

# Machine Learning: Exercise 0

Amélie Assmayr (12007770)

Konstantinos Damanakis (12106343)

Teresa Schuch (12007762)

## 1 Introduction

In this Exercise we have chosen two datasets: Laptop Prices and Road Traffic Accidents. The Laptop Prices dataset will be utilized to develop regression models for predicting laptop prices, while the Road Traffic Accidents dataset will be employed to classify the severity of accidents. We chose these datasets because they enable us to try various techniques as they differ in important aspects. The Laptop Prices dataset has 2275 instances, 15 features and no missing values. In contrast, the Road Traffic Accidents dataset contains 12316 instances, 32 features and missing values in 16 features, making the preprocessing more important and complex.

## 2 Dataset 1: Laptop Prices

The Laptop Price dataset contains information about various laptops and their price in euros. We chose this dataset because of its importance in real-life situations. Most people today rely on laptops for work, education and entertainment. Predicting these prices can provide retailers with crucial insights into expected revenues before official pricing is announced. Additionally, the variety of features in the dataset, including both categorical and numeric data makes it ideal for experimenting with feature engineering.

The variables, their data types and the number of unique values for nominal and ordinal variables can be seen below.

Variable	Datatype	Unique	Variable	Datatype	Unique
<i>Company</i>	nominal	19	<i>RAM (GB)</i>	ratio quantity	/
<i>Product</i>	nominal	618	<i>Memory</i>	ordinal	39
<i>Type Name</i>	nominal	6	<i>GPU Company</i>	nominal	4
<i>Inches</i>	ratio quantity	/	<i>GPU Type</i>	nominal	106
<i>Screen Resolution</i>	ordinal	40	<i>Operating System</i>	nominal	9
<i>CPU Company</i>	nominal	3	<i>Weight (kg)</i>	ratio quantity	/
<i>CPU Type</i>	nominal	106	<i>Price (Euro)</i>	ratio quantity	/
<i>CPU Frequency (GHz)</i>	ratio quantity	/			

### 2.1 Target attribute

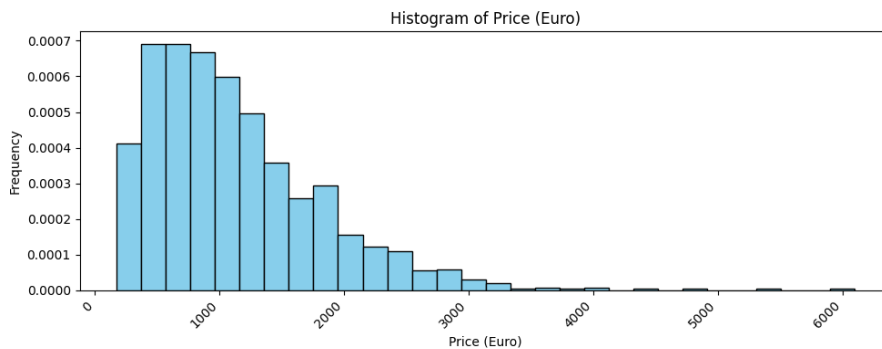


Figure 1: Histogram of Laptop Prices

The target variable is the price of the laptop. The prices range between 174€ and 6099€. The histogram in Figure 1 illustrates the distribution of the prices, revealing that most laptops are priced between 400 and 800€. The mean price is 1135€ and the median price 989 €. There are a few outliers that cost over 4000€, but these outliers are valid values.

### 2.2 Additional Attribute Insights

The screensize ranges from 10.10 to 18.40 inches, with a mean of 15.02 inches. The weight varies between 0.690 kg and 4.7 kg, with an average weight of 2.041 kg. The CPU frequency spans from 0.9 to 3.6 GHz, with a mean value of 2.5 GHz and RAM takes integer values between 2 and 64 GB, with a median value of 8 GB. The variable product may require preprocessing and summarizing due to its numerous unique values.

Due to limited space not all variable plots could be displayed. In Figure 2 the distribution of RAM and the Top 10 Companies can be seen as an example.

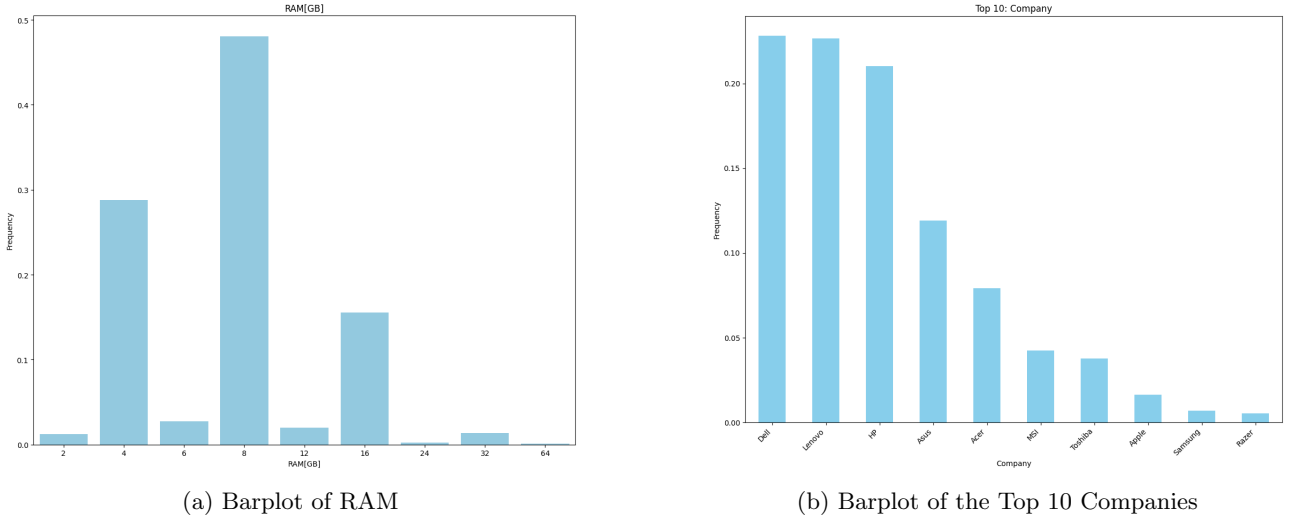


Figure 2: Plots from Laptop dataset

### 3 Dataset 2: Road Traffic Accidents

The second dataset includes instances of road traffic accidents recorded in the years 2017 to 2020 in Addis Ababa in Ethiopia. This dataset contains information regarding severity of the accident, the number and the characteristics of the involved vehicles and people, area etc. We chose this dataset because it addresses a critical public safety issue. A study of this data could help predict the severity of the accidents which can help identify causes and factors contributing to road safety.

Missing data can be observed in 16 of the 32 features. Overall there are 20,057 missing values. In the following table all attributes, their datatype (r=ratio quantity, o=ordinal, n=nominal), the number of unique elements for categorical variables and the number of missing values can be seen :

Variable	Type	Unique	Missing	Variable	Type	Unique	Missing
<i>Time</i>	r	/	/	<i>Road Surface Condition</i>	n	4	172
<i>Day of Week</i>	n	7	/	<i>Light Condition</i>	n	4	/
<i>Age of Driver</i>	o	/	/	<i>Weather Condition</i>	n	9	/
<i>Sex of Driver</i>	n	3	/	<i>Collision Type</i>	n	10	155
<i>Educational Level</i>	n	7	741	<i>Number of Vehicles</i>	r	/	/
<i>Vehicle-Driver Relation</i>	n	4	579	<i>Number of Casualty</i>	r	/	/
<i>Driving Experience</i>	o	/	829	<i>Vehicle Movement</i>	n	13	308
<i>Vehicle Type</i>	n	17	950	<i>Casualty Class</i>	n	4	/
<i>Vehicle Owner</i>	n	4	482	<i>Sex of Casualty</i>	n	3	/
<i>Vehicle Service Year</i>	r	/	3928	<i>Age of Casualty</i>	r	/	/
<i>Defect of Vehicle</i>	n	3	4427	<i>Casualty Severity</i>	o	4	/
<i>Accident Area</i>	n	14	239	<i>Work of Casualty</i>	n	7	3198
<i>Lanes or Medians</i>	n	7	385	<i>Fitness of Casualty</i>	n	5	2635
<i>Road Alignment</i>	n	9	142	<i>Pedestrian Movement</i>	n	9	/
<i>Junction Type</i>	n	8	887	<i>Cause of Accident</i>	n	20	/
<i>Road Surface Type</i>	n	4	172	<i>Severity of Accident</i>	n	3	/

#### 3.1 Target Attribute

The target variable is the severity of the accident. It is classified into three classes: Light Injury, Serious Injury or Fatal Injury. While the severity could be considered an ordinal variable, we will treat it as nominal, framing the task as a multiclass classification problem. The following plot shows the distribution of these classes. It can be seen that the classes are unevenly distributed. Most accidents fall under the "Slight Injury" category, followed by "Serious Injury" and only a few accidents are classified as "Fatal Injury." This imbalance is a critical factor to consider. Methods like bootstrapping might be helpful to employ.

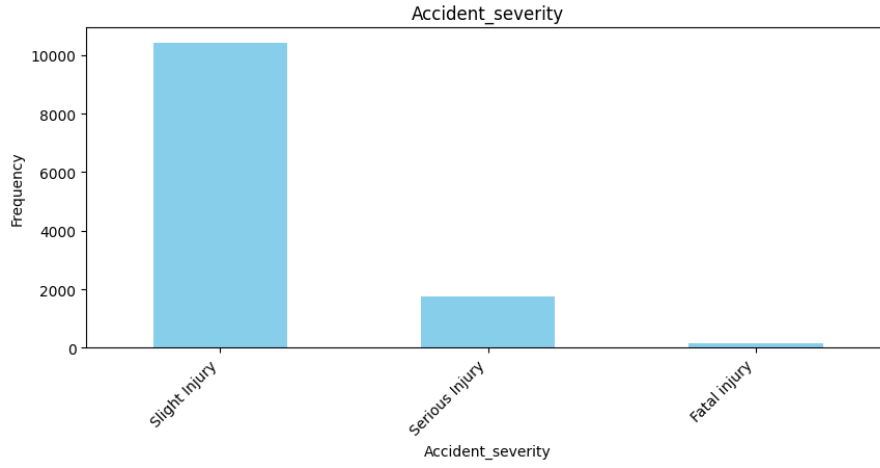
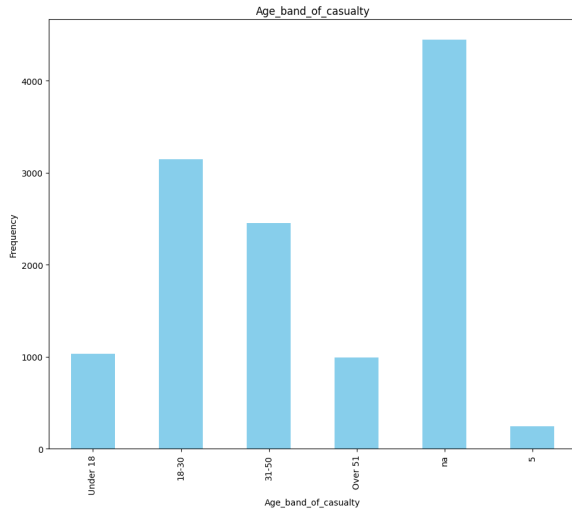


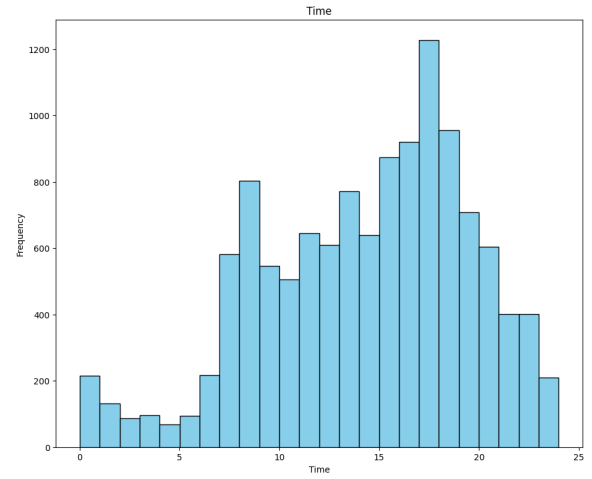
Figure 3: Barplot of Accident Severity

### 3.2 additional attribute insights

Most of the numeric variables in the dataset are grouped into intervals, effectively converting them into ordinal data. Besides regular missing values, the dataset also has other entries like 'Unknown' or 'na' that should be handled during preprocessing. As an example the histogram of the time of the accident and the barplot of the age classes of the driver are shown below in Figure 4:



(a) Barplot of the age classes



(b) Histogram of the Time

Figure 4: Plots from Road Traffic Accident dataset