



Project Presentation

Data Management

Florian Deroua and Eve Schmitz Schair

TABLE OF CONTENTS

01

INTRODUCTION

02

WEBCRAPING

03

VISUALIZATION

04

**MACHINE
LEARNING**

05

CONCLUSION



01

INTRODUCTION



Motivation

Question: can we predict the expected price for adult tickets for attractions in a specific city based on the type of attraction and historical pricing data ?

To respond to the question, we used a website with all the attractions of Wallonia and Brussel.

02

WEBCRAPING



DATA

Name	Type	City	Price adult	Price child
La boverie	Art	Liège	5€	0€

Columns: 5

Lines: 289

We made a list that stores the information about attractions extracted from the website. Each code is for a different type of attraction. Each attraction's information is stored as a dictionary containing the name of the attraction, the type of attraction, the city where the attraction is and the price for each group of persons.

03

VISUALIZATION

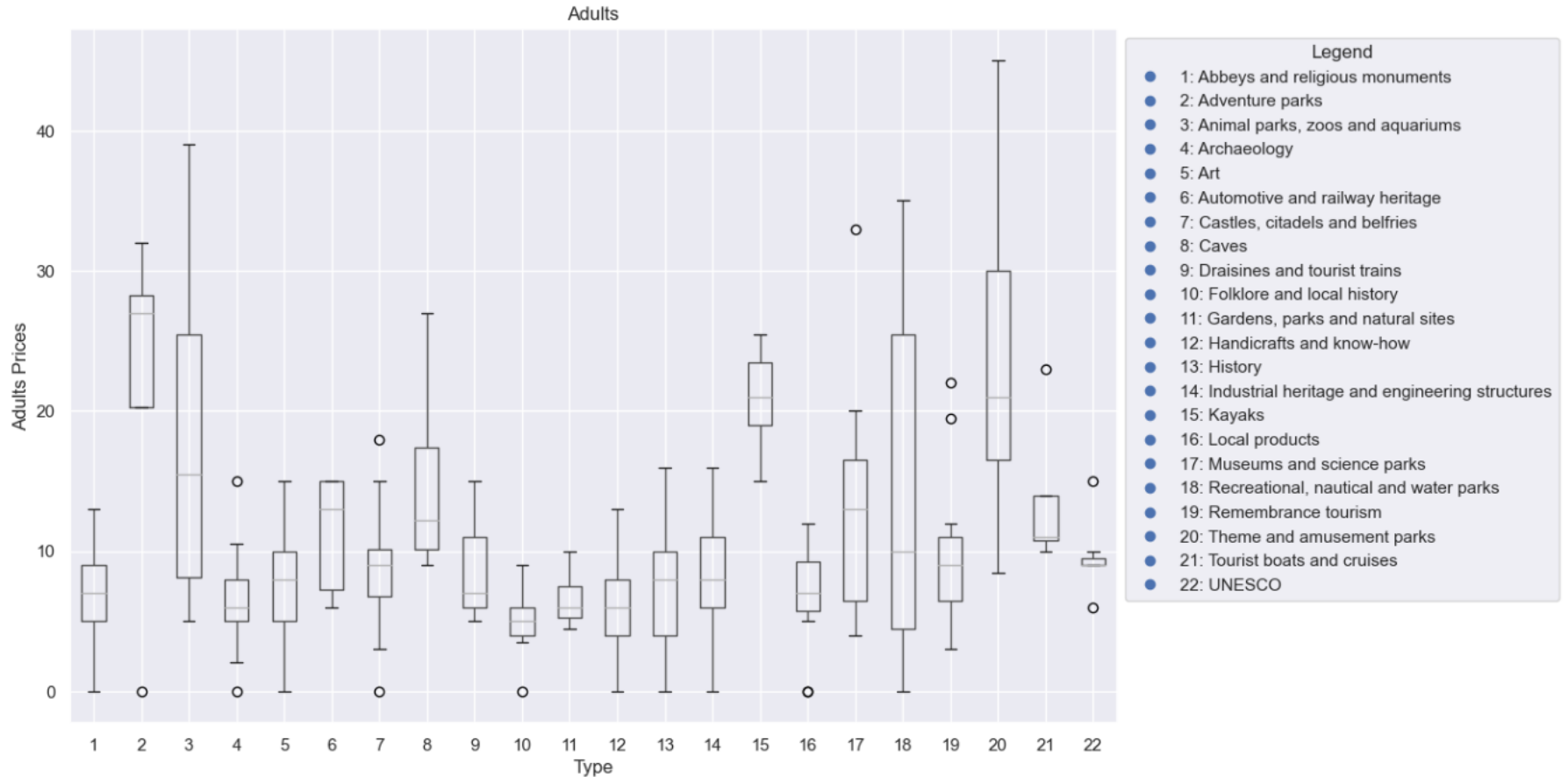


BAR CHART

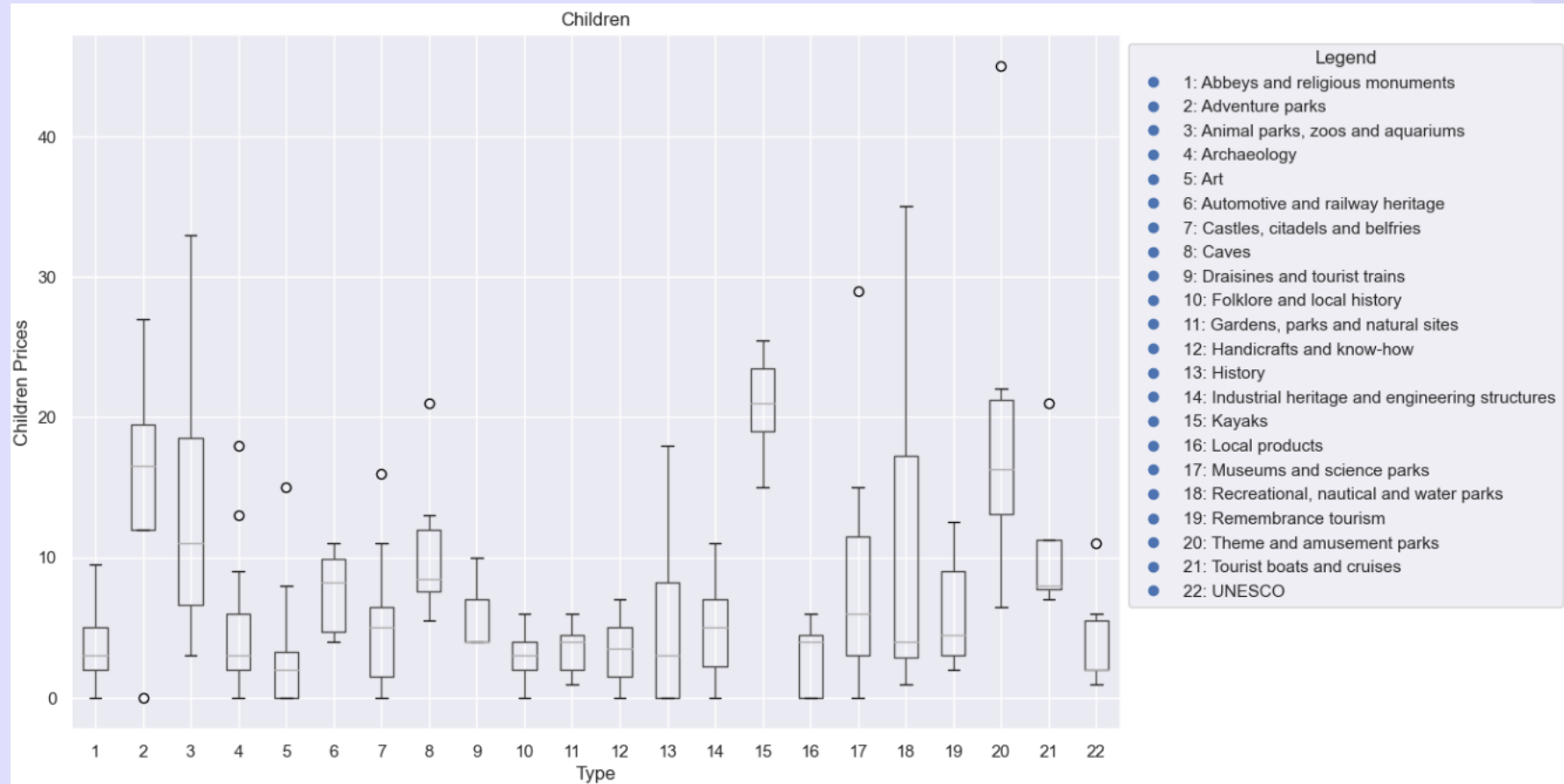
We can observe that the prices for children are very often lower than the prices of adults.



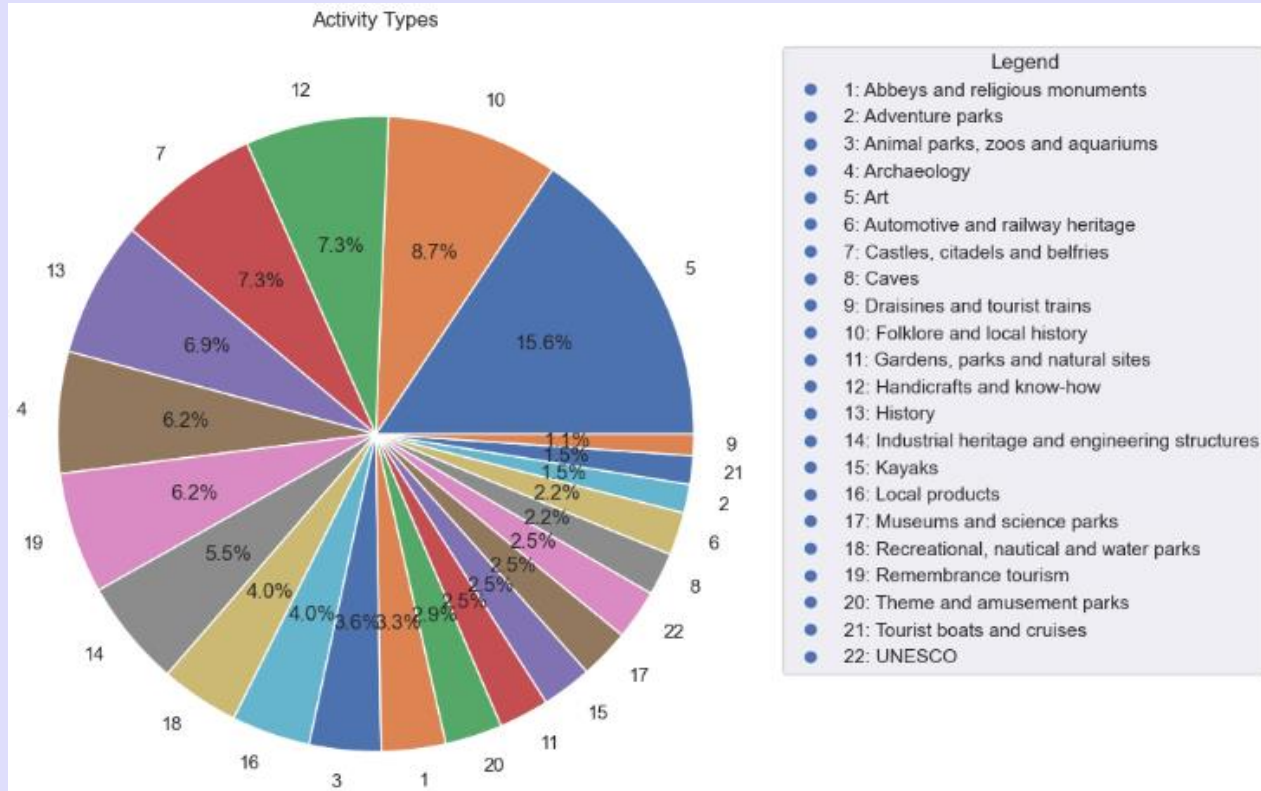
BOXPLOT GROUPED BY TYPE



BOXPLOT GROUPED BY TYPE



PIE CHART





04

MACHINE LEARNING

LINEAR REGRESSION



R-squared: 0.338

Test RMSE: 7.148

Train RMSE: 2.617

MODEL REGRESSION 1

$$\text{Adult Price} = \beta_0 + \beta_1 * \text{Activity Type} + \beta_2 * \text{City} + \epsilon$$

Where:

- Adult Price represents the dependent variable, which is the price for adults.
- Activity Type and City are the independent variables.
- β_0 , β_1 , and β_2 are the coefficients or parameters that need to be estimated.
- ϵ represents the error term, which accounts for unexplained variation in the dependent variable.

LINEAR REGRESSION



R-squared: 0.749

Test RMSE: 4.401

Train RMSE: 1.918

MODEL REGRESSION 2

$$\text{Adult Price} = \beta_0 + \beta_1 * \text{Activity Type} + \beta_2 * \text{City} + \beta_3 * \text{Children Price} + \epsilon$$

Where:

- Adult Price represents the dependent variable, which is the price for adults.
- Activity Type, City, and Children price are the independent variables.
- β_0 , β_1 , β_2 and β_3 are the coefficients or parameters that need to be estimated.
- ϵ represents the error term, which accounts for unexplained variation in the dependent variable.

DECISION TREE REGRESSION



R-squared: 0.789

Test RMSE: 4.034

Train RMSE: 0.783

MODEL REGRESSION 3

$$\text{Adult Price} = f(\text{Activity Type, City, Children Price}) + \epsilon$$

f represents the decision tree regression model, which learns how to make predictions based on the values of the independent variables. The ϵ term represents the error or residual, which captures the unexplained variation in the dependent variable.

In summary, the second model outperforms the first model in terms of train RMSE, test RMSE, and R-squared. Therefore, the second model is likely to be a better choice for predicting the target variable based on the given results. However, we could try in further research to add more data to the dataframe or to test other models to have better predictions.



05 CONCLUSION



The background of the slide features several light blue geometric shapes, including triangles and polygons, arranged in a modern, abstract pattern. A central light blue rectangular box contains the main text.

THANK YOU

Questions ?