parking, Paid parking lot on premises, Paid parking on premises u2013 1 space, Self-parking u2014 \$35/day, Valet parking u2014 \$65/day, Valet parking u2014 \$75/day, Paid street parking off premises, Valet parking u2014 \$80/day, Paid parking garage on premises u2013 1 space, Valet parking u2014 \$45/day, Paid parking on premises, Paid parking off premises, Self-parking u2014 \$19/day, Free driveway parking on premises, Paid parking lot on premises u2013 1 space, Free driveway parking on premises u2013 1 space, Self-parking u2014 \$40/stay, Valet parking u2014 \$85/day, Paid parking garage on premises, Valet parking u2014 \$70/day, Free parking on premises, Paid valet parking on premises, Self-parking u2014 \$50/day, Paid parking lot off premises, Paid parking garage off premises, Valet parking u2014 \$40/day*
- **Wifi** → *Wifi u2013 500 Mbps, Pocket wifi, Wifi u2013 60 Mbps, Wifi u2013 200 Mbps, Wifi u2013 100 Mbps, Wifi u2013 24 Mbps, Free wifi, Wifi u2013 870 Mbps, Wifi, Wifi u2013 400 Mbps*
- **Oven** → *Frigedare stainless steel oven, Frigidaire stainless steel oven, Oven, GE oven, Stainless steel oven, Toaster oven oven, Samsung stainless steel oven, Fridgedare oven, Power Airfryer 360 stainless steel oven, Small portable oven oven, Wolf stainless steel oven, Frigidaire oven, Viking stainless steel oven, electric stainless steel oven, large toaster oven oven*
- **Garden** → *Onsite bar u2014 Clinton Hall & Rooftop Beer Garden, Garden, Onsite restaurant u2014 Clinton Hall & Rooftop Beer Garden, Garden or backyard*
- **Heating** → *Heating, Radiant heating, Central heating*
- **Kitchen** → *Kitchenette, Kitchen*
- **Onsite bar** → *Onsite bar u2014 Gleason's Tavern, Onsite bar u2014 Crown Shy, Onsite bar u2014 Molyvos Restaurant - Bar, Onsite rooftop bar u2014 Make Believe, Onsite bar u2014 The National, Onsite bar, Minibar, face&body bar Bergman Kelly body soap, Onsite bar u2014 The Seville, Barbecue utensils*
- **Onsite restaurant** → *Onsite restaurant u2014 Above SIXTY SoHo, Onsite restaurant u2014 Gleason's Tavern, Onsite restaurant u2014 Butter, Onsite restaurant u2014 Parker & Quinn, Onsite restaurant u2014 Blue Ribbon Sushi Izakaya, Onsite restaurant u2014 Caf u00e9 Hugo, Onsite restaurant u2014 Scarpetta, Onsite restaurant u2014 Park Cafe, Restaurant, Onsite restaurant u2014 Caf u00e9 Boulud, Onsite restaurant u2014 Broome Caf u00e9, Onsite restaurant u2014 The National, Onsite restaurant u2014 Blue Park Kitchen*

- **Pool** → *Pool, Outdoor pool*
- **Washer** → *Washer u2013 u00a0In building, Washer, Dishwasher, Washer u2013 u00a0In unit*
- **Hot water** → *Hot water, Hot tub*
- **Gym** → *24-hour fitness center, Gym, Fitness center*
- **Coffee machine** → *Keurig coffee machine, Coffee maker, Nespresso machine, Pour Over Coffee*
- **Clothing storage** → *Clothing storage, Clothing storage-closed*

## 2.5 Preprocessing Implementation

As already mentioned before, the preliminary operations needed to make the *csv* readable by **Weka** have been performed using spreadsheet software. The resulting cleansed file still needs some modifications, *i.e.*:

- As reported in par. 2.3, missing values need to be managed
- As widely discussed in par. 2.4, the *amenities* must be converted in a **One Hot**
- On the *price* attribute, we get rid of the \$ sign and we convert it into numeric
- The *bathroom* attribute needs its format to be changed

After the loading of the *CSV* file, in order to speed up the operations we parallelized them using 4 different threads: the dataset is vertically split and every thread works only on its related partition, using *Java Weka API* or working directly on the text according to what was the fastest approach in every single scenario; every execution flow, before ending, writes results in a *CSV* file. The main thread spawns a thread for each of the tasks listed before, and then waits the termination of all of them before merging the results in a single file.

In addition to the benefits of parallelism, this approach promotes the separation of concerns of the tasks: the main method defines the preprocessing pipeline, but the actual operations on data are performed by separated and independent components. All the classes implemented at this point are collected in the *com.unipi.dmaml.airbnbpriceestimator.preprocessing* package.

## 3 — Classification

After the preprocessing phase, the dataset is ready to be used to learn regression models that will be used in the final application. In the following chapter all the chosen strategies are discussed, and the results are compared.

### 3.1 Strategies

#### 3.1.1 Train and Test Splitting

All the classifiers have been evaluated with the same strategy: *10-fold cross validation*. This means that a model is learned from the folds of the training set (90%) and they are evaluated against the corresponding test fold (10%); then the operation is repeated 10 times, and at each iteration a different fold acts as test set.

Despite the fact that Weka provides the **Evaluation.crossValidateModel()** method, we wanted to save the results of every fold and to have access to the selected attributes in case of algorithm with feature selection.

For these reasons, instead of using a meta classifier, we manually split training set and test set through the *Weka* provided functions **Instances.trainCV(numFolds, currentFold)** and **Instances.testCV(numFolds, currentFold)** according to the *WekaWiki tutorial* on how to manually construct a *k-fold cross validation*; then the *attribute selection* algorithms (when present) have been built exclusively on the training set but applied both to training and to test set.

Finally, we saved the results of each fold, but since we need only a model of the classifier while this procedure generates 10 different models, we picked the one with the lowest *root mean squared error*, aware that this does not necessarily means that we chose the best one

#### 3.1.2 Chosen Classifiers

According to our application domain, the classifiers needed are actually regression algorithms capable of handling our non-nominal class. The selected algorithms have been tested both on the full dimensionality of the dataset (138 features) and on the reduced subspaces identified by the supervised attribute selection methods. Once again, we underline the fact that those features selection algorithms have been modeled exclusively on the bases of the training set. They are *CfsSubsetEval+BestFirst* and *CfsSubsetEval+GreedyStepwise*.

We couldn't exploit the *InfoGain* evaluator since it is not capable of working with numerical classes; we discarded *PCA* since it works transforming dimensions, thus it would be difficult to understand if we could have afforded to ask only for a limited number of parameters to input from the user in the final application (e.g., if *CfsSubsetEval+BestFirst* selects only 3 attributes, we can create an application that requires only 3 parameters as input. With *PCA* we don't know which original features has been chosen in order to generate the new space, so we are forced to ask the user to input alle the 138 parameters); eventually, we discard the use of a wrapped classifier in order not to deteriorate much the performance of the application.

To summarize, the tested classifiers are:

- *Linear Regression, Linear Regression with attribute selection*
- *Random Forest, Random Forest with attribute selection*
- *5-NN, 5-NN with attribute selection*
- *M5Rules with and without attribute selection*

## 3.2 Building Classification Models

On the following pictures the implementation of the classifiers is reported. The main method is not here shown, it simply loads data and triggers the algorithm. A multithreaded approach has been tested but not implemented since the huge amount of principal memory consumption due to the model building made the system fail more than once.

### 3.2.1 Data Loading

```
1 public Instances loadData() {
2     Instances data=null;
3     try {
4         ConfigurationData configurationData = getConfigData();
5         CSVLoader loader= new CSVLoader();
6         loader.setDateAttributes(configurationData.dateAttributes);
7         loader.setDateFormat(configurationData.dateFormat);
8         loader.setEnclosureCharacters("\\");
9         loader.setFieldSeparator(",");
10        loader.setMissingValue("?");
11        loader.setNominalAttributes(configurationData.nominalAttributes);
12        loader.setNumericAttributes(configurationData.numericAttributes);
13        loader.setSource(new File(datasetPath));
14        data=loader.getDataSet();
15    } catch (Exception e){
16        e.printStackTrace();
17    }
18    return data;
19 }
20
21 private ConfigurationData getConfigData() throws JAXBException {
22     JAXBContext jaxbContext = JAXBContext.newInstance(ConfigurationData.
23         class);
24     Unmarshaller jaxbUnmarshaller = jaxbContext.createUnmarshaller();
25     ConfigurationData conf = (ConfigurationData) jaxbUnmarshaller.unmarshal
26     (new File(confDataPath));
27     System.out.println(conf.dateAttributes + "\\n" + conf.dateFormat + "\\n"
28     + conf.numericAttributes);
29     return conf;
30 }
```

**DatasetFromCsvLoader.java** is the class in charge of loading the dataset stored as a *CSV file*. It simply uses the *Weka CSVLoader* but the options have been saved in a configuration file, so that to maximize the separation of concerns between the *Weka* internal representation of the data and our code.

### 3.2.2 Classifier Definition

```
1 private void buildLinearRegressionWithAttributeSelection() {
2     try{
3         Instances randData = new Instances(dataset);
4         randData.randomize(new Random(1));
5
6         for(int i=0; i<numFolds; i++){
7             AttributeSelection filter = new AttributeSelection();
8             filter.setEvaluator(new CfsSubsetEval());
9             filter.setSearch(new BestFirst());
10            executeCV(randData, i, filter, "CfsSubsetEval+BestFirst");
11        }
12
13        randData = new Instances(dataset);
```

```

14         randData.randomize(new Random(1));
15
16         for(int i=0; i<numFolds; i++){
17             AttributeSelection filter = new AttributeSelection();
18             filter.setEvaluator(new CfsSubsetEval());
19             filter.setSearch(new GreedyStepwise());
20             executeCV(randData, i, filter, "CfsSubsetEval+GreedyStepwise");
21         }
22
23         System.out.println("linear regression with attribute selection
terminated");
24
25     } catch (Exception e) {
26         e.printStackTrace();
27     }
28 }

```

Every algorithm is implemented in a separate class, which contains different methods based on the way we want to test it (*with* or *w/o* attribute selection). In the code shown the classifier definition is reported: after a randomization of the dataset, we call *numFolds* time the `executeCV()` that performs a single round of the *cross-validation*.

### 3.2.3 Classifier Implementation

```

1 private void executeCV(Instances dataset, int currentFold, AttributeSelection
   filter, String filterName) throws Exception {
2     Instances train = dataset.trainCV(numFolds, currentFold);
3     Instances test = dataset.testCV(numFolds, currentFold);
4     List<Attribute> chosen=new ArrayList<>();
5     if(filter!=null) {
6         filter.setInputFormat(train);
7         Instances trainReduced = Filter.useFilter(train, filter);
8         Instances testReduced = Filter.useFilter(test, filter);
9         train = new Instances(trainReduced);
10        test = new Instances(testReduced);
11        for(int i=0; i<train.numAttributes(); i++)
12            chosen.add(train.attribute(i));
13    }
14    weka.classifiers.functions.LinearRegression classifier = new weka.
classifiers.functions.LinearRegression();
15    classifier.buildClassifier(train);
16    Evaluation evaluation = new Evaluation(train);
17    evaluation.evaluateModel(classifier, test);
18    new FileSaver(evaluation, "LinearRegression", filterName, currentFold,
chosen).save();
19    if(currentFold==8)
20        saveModel(filterName, classifier, test);
21 }
22
23 private void saveModel(String filterName, weka.classifiers.functions.
   LinearRegression classifier, Instances dataFormat) throws Exception{
24     SerializationHelper.write("models/LinearRegression_" + filterName + ".
model", classifier);
25     ArffSaver saver = new ArffSaver();
26     saver.setInstances(dataFormat);
27     saver.setFile(new File("models/LinearRegression_" + filterName + "_data
.arff"));
28     saver.writeBatch();
29 }

```

In this method we perform a round of the **10-fold cross validation**, in the way suggested by the *WekaWiki tutorial*. At each iteration, the *trainCV* and *testCV* return non-overlapping

training sets and test sets; if an attribute selection method has been defined, it is built on top of the training set and applied to both of them and the list of the chosen attributes is stored. In the end, the model built on the **fold n.8** is stored in a *.model* file.

### 3.3 Performance Evaluation and Effects of Attribute Selection

Once the classifiers have been built, as anticipated in the paragraph 3.1 we tested them using a **10-fold cross validation**. In the following table we report the average values respectively of *correlation coefficient (CC)*, *mean absolute error (mae)*, *root mean squared error (rmse)*, *relative absolute error (rae)* and *root relative squared error (rrse)* for each tested classifier (note that the chosen model, *i.e.* the one built in the 8<sup>th</sup> fold, shows better measures than the average for every tested classifier). The green-colored values are the best ones in absolute.

| ATTRIBUTE SELECTION $\Rightarrow$<br>REGRESSION ALGORITHM $\Downarrow$ | None   | CfsSubsetEval<br>+ BestFirst   | CfsSubsetEval<br>+ GreedyStep-<br>wise             |
|--|--|--|--|
| <b>Linear Regression</b>   | Out of<br>memory<br>error                                    | $\overline{CC} = 0,6793$<br>$\overline{Mae} = 40,0323$<br>$\overline{Rmse} = 57,8399$<br>$\overline{Rae} = 66,9399$<br>$\overline{Rrse} = 73,3966$ | 0,6794<br>40,0308<br>57,8339<br>66,9375<br>73,3891 |
| <b>Random Forest</b>   | 0.7495<br>36,0881<br>52,8727<br>60,3423<br>67,0864           | 0.6988<br>37,9716<br>56,3332<br>63,4916<br>71,7302   | 0,7000<br>37,9451<br>56,5060<br>63,4470<br>71,6955 |
| <b>5-NN</b>  | 0,5413<br>46,4003<br>67,0074<br>77,2522<br>84,3622           | 0,6772<br>39,0709<br>58,0759<br>65,3630<br>73,6845   | 0,6772<br>39,0709<br>58,0759<br>65,3630<br>73,6845 |
| <b>M5Rules</b>   | 0,6589<br>567,0975<br>31.491,7374<br>942,2971<br>39.372,3750 | 0,6942<br>38,5873<br>56,7530<br>64,5234<br>72,0102   | 0,6942<br>38,5872<br>56,7498<br>64,5211<br>72,0100 |

As shown in the table, *Random Forest* without attribute selection is the algorithm that performs best **w.r.t.** all the parameters. Note that the attribute selection improved the outcomes of certain algorithms, anyway, deteriorating the ones of *Random Forest*. Now we will discuss the pros and cons of each of them, justifying the final choice for our application:

1. All the classifiers without attribute selection are, unfortunately, not very suitable for the usage of the application: they would require the user to input 288 parameters, which is a huge amount of work to do for him/her. On the contrary, all the attribute selected classifier needs less than 30 features.
2. *Linear Regression* shows good but not awesome results, it requires not much memory for the model loading and it is very fast at prediction time. However, it showed very long times and high memory consumption at training time.

3. *Random Forest* is the one with the best results in absolute, both with and without attribute selection. It is quite fast both to train and to use for prediction, anyway it needs a lot of time and memory for the model to be saved and loaded. Since the considerations at the point 0), we discarded the pure *Random Forest* and we chose *Random Forest with CfsSubsetEval+GreedyStepwise* as the default regression model for our application.
4. *K-NN* is definitely inadequate, although it requires very little time and memory for the model to be built (and then also to be loaded), the computational effort is concentrated at prediction time (thus increasing the application response time), and it also shows quite poor results.
5. *M5Rules* has very good performances and it needs very little time and memory to save and load the model. The prediction is not so fast, the time required to build the model is terribly long (up to 3 or 4 times more than *Random Forest*). Anyway, since the majority of the effort is at training time, it is suitable for our application, thus *M5Rules with CfsSubsetEval+GreedyStepwise* has been chosen as backup regression model.

### 3.4 Evaluation of Significance of the Classifiers' Results

Once evaluated the performance of the algorithms, the main question is if we can consider those differences statistically significant or not. In order to make such a computation, we used the paired **Student's t-test** on the *RMSE* considering all the outcomes of all the folds (of every algorithm we chose only the *CfsSubsetEval+ GreedyStepwise* version). The results are the following (*sig\_rate* = 5%, *deg* = 9):

- *Random Forest w.r.t. Linear Regression*:  $t=5.5342$ ;  $p=0.004 \Rightarrow$  stat. sign.
- *Random Forest against M5Rules*:  $t=0.805$ ;  $p=0.4413 \Rightarrow$  non stat. sign.
- *M5Rules against Linear Regression*:  $t=5.5104$ ;  $p=0.032 \Rightarrow$  stat. sign.
- *Linear Regression against KNN*:  $t=0.6601$ ;  $p=0.5257 \Rightarrow$  non stat. sign

*Random Forest* and *M5Rules* are hence comparable, so they are equally valid for our application. On the contrary, the **t-test** shows that *Linear Regression* 's results are much more similar to *KNN* ones rather than to the outcomes of the other two algorithms, thus it has been discarded.

## 4 — AirBnb Price Estimator App

### 4.1 Introduction

**AirBnB Price Estimator** is a real estate cost estimation application in which the Users have the possibility to worth their own Bed and Breakfast simply by inserting some information related to it, in a very simple and fast way.

### 4.2 Use Cases List

There is no need for registration or authentication, every user can operate with the applicative as they open it. The can:

- Select the classifier they want to predict the price of their B&B between:
  - *Random Forest*
  - *M5Rules*
- Insert the requested values for B&B related characteristics
- Estimate the most suitable price using the algorithm previously selected

### 4.3 Application Guide

When the application is launched a window is opened telling the **User** to choose between the two classifiers previously discussed. The User cannot undo the choose, once selected the classifier, in order to change it he/she must reopen the application. Selected the classifier, a new scene is loaded. This scene displays some *Text Fields* and *Radio Buttons*, corresponding to the subset of the overall features in our dataset that have been chosen by the **Attribute Selection process** (see chapter 3.3). If a *Radio Button* is not selected the system will understand that as the absence of the particular amenity the *Radio Button* refers to. After the **User** inputs all the values, a *Button* called "*PREDICT PRICE*" is located at the bottom of the window (scroll down if you don't see it). A click on it will display in a label the predicted price, just below the *Button*. The user can perform multiple predictions sequentially, just changing the contents of the *Text Fields* or the selection of the *Radio Buttons*.

### 4.4 Implementation Details

#### 4.4.1 Technologies

In order to ensure a higher level of accessibility to the users, we decided to adopt the **JavaFx** framework to create a very simple and intuitive *Graphical User Interface*. For what concerns the prediction system, we used again the *Java Weka API* so that to make use of the model stored on the previous phases, with no need to compute it again.

#### 4.4.2 Input Validation

To make the usage of the application easier, the User is not required to know the exact format of the nominal attributes, indeed such a property is ensured by an input validation system. In fact, to standardize the nominal attributes, in the preprocessing step we unified the format that is:

- First letter uppercase
- No other letter can be uppercase



- All the special characters, included the blank space, are replaced by the underscore

After the user clicks the “*PREDICT PRICE*” button, those operations are performed again on the input data. Anyway, the system is not capable of recognizing misspelled words and replacing it (e.g, If the use writes *Alerton* instead of *Allerton*, the system will recognize it as a value out of the admitted range of the nominal attribute).

#### 4.4.3 Model Loading and Prediction

When the user clicks on the classifier he/she wants to use, automatically the system loads it, through the *Weka Api* function **SerializationHelper.read(File)**. Once loaded, it is kept in memory for the application lifespan, so that to speed up multiple predictions performed one after the other. After that, all the attributes that the classifier needs are passed to the GUI, which generates a new input field for each of them (a Radio Button for the binary ones, a *TextField* for all the others). The prediction of the price is made through the function **Classifier.classifyInstance(Instance)** that returns a double value. After the 2-decimal rounding, the value is output to the user.

## 5 — References

1. Y. Ma, Z. Zhang, A. Ihler, B. Pan - Estimating Warehouse Rental Price using Machine Learning Techniques
2. H. Yu, J. We - Real Estate Price Prediction with Regression and Classification