

# Intro

## *Project Description*

The Project's objective is to scrape products from Amazon using the **Robot Framework**. It involves data processing and analysis, as well as implementing **machine learning** techniques to predict prices based on review counts and rating stars. A user interface has been provided.

**This manual consists of six slides:**

1. **Intro** (current): Manual Structure, Project Files Description
2. **User Interface**: Showcasing initial user input and the product dashboard (the end product)
3. **Sequential Explanation of the Code Execution Steps**: shows the functions used and provides pseudo code, serving as a general guide to understanding the project's execution process.
4. **Robotic Part**: Further explanation of the robotic part of the project and the associated code.
5. **Machine Learning Part**: Further explanation of the machine learning part of the project and the relevant code.
6. **Full System Representation**: Provides a detailed understanding of the entire project process and its interdependencies, should you desire a comprehensive view.

The Project Manual abstracts the code, focusing on the relevant parts to aid comprehension. However, the full code can be found in the respective files. The goal is to present the information in an easy-to-understand manner. **PS**: Zooming might be required.

**The project folder contains the following files:**

- **README\_Project\_Manual.pdf** (current)
- **front\_interface.py**: main file that integrates all the components of the project, including the execution of machine learning algorithms
- **amzn\_product\_scraping.txt**: robot script responsible for scraping products from Amazon
- **amzn\_product\_scraping\_modified.txt**: modified version of the robot script, explained in detail in the project manual
  - **udf\_robot.py**: A user-defined function required for time randomization
  - **output.xml**: The log file generated by executing the robot script
- **amzn\_data\_transformation.py**: handles data processing and transformation
- **streamlit\_subprocess\_amzn\_scraping\_script.sh**: shell script that sequentially executes the modified robot script and the python data transformation script
- **product\_table.csv**: A CSV file containing the scraped product data

## Product Scraping

With product do you want to scrap? Sample: headphones

usb c cable

Scrape Data

## Product Dashboard

Searched Product

usb c cable

Products Scraped

466

Average Price

13.24€

Most Expensive

98.49€

Cheapest

2.92€

## Product Table

Update Data

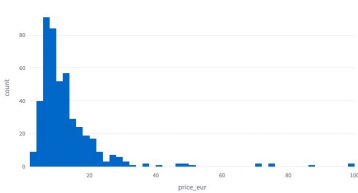
	product_id	title	price_eur	review_count	rating_stars
24	product_25	AINOPE USB C Kabel, 2 S	8.05	88,296	4.6
311	product_312	Amazon Basics USB Typ	6.7	86,563	4.6
6	product_7	Amazon Basics USB Typ	13.26	86,563	4.6
71	product_72	Anker Premium 180 cm l	10.07	82,896	4.7
53	product_54	Amazon Basics Verbindu	7.25	37,810	4.6
0	product_1	Amazon Basics Verbindu	7.25	37,810	4.6
20	product_21	INIU USB C Kabel, Ladek	10.07	31,006	4.6
416	product_417	USB C Kabel RAMPOW, if	5.09	29,939	4.6
26	product_27	USB C Kabel RAMPOW, if	6.99	29,939	4.6
21	product_22	USB C Kabel RAMPOW, if	5.09	29,939	4.6

### Trends in Product Data

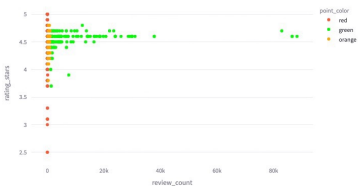
- The majority of products in the data have a high review count, indicating that they are popular among customers.
- The products have a high average rating of 4.6 stars, suggesting that they are well-received by customers.
- There is a wide range of price points, with some products priced as low as 4.63 EUR and others as high as 13.26 EUR, providing options for different budget ranges.

Ask AI anything about your scraped products (first 30 rows)

### Price Distribution



### Correlation Scatter Plot



## Machine Learning Performance

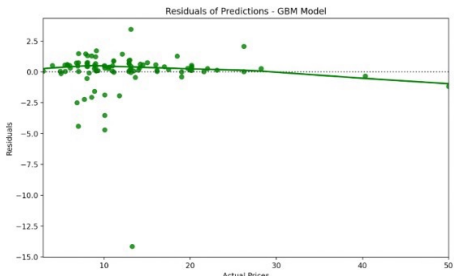
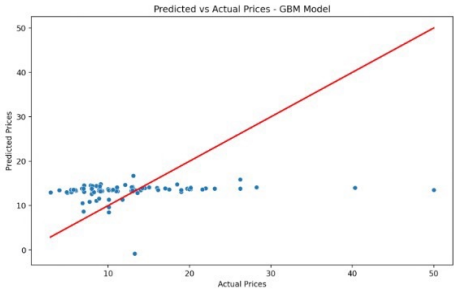
### Model: Linear Regression

Mean Squared Error

51.6017

R-Squared

-0.003



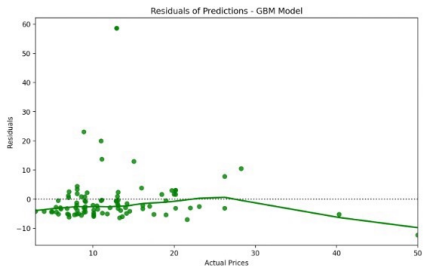
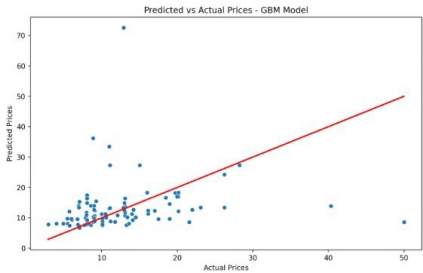
### Model: Gradient Boosting

Mean Squared Error

137.404

R-Squared

-1.6709



```

def welcome():...
def scrape_products(search_term):...
def product_dashboard(search_term):...
def perform_machine_learning(df):...    execution inside of product_dashboard fct.
def perform_linear_regression(X_train, X_test, y_train, y_test):...
def perform_gradient_boosting(X_train, X_test, y_train, y_test):...
def ask_ai_for_trends(product_data_string):...
def ask_ai_anything(product_data_string, user_prompt):...

```

Pseudo code of the project execution process

```

2  # 1. User types in a product that he wants to search
3  search_term = st.text_input('Which product do you want to scrape? Example: headphones')
4
5  # 2. product is being scraped from Amazon (with the help of Robot Framework)
6  scrape_products(search_term)    Further explanation below (Robotic Part)
7
8  # 3. after scraping is completed, dashboard is being created with the given searched product
9  product_dashboard(search_term)
10
11 # 4. machine learning is then executed inside of the product_dashboard
12 perform_machine_learning(df)
13
14 # specific algorithms being applied    Further explanation below (Machine Learning Part)
15 perform_linear_regression()
16 perform_gradient_boosting()
17
18 # additional functionality and features (inside the product dashboard)
19
20 # AI Section
21 ask_ai_for_trends() # Generative Pretrained Transformer returns some trends in the product set
22 ask_ai_anything()  # user AI prompt input related to the scraped product set
23                   # example: list the companies of the products listed
24
25 # Price Distribution Histogram
26 fig = px.histogram(df, x='price_eur')
27
28 # Correlation Scatter Plot
29 fig2 = px.scatter(df, x='review_count', y='rating_stars', color='point_color', color_discrete_map=color_scale)
30
31 # Sidebar with DataFrame Filtering Options
32 st.sidebar.header('Filters')
33 min_price = st.sidebar.number_input('Min Price', min_value=0, key=1, on_change=filter_data)
34 max_price = st.sidebar.number_input('Max Price', min_value=0, key=2, on_change=filter_data)
35 min_review_count = st.sidebar.number_input('Min Review Count', min_value=0, key=3)
36 min_stars = st.sidebar.number_input('Min Stars', min_value=1, max_value=5)
37 max_stars = st.sidebar.number_input('Max Stars', min_value=1, max_value=5)

```

# Robotic Part

When welcome() function is executed, this is what happens:

```
st.title('Product Scraping')
search_term = st.text_input('Which product do you want to scrape? Example: headphones')
st.session_state.search_term = search_term

if st.button("Scrape Data"):
    if len(search_term) != 0:
        st.write(f'Started scraping data for {search_term}...')
        scrape_products(search_term) # running the shell script, robot framework, and python data processing script
        st.write('Scraping Completed ✅')

        time.sleep(2)

    # creating dashboard
    product_dashboard(search_term)
    else:
        st.write('Please enter a search term.')
```

```
def scrape_products(search_term):
    # opening the robot script to replace <<replace_string>> with the actual searched product
    robot_script = open('amzn_product_scraping.txt', 'r').read()
    robot_script = robot_script.replace('<<replace_string>>', search_term)

    # writing modified script to a temporary file
    with open('amzn_product_scraping_modified.txt', 'w') as file:
        file.write(robot_script)

    # running the script which executes the modified robot script and python script
    subprocess.run("./streamlit_subprocess_amzn_scraping_script.sh", shell=True)
```

1. robot amzn\_product\_scraping\_modified.txt
2. python amzn\_data\_transformation.py
- # first it runs the modified robot script with the actual searched product which scrapes product data from Amazon
- # when scraping is completed it executes the python script which grabs the scraped product data from the output.xml log
- # in the log file it finds title, price, review count and star rating information for each product
- # then data processing is performed, saved into a dictionary which is transformed to a pandas dataframe
- # finally the dataframe is saved as a CSV file named product\_table.csv
- # only after this step the final dashboard is created and machine learning on the dataset is performed and results shown

## Robot Script - amzn\_product\_scraping.txt

```
*** Settings ***
Documentation
Library SeleniumLibrary
Test Setup
Test Teardown
Library

*** Variables ***
${SEARCH_TERM}
${URL}
${BROWSER}
${NUM_PAGES}

*** Keywords ***
Pause Inbetween
    Sleep Random Pause Length
    Log Pause length: ${PAUSE_LENGTH} #Pauses between Page Flips to mimic human behaviour

Go To Page
    [Arguments]
    Go To ${page_number}
    Log ${URL}&page=${page_number}
    Page flipped successfully

*** Test Cases ***
Perform Scraping
    Open Browser ${URL} ${BROWSER}
    Maximize Browser Window
    Pause Inbetween

    FOR ${page_num} IN RANGE 2 ${NUM_PAGES} # iterate through amazon pages
        Pause Inbetween-
        Wait Until Element Is Visible xpath=//div[starts-with(@el_widget_id, 'MAIN-SEARCH_RESULTS')] #ELEMENT VISIBLE # this element is a container of each individual product listing
        @element_list = Get WebElements xpath=//div[starts-with(@el_widget_id, 'MAIN-SEARCH_RESULTS')] # each product container is identified by a @el_widget_id starting with 'MAIN-SEARCH-RESULTS'
        FOR ${element} IN @element_list
            ${html}= Get Element Attribute ${element} outerHTML # we save each found product container as an element into our element list
            Log ${html} # every box container containing product listing information is logged to output.xml which is later handed over to python script to perform data processing
        END
    END

    Go To Page ${page_num}

END

[Teardown] Close Browser
```

→ amzn\_product\_scraping\_modified.txt

```
*** Variables ***
${SEARCH_TERM} usb c cable
${URL} https://www.amazon.de/s?k=${SEARCH_TERM}
```

# Machine Learning Part # Goal is to predict the price based on the review count and rating stars of a product

```
st.title('Product Scraping')
search_term = st.text_input('Which product do you want to scrape? Example: headphones')
st.session_state.search_term = search_term

if st.button('Scrape Data'):
    if len(search_term) != 0:
        st.write(f'Started scraping data for {search_term}...')
        scrape_products(search_term) # running the shell script, robot framework, and python data processing script
        st.write('Scraping Completed ✅')

        time.sleep(2)

    # creating dashboard
    product_dashboard(search_term)
    else:
        st.write('Please enter a search term.')
```

def product\_dashboard(search\_term):

```
def load_data():
    data = pd.read_csv('product_table.csv')
    return data
```

df = load\_data()

⋮  
⋮ → Code inbetween is hidden for demonstration purposes (focus on Machine Learning)  
⋮

# -----MACHINE LEARNING -----

perform\_machine\_learning(df)

```
def perform_machine_learning(df):
    """Goal is to predict the price based on the review count and rating stars of a product"""
    st.header('Machine Learning Performance')

    # subset dataframe for machine learning
    df_for_ml = df[['price_eur', 'review_count', 'rating_stars']]

    # feature and target variable selection
    X = df_for_ml[['review_count', 'rating_stars']]
    y = df_for_ml['price_eur']

    X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=0)
```

## ⋮ Linear Regression

```
# LINEAR REGRESSION
st.subheader('Model: Linear Regression')
mse_linreg, r2_linreg, y_pred_linreg = perform_linear_regression(X_train, X_test, y_train, y_test)

# show model performance metrics
st.metric('Mean Squared Error', round(mse_linreg, 4))
st.metric('R-Squared', round(r2_linreg, 4))

# embedding predicted vs actual price plot for linear regression
plt.figure(figsize=(10, 6))
sns.scatterplot(x=y_test, y=y_pred_linreg)
plt.title('Predicted vs Actual Prices - Linear Regression Model')
plt.xlabel('Actual Prices')
plt.ylabel('Predicted Prices')
plt.plot([y_test.min(), y_test.max()], [y_test.min(), y_test.max()], color='red',
         lw=2)
st.pyplot(plt)

# embedding residual plot
plt.figure(figsize=(10, 6))
sns.residplot(x=y_test, y=y_pred_linreg, lowess=True, color='g')
plt.title('Residuals of Predictions - Linear Regression Model')
plt.xlabel('Actual Prices')
plt.ylabel('Residuals')
st.pyplot(plt)
```

```
def perform_linear_regression(X_train, X_test, y_train, y_test):
    linreg_model = LinearRegression()

    # train the model
    linreg_model.fit(X_train, y_train)

    # predict prices
    y_pred_linreg = linreg_model.predict(X_test)

    # evaluate the model
    mse_linreg = mean_squared_error(y_test, y_pred_linreg)
    r2_linreg = r2_score(y_test, y_pred_linreg)

    return [mse_linreg, r2_linreg, y_test, y_pred_linreg]
```

MSE (Mean Squared Error):  
Measures average squared prediction error magnitude.

R2 (R-Squared):  
coefficient of determination, how close the data is to the fitted regression line, helps assess goodness of fit

Linear Regression:  
MSE: 81.60, R-squared: -0.003  
Moderately inaccurate, poor explanatory power

Gradient Boosting:  
MSE: 137.40, R-squared: -1.6709  
Highly inaccurate, very poor fit.

## Packages for ML

```
# modules for machine learning
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression
from sklearn.ensemble import GradientBoostingRegressor
from sklearn.metrics import mean_absolute_error, mean_squared_error, r2_score
```

## ⋮ Gradient Boosting

```
# GRADIENT BOOSTING
st.subheader('Model: Gradient Boosting')
mse_gbm, r2_gbm, y_test, y_pred_gbm = perform_gradient_boosting(X_train, X_test, y_train, y_test)

# show model performance metrics
st.metric('Mean Squared Error', round(mse_gbm, 4))
st.metric('R-Squared', round(r2_gbm, 4))

# embedding predicted vs actual price plot for gradient boosting
plt.figure(figsize=(10, 6))
sns.scatterplot(x=y_test, y=y_pred_gbm)
plt.title('Predicted vs Actual Prices - GBM Model')
plt.xlabel('Actual Prices')
plt.ylabel('Predicted Prices')
plt.plot([y_test.min(), y_test.max()], [y_test.min(), y_test.max()], color='red',
         lw=2)
st.pyplot(plt)

# embedding residual plot
plt.figure(figsize=(10, 6))
sns.residplot(x=y_test, y=y_pred_gbm, lowess=True, color='g')
plt.title('Residuals of Predictions - GBM Model')
plt.xlabel('Actual Prices')
plt.ylabel('Residuals')
st.pyplot(plt)
```

```
def perform_gradient_boosting(X_train, X_test, y_train, y_test):
    gbm_model = GradientBoostingRegressor()

    # training the model
    gbm_model.fit(X_train, y_train)

    # predict prices
    y_pred_gbm = gbm_model.predict(X_test)

    # model evaluation
    mse_gbm = mean_squared_error(y_test, y_pred_gbm)
    r2_gbm = r2_score(y_test, y_pred_gbm)

    return [mse_gbm, r2_gbm, y_test, y_pred_gbm]
```

## Product Scraping

Which product do you want to scrape? Example: headphones

usb c cable

Scrape Data

```
27 st.title('Product Scraping')
28 search_term = st.text_input('Which product do you want to scrape? Example: headphones')
29 st.session_state.search_term = search_term
30
31 if st.button('Scrape Data'):
32     if len(search_term) != 0:
33         st.write('Started scraping data for {search_term}...')
34         scrape_products(search_term) # running the script, robot framework and scraping products
35         st.write('Scraping Completed')
36
37         time.sleep(2)
38
39         # creating dashboard
40         product_dashboard(search_term)
41     else:
42         st.write('Please enter a search term.')
```

Scraping Completed

## Product Dashboard

Searched Product  
usb c cable

Products Scraped  
466

Average Price  
13.24€

Most Expensive  
98.49€

Cheapest  
2.92€

product_id	name	price_eur	review_count	rating_star
1	AmazonBasics USB Type-C to USB-A 3.1 Gen 1 Cable, 6.6 ft (2.0 m), Black, AmazonBasics	6.99	86,236	4.6
2	AmazonBasics USB Type-C to USB-A 3.1 Gen 1 Cable, 6.6 ft (2.0 m), Black, AmazonBasics	6.7	86,043	4.6
3	AmazonBasics USB Type-C to USB-A 3.1 Gen 1 Cable, 6.6 ft (2.0 m), Black, AmazonBasics	13.24	86,043	4.6
4	AmazonBasics USB Type-C to USB-A 3.1 Gen 1 Cable, 6.6 ft (2.0 m), Black, AmazonBasics	10.87	82,096	4.7
5	AmazonBasics USB Type-C to USB-A 3.1 Gen 1 Cable, 6.6 ft (2.0 m), Black, AmazonBasics	1.29	87,020	4.6
6	AmazonBasics USB Type-C to USB-A 3.1 Gen 1 Cable, 6.6 ft (2.0 m), Black, AmazonBasics	7.25	37,046	4.6
7	AmazonBasics USB Type-C to USB-A 3.1 Gen 1 Cable, 6.6 ft (2.0 m), Black, AmazonBasics	10.87	31,006	4.6
8	AmazonBasics USB Type-C to USB-A 3.1 Gen 1 Cable, 6.6 ft (2.0 m), Black, AmazonBasics	9.99	25,999	4.6

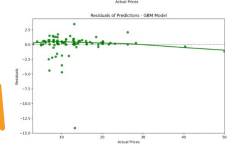
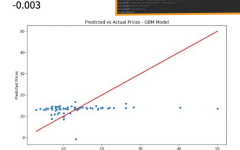
```
def product_dashboard(search_term):
    def load_data():
        data = pd.read_csv('product_table.csv')
        return data
    # Reading in the dataframe
    df = load_data()
    def perform_machine_learning(df):...
```

### Machine Learning Performance

Model: Linear Regression

R-squared  
51.6017

Adjusted  
-0.003



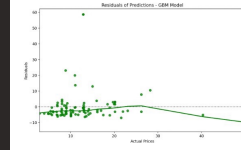
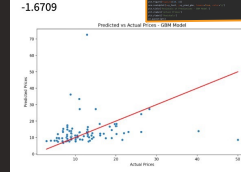
```
# subset dataframe for machine learning
df_for_ml = df[['price_eur', 'review_count', 'rating_star']]
# feature and target variable selection
X = df_for_ml[['review_count', 'rating_star']]
y = df_for_ml['price_eur']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=0)
```

```
from sklearn.linear_model import LinearRegression
# train to predict the price based on the review count and rating stars of a product
linear_model = LinearRegression()
# train the model
linear_model.fit(X_train, y_train)
# Predict prices
y_pred_linear = linear_model.predict(X_test)
# Evaluate the model
mse_linear = mean_squared_error(y_test, y_pred_linear)
r2_linear = r2_score(y_test, y_pred_linear)
return [mse_linear, r2_linear, y_test, y_pred_linear]
```

### Model: Gradient Boosting

Mean Squared Error  
137.404

R-squared  
-1.6709



```
from sklearn.ensemble import GradientBoostingRegressor
# train to predict the price based on the review count and rating stars of a product
gbm_model = GradientBoostingRegressor()
# train the model
gbm_model.fit(X_train, y_train)
# Predicting prices
y_pred_gbm = gbm_model.predict(X_test)
# Model evaluation
mse_gbm = mean_squared_error(y_test, y_pred_gbm)
r2_gbm = r2_score(y_test, y_pred_gbm)
return [mse_gbm, r2_gbm, y_test, y_pred_gbm]
```

```
def scrape_products(search_term):
    robot_script = open('amazon_product_scraping.txt', 'r').read()
    robot_script = robot_script.replace('<<replace_string>>', search_term)
    print(robot_script)
    # writing modified script to a temporary file
    with open('amazon_product_scraping_modified.txt', 'w') as file:
        file.write(robot_script)
    # running the script, robot framework and scraping products
    subprocess.run('./streamlit subprocess_amzn_scraping_script.sh', shell=True)
```

robot\_amzn\_product\_scraping\_modified.txt  
python amzn\_data\_transformation.py

```
import xml.etree.ElementTree as ET
from lxml import html
import pandas as pd
import json
# Parse the XML file
tree = ET.parse('output.xml')
root = tree.getroot()
product_dict = {}
# Loop through the product elements and extract data
for product in root.findall('.//product'):
    # Extract product details
    product_id = product.get('id')
    name = product.get('name')
    price_eur = product.get('price_eur')
    review_count = product.get('review_count')
    rating_star = product.get('rating_star')
    product_dict[product_id] = {'name': name, 'price_eur': price_eur, 'review_count': review_count, 'rating_star': rating_star}
```

```
df = pd.DataFrame(product_dict)
df = df.reset_index()
df = df.rename(columns={'index': 'product_id'})
df.to_csv('product_table.csv', index=False)
```

## Additional Functionality and Features

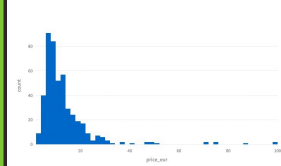
### Trends in Product Data

- USB C cables are popular in the product dataset, with multiple listings featuring USB C cables from different brands.
- The average price of the products is around 10.0 EUR, with some outliers priced higher or lower.
- The majority of the products have high review counts, indicating a high level of customer engagement and satisfaction.

Ask AI anything about your scraped products (first 30 rows)

test the companies of the product

### Price Distribution



AI Prompting

Histogram

Correlation Scatter Plot