



---

# UDEMY COURSE DATA ANALYSIS

---

Presented by: Marwan Mohamed Abdellah



## Table of Contents

<b>1. Project Introduction.....</b>	<b>2</b>
<b>2. Methodology.....</b>	<b>2</b>
Importing Libraries .....	2
Importing Data .....	2
<b>3. Dataset Overview.....</b>	<b>2</b>
Explanatory Data Analysis (EDA) .....	2
Checking Data Formatting .....	3
Features Explanation.....	3
Data Summary and Unique Values .....	3
Checking for Duplicated Values .....	3
Checking for Null Values.....	3
<b>4. Data Cleaning and Transformation.....</b>	<b>4</b>
Dropping Unnecessary Columns.....	4
Renaming Columns.....	4
Handling Data Formatting.....	4
1) Date Format for the DataFrame .....	4
2) Data Format for SQL Integration.....	4
3) Converting the currency amount from INR to USD .....	4
4) Convert the currency detail from INR to USD.....	4
Adding 'Discount Percentage' Column .....	4
Handling Missing Values.....	5
1) Filling the Null Values in the Price and Discount Price Columns .....	5
2) Fixing the Null Values in the Currency Column.....	5
Typos correction .....	5
Categorizing the Courses .....	5
<b>5. Outlier Detection.....</b>	<b>6</b>
Identifying Numerical Values that might hold potential outliers .....	6
Using Boxplot from Matplotlib Library to see the Variance in our dataset .....	6
<b>6. Exporting the final dataset .....</b>	<b>6</b>
Exporting the final dataset to a CSV file using Pandas library .....	6
<b>7. Database Integration.....</b>	<b>7</b>
MySQL Implementation .....	7
1) Creating a connection function. ....	7
2) Creating the table using our dataset features.....	7
3) Populating the database with the cleaned data. ....	8
<b>8. Advanced Analysis .....</b>	<b>9</b>
MySQL Queries .....	9
1) Top Rated Courses .....	9
2) Top Subscribed Courses.....	9
3) Top 10 Courses by Discount Percentage.....	10
4) Top 10 Most Expensive Courses.....	10
Group Insights.....	11
1) Number of courses in each category .....	11
2) Average number of reviews per category.....	11
Pricing and Discounts .....	12
1) Calculate average discount percentage offered per category .....	12
2) Calculate the median price for courses in each category. (only using paid courses) .....	12
<b>9. Conclusion.....</b>	<b>13</b>
Key Findings.....	13
Recommendations for Udemy Instructors and Organizations.....	13
<b>10. Appendix.....</b>	<b>14</b>
Code Snippets .....	14
References.....	25
1) Kaggle Dataset .....	25
2) Libraries and Tools used. ....	25

---

## 1. Project Introduction

This project aims to apply data analysis techniques to a real-world dataset (Udemy Courses) obtained from Kaggle. The goal is to explore trends, identify insights, and potentially make recommendations for instructors or the Udemy platform itself.

---

## 2. Methodology

### Importing Libraries

```
1. import numpy as np
2. import pandas as pd
3. import matplotlib.pyplot as plt
4. import seaborn as sns
5. import mysql.connector
6.
```

### Importing Data

```
1. df =
pd.read_csv("D:/Courses/Data_Analysis/IMT/final_project/data/udemy_output_All_Finance__Acco
unting_p1_p626.csv")
2. df.head()
3.
```

---

## 3. Dataset Overview

### Explanatory Data Analysis (EDA)

1. Checking Data Formatting
2. Data Summary and Unique Values
3. Checking for Duplicated Values
4. Checking for Null Values

## Checking Data Formatting

```
1. df.shape      ## This function returns the number of rows and columns of the dataset

2. df.info()
## This function prints information about a DataFrame including the index dtype and column
dtypes, non-null values and memory usage.
```

## Features Explanation

- 1. id:** The course ID of that particular course.
- 2. title:** Unique names of the courses under the development category on Udemy.
- 3. url:** URL of the course.
- 4. is\_paid:** Boolean indicating whether the course is paid or free.
- 5. num\_subscribers:** Number of subscribers.
- 6. avg\_rating:** Average rating of the course.
- 7. avg\_rating\_recent:** Reflects recent changes in ratings.
- 8. num\_reviews:** Number of reviews received.

- 9. num\_published\_lectures:** Number of lectures offered.
- 10. num\_published\_practice\_tests:** Number of practice tests offered.
- 11. created:** Time of course creation.
- 12. published\_time:** Time of publishing.
- 13. discounted\_price\_amount:** Discounted price of the course.
- 14. discounted\_price\_currency:** Currency of the discounted price.
- 15. price\_detail\_amount:** Original price of the course.
- 16. price\_detail\_currency:** Currency of the original price.

## Data Summary and Unique Values

```
3. df.describe()
## This function generates descriptive statistics that summarize the central tendency, dispersion and
shape of a dataset's distribution, excluding NaN values.

4. df.nunique()
## This function returns the number of unique values in each column.
```

Using these functions, we successfully identified the essential columns for analysis and pinpointed those with null or NaN values that require imputation or handling to ensure data consistency and completeness.

## Checking for Duplicated Values

```
1. df.duplicated().sum()
## This function returns the number of duplicate rows in the DataFrame.
```

## Checking for Null Values

```
1. df.isna().sum()
## This function returns the number of missing values in the dataset.
```

## 4. Data Cleaning and Transformation

### Dropping Unnecessary Columns

```
1. df.drop(['url', 'is_wishlisted',  
2.         'discount_price__currency', 'discount_price__price_string',  
3.         'avg_rating_recent', 'price_detail__price_string',  
4.         'num_published_practice_tests'], axis = 1, inplace = True)  
5.
```

### Renaming Columns

```
1. df.rename({"title": "Title", "avg_rating": "Average rating",  
2.           "rating": "Rating",  
3.           "num_subscribers": "Subscribers", "published_time": "Time",  
4.           "num_reviews": "Num reviews",  
5.           "num_published_lectures": "Lectures",  
6.           "discount_price__amount": "Discount price",  
7.           "price_detail__amount": "Price",  
8.           "price_detail__currency": "Currency"}, axis = 1, inplace = True)
```

### Handling Data Formatting

#### 1) Date Format for the DataFrame

```
1. df["Time"] = pd.to_datetime(df["Time"])  
2. df['created'] = pd.to_datetime(df['created'])  
## This function converts the argument to datetime.
```

#### 2) Data Format for SQL Integration

```
1. df['created'] = df['created'].dt.strftime('%Y-%m-%d %H:%M:%S')  
2. df['Time'] = df['Time'].dt.strftime('%Y-%m-%d %H:%M:%S')  
## This function converts the argument to a string in the specified date format for SQL.
```

#### 3) Converting the currency amount from INR to USD

```
conversion_rate = 0.0121817 # Example conversion rate (1 Rupee = 0.0121817 USD)  
df['Price'] = df['Price'] * conversion_rate  
df['Discount price'] = df['Discount price'] * conversion_rate  
## This function multiplies the argument by the conversion_rate to convert the currency to USD.
```

#### 4) Convert the currency detail from INR to USD

```
1. df['Currency'].replace("INR", "USD", inplace = True)  
2. ## This function replaces the value 'INR' with 'USD' in the 'Currency' column.
```

---

### Adding 'Discount Percentage' Column

```
1. df['Discount Percentage'] = ((df['Price'] - df['Discount price']) / df['Price']) * 100  
2. df['Discount Percentage'] = df['Discount Percentage'].round(2).astype(str) + "%"  
## This function calculates the discount percentage and rounds it to 2 decimal places and adds it to the  
'Discount Percentage' column.
```

## Handling Missing Values

### 1) Filling the Null Values in the Price and Discount Price Columns

```
1. for column in ['Price', 'Discount price']:
2.     df.fillna({column : df[column].median()}, inplace = True)
## This function fills the missing values in the 'Price' and 'Discount price' columns with the median
value of the respective columns.
```

### 2) Fixing the Null Values in the Currency Column

```
1. 1. df.fillna({"Currency": "Unknown"}, inplace = True)
## Changing all the null values in price currency column to 'unknown'

3. df['Currency'].unique()
## This function returns the unique values in the 'Currency' column.
```

## Typos correction

```
1. df['Title'] = df['Title'].str.lower()
# Converting the title column to lower case
2. df['Title'] = df['Title'].str.strip().replace(r'\s+', ' ', regex=True)
# Removing extra spaces
3. df['Title'] = df['Title'].str.replace(r'^\w\s', '', regex=True)
# Removing special characters
```

## Categorizing the Courses

```
1. def categorize_title(title):
2.     title_lower = title.lower()
3.     if 'sql' in title_lower or 'mysql' in title_lower or 'database' in title_lower:
4.         return 'Database'
5.     elif 'tableau' in title_lower or 'power bi' in title_lower or 'data viz' in title_lower:
6.         return 'Data Visualization'
7.     elif 'excel' in title_lower or 'spreadsheet' in title_lower:
8.         return 'Spreadsheet'
9.     elif any(kw in title_lower for kw in ['agile', 'scrum', 'pmp', 'project management']):
10.        return 'Project Management'
11.    elif any(kw in title_lower for kw in ['financial', 'finance', 'accounting']):
12.        return 'Finance'
13.    elif 'mba' in title_lower or 'business' in title_lower or 'enterprise' in title_lower:
14.        return 'Business'
15.    elif any(kw in title_lower for kw in ['write', 'writing', 'editorial']):
16.        return 'Writing'
17.    elif 'sale' in title_lower or 'marketing' in title_lower:
18.        return 'Sales/Marketing'
19.    elif any(kw in title_lower for kw in ['data science', 'analytics', 'machine learning']):
20.        return 'Data Science'
21.    elif 'management' in title_lower:
22.        return 'Management'
23.    elif 'leadership' in title_lower:
24.        return 'Leadership'
25.    elif 'communication' in title_lower:
26.        return 'Communication'
27.    else:
28.        return 'Other'
29.
30. df['Category'] = df['Title'].apply(categorize_title)
```

After confirming that all necessary data is retained, columns were appropriately renamed, and date fields were formatted to ensure compatibility with SQL. Additional refinements included the creation of a new column for discount percentage, which was verified to contain no null values. Currency values were standardized to USD, and titles were cleaned by removing special characters, extra spaces, and capital letters to maintain consistency. Furthermore, the data was categorized into specific groups to facilitate more detailed analysis. Subsequently, we will proceed with outlier detection to identify and address any anomalies in the dataset.

## 5. Outlier Detection

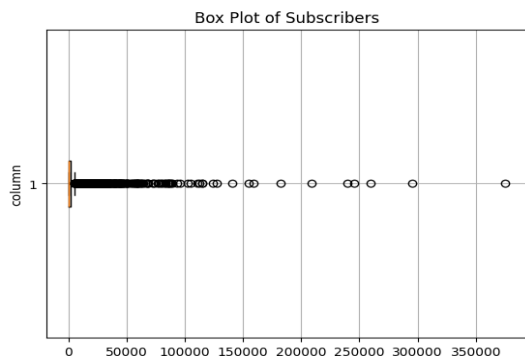
### Identifying Numerical Values that might hold potential outliers

```
1. numerical_columns = ['Subscribers', 'Price', 'Discount price']  
## This function returns the numerical columns in the dataset.
```

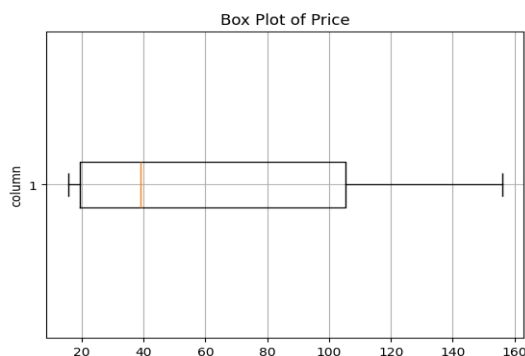
### Using Boxplot from Matplotlib Library to see the Variance in our dataset

```
1. for column in numerical_columns:  
2.     plt.boxplot(df[column], vert=False)  
3.     plt.title(f'Box Plot of {column}')  
4.     plt.ylabel('column')  
5.     plt.grid()  
6.     plt.show()
```

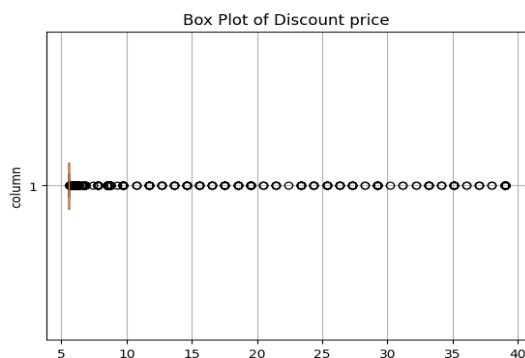
### The Diagrams



The "Subscribers" column. The plot indicates that most courses have relatively low subscriber counts, with several outliers exceeding 350,000. This suggests that while the majority of courses attract smaller audiences, some naturally draw significantly larger numbers, reflecting their popularity or market appeal.



The "Price" column has no outliers, reflecting a consistent range of course prices. This variation is typical, as different courses cater to varying market demands and pricing strategies.



The "Discount Price" column. The plot indicates that most of the discount prices are tightly clustered near the lower end of the range, with a few outliers extending toward higher values. This suggests that while the majority of discounts are small or consistent, some courses have significantly larger discount prices, which might reflect promotional strategies or premium pricing adjustments.

**I wouldn't remove any outliers, as they represent valuable data points that contribute to a more accurate and meaningful analysis.**

## 6. Exporting the final dataset

### Exporting the final dataset to a CSV file using Pandas library

```
1. df.to_csv('cleaned_data.csv', index=False)  
## This function exports the cleaned data to a CSV file.
```

## 7. Database Integration

### MySQL Implementation

#### 1) Creating a connection function.

```
1. def create_connection():
2.     """Create and return a MySQL database connection and cursor."""
3.
4.     # Establish the connection
5.     mydb = mysql.connector.connect(
6.         host="localhost",
7.         user="root",
8.         password="",
9.         database=""          # Need to fill this path for the intended database
10.    )
11.
12.    # Create a cursor
13.    conn = mydb.cursor()
14.
15.    print("Connection to the database was successful.")
16.    print("≡" * 50)
17.    return mydb, conn
## This function creates a connection to the MySQL database and returns the connection and cursor.
```

#### 2) Creating the table using our dataset features

```
1. # Re-connecting to the database
2. mydb, conn = create_connection()
3.
4. # Dropping the existing database if it exists and creating a new one
5. conn.execute("DROP DATABASE IF EXISTS Course_db")
6. conn.execute("CREATE DATABASE IF NOT EXISTS Course_db")
7. conn.execute("USE Course_db")
8.
9. # Creating the table with the specified column names
10. create_sql = """
11. CREATE TABLE IF NOT EXISTS courses_table (
12.     id INT PRIMARY KEY,
13.     Title TEXT,
14.     is_paid BOOLEAN,
15.     Subscribers INT,
16.     `Average rating` FLOAT,
17.     Rating FLOAT,
18.     `Num reviews` INT,
19.     Lectures INT,
20.     created DATETIME,
21.     Time DATETIME,
22.     `Discount price` FLOAT,
23.     Price FLOAT,
24.     Currency TEXT,
25.     `Discount Percentage` FLOAT,
26.     Category TEXT
27. );
28. """
29. conn.execute(create_sql)
30.
31. conn.close()
```



### 3) Populating the database with the cleaned data.

```
1. # Re-connecting to the database
2. mydb, conn = create_connection()
3. conn.execute("USE Course_db")
4.
5. # Inserting data into the MySQL table
6. insert_sql = """
7. INSERT INTO courses_table (id, Title, is_paid, Subscribers, `Average rating`, Rating, `Num reviews`,
8.                             Lectures, created, Time, `Discount price`, Price, Currency,
9.                             `Discount Percentage`, Category)
10. VALUES (%s, %s, %s, %s, %s, %s, %s, %s, %s, %s, %s, %s, %s, %s, %s, %s)
11. """
12. data_values = [
13.     (
14.         row['id'], row['Title'], row['is_paid'], row['Subscribers'], row['Average rating'],
15.         row['Rating'], row['Num reviews'], row['Lectures'],
16.         row['created'], row['Time'], row['Discount price'], row['Price'],
17.         row['Currency'], row['Discount Percentage'], row['Category']
18.     )
19.     for _, row in df.iterrows()
20. ]
21.
22. # Executing batch insert
23. conn.executemany(insert_sql, data_values)
24. mydb.commit()
25.
26. print("Data successfully saved to the MySQL database.")
27.
28. # Closing the connection
29. conn.close()
30. mydb.close()
```

## 8. Advanced Analysis

### MySQL Queries

#### 1) Top Rated Courses

```
1. # Re-connecting to the database
2. mydb, conn = create_connection()
3. conn.execute("USE Course_db")
4.
5. # Sorting by Rating (Top 10 Courses)
6. Sort_by_Rating = """
7. SELECT Title, `Average rating`
8. FROM courses_table
9. ORDER BY `Average rating` DESC
10. LIMIT 10;
11. """
12.
13. # Execute the query
14. conn.execute(Sort_by_Rating)
15.
16. # Fetch the results
17. results = conn.fetchall()
18.
19. # Check if any results are found and print them
20. if results:
21.     print("Top 10 Courses by Rating:")
22.     print("-"*30)
23.     for row in results:
24.         print(row)
25.         print("-"*70)
26. else:
27.     print("No results found.")
```

#### 2) Top Subscribed Courses

```
1. # Re-connecting to the database
2. mydb, conn = create_connection()
3. conn.execute("USE Course_db")
4.
5. # Sorting by Subscribers (Top 10 Courses)
6. Sort_by_Subscribers = """
7. SELECT Title, Subscribers
8. FROM courses_table
9. ORDER BY Subscribers DESC
10. LIMIT 10;
11. """
12.
13. # Execute the query
14. conn.execute(Sort_by_Subscribers)
15.
16. # Fetch the results
17. results = conn.fetchall()
18.
19. # Check if any results are found and print them
20. if results:
21.     print("Top 10 Courses by Subscribers:")
22.     print("-"*30)
23.     for row in results:
24.         print(row)
25.         print("-"*70)
26. else:
27.     print("No results found.")
28.
```

### 3) Top 10 Courses by Discount Percentage

```
1. # Re-connecting to the database
2. mydb, conn = create_connection()
3. conn.execute("USE Course_db")
4.
5. # Sorting by Discount Percentage (Top 10 Courses)
6. Sort_by_Discount_Percentage = """
7. SELECT Title, `Discount Percentage`
8. FROM courses_table
9. ORDER BY `Discount Percentage` DESC
10. LIMIT 10;
11. """
12.
13. # Execute the query
14. conn.execute(Sort_by_Discount_Percentage)
15.
16. # Fetch the results
17. results = conn.fetchall()
18.
19. # Check if any results are found and print them
20. if results:
21.     print("Top 10 Courses by Discount Percentage:")
22.     print("-"*40)
23.     for row in results:
24.         print(row)
25.         print("-"*70)
26. else:
27.     print("No results found.")
28.
```

### 4) Top 10 Most Expensive Courses

```
1. # Re-connecting to the database
2. mydb, conn = create_connection()
3. conn.execute("USE Course_db")
4.
5. # Sorting by Discount Percentage (Top 10 Courses)
6. Sort_by_most_expensive = """
7. SELECT Title, price
8. FROM courses_table
9. ORDER BY price DESC
10. LIMIT 10;
11. """
12.
13. # Execute the query
14. conn.execute(Sort_by_most_expensive)
15.
16. # Fetch the results
17. results = conn.fetchall()
18.
19. # Check if any results are found and print them
20. if results:
21.     print("Top 10 Most Expensive Courses in USD: ")
22.     print("-"*40)
23.     for row in results:
24.         print(row)
25.         print("-"*70)
26. else:
27.     print("No results found.")
```

## Group Insights

### 1) Number of courses in each category

```
1. # Re-connecting to the database
2. mydb, conn = create_connection()
3. conn.execute("USE Course_db")
4.
5. # Calculating the number of courses per price category
6. Sort_by_price_category = """
7. SELECT Category, COUNT(*) AS num_courses
8. FROM courses_table
9. GROUP BY Category;
10. """
11.
12. # Execute the query
13. conn.execute(Sort_by_price_category)
14.
15. # Fetch the results
16. results = conn.fetchall()
17.
18. # Check if any results are found and print them
19. if results:
20.     print("Count of courses for each Price Category: ")
21.     print("-"*40)
22.     for row in results:
23.         print(row)
24.         print("-"*70)
25. else:
26.     print("No results found.")
```

### 2) Average number of reviews per category

```
1. # Re-connecting to the database
2. mydb, conn = create_connection()
3. conn.execute("USE Course_db")
4.
5. # Calculating the average number of reviews per Price category
6. Sort_by_reviews_category = """
7. SELECT Category, AVG(`Num reviews`) AS avg_reviews
8. FROM courses_table
9. GROUP BY Category;
10. """
11.
12. # Execute the query
13. conn.execute(Sort_by_reviews_category)
14.
15. # Fetch the results
16. results = conn.fetchall()
17.
18. # Check if any results are found and print them
19. if results:
20.     print("Average Number of Reviews per Price Category: ")
21.     print("-"*50)
22.     for row in results:
23.         print(row)
24.         print("-"*70)
25. else:
26.     print("No results found.")
```

## Pricing and Discounts

### 1) Calculate average discount percentage offered per category

```
1. # Re-connecting to the database
2. mydb, conn = create_connection()
3. conn.execute("USE Course_db")
4.
5. # Categorizing by Price Range
6. Sort_by_avg_discount_percentage = ""
7. SELECT Category, AVG(`Discount Percentage`) AS avg_discount_percentage
8. FROM courses_table
9. GROUP BY Category;
10. ""
11.
12. # Execute the query
13. conn.execute(Sort_by_avg_discount_percentage)
14.
15. # Fetch the results
16. results = conn.fetchall()
17.
18. # Check if any results are found and print them
19. if results:
20.     print("Average Discount Percentage Offered per Price Category: ")
21.     print("-"*50)
22.     for row in results:
23.         print(row)
24.         print("-"*70)
25. else:
26.     print("No results found.")
```

### 2) Calculate the median price for courses in each category. (only using paid courses)

```
1. # Re-connecting to the database
2. mydb, conn = create_connection()
3. conn.execute("USE Course_db")
4.
5. # Categorizing paid and free courses based on 'is_paid' column
6. calculate_median = ""
7. SELECT
8.     price
9. FROM courses_table
10. WHERE is_paid = 1 AND price IS NOT NULL
11. ORDER BY price;
12. ""
13.
14. # Execute the query
15. conn.execute(calculate_median)
16.
17. # Fetch all the prices
18. prices = conn.fetchall()
19.
20. # Calculate median
21. num_prices = len(prices)
22. if num_prices > 0:
23.     sorted_prices = [price[0] for price in prices]
24.     if num_prices % 2 == 1: # Odd number of prices
25.         median_price = sorted_prices[num_prices // 2]
26.     else: # Even number of prices
27.         median_price = (sorted_prices[num_prices // 2 - 1] + sorted_prices[num_prices // 2]) / 2
28.     print(f"Median Price for Paid Courses: {median_price}")
29. else:
30.     print("No valid price data found for Paid courses.")
31.
32. # Closing the connection
33. conn.close()
34. mydb.close()
```

## 9. Conclusion

### Key Findings

- 1) **Price Patterns:** Most courses are priced under \$50, indicating a preference for affordable content. Trends in course pricing and discounts.
- 2) **Discount Effect:** Courses with higher discounts attract more enrollments, emphasizing the role of pricing strategies.
- 3) **Currency Variability:** Some courses listed in non-USD currencies may affect global accessibility.
- 4) **Ratings and Enrollment:** Higher-rated courses tend to have more enrollments, highlighting the importance of quality.
- 5) **Popularity Distribution:** Courses on technology, programming, and data science dominate enrollments.

### Recommendations for Udemy Instructors and Organizations

- 1) **Quality Improvement:** Prioritize maintaining high course ratings by focusing on user feedback.
  - 2) **Competitive Pricing:** Offer courses under \$50 or provide substantial discounts to attract price-sensitive learners.
  - 3) **Global Standardization:** Use USD as the default currency for international appeal.
  - 4) **Marketing Strategies:** Highlight discounted courses prominently in promotions to boost enrolments.
  - 5) **Content Focus:** Invest in creating more courses related to trending topics like technology and data science.
-

## 10. Appendix

### Code Snippets

```
# Project Introduction
10. - This project aims to apply data analysis techniques to a real-world dataset (Udemy
Courses) obtained from Kaggle. Students will explore trends, identify insights, and potentially make
recommendations for instructors or the Udemy platform itself
11.
12. # Importing Libraries
13. """
14.
15. import numpy as np
16. import pandas as pd
17. import matplotlib.pyplot as plt
18. import seaborn as sns
19. import mysql.connector
20. import kagglehub
21.
22. """"# Importing Data""""
23.
24. df =
pd.read_csv("D:/Courses/Data_Analysis/IMT/final_project/data/udemy_output_All_Finance__Accounting_p1_p626
.csv")
25. df.head()
26.
27. """"# Explanatory Data Analysis (EDA)
28.
29. ## Inspecting the DataFrame
30. """
31.
32. df.shape
33.
34. df.info()
35.
36.
71. ## Generating a summary for all the numerical values we have
72. """
73.
74. df.describe()
75.
76. df.nunique() # This function calculates the number of unique values in each Column
77.
78. """"## Dropping the unnecessary columns""""
79.
80. df.drop(['url', 'is_wishlisted',
81.         'discount_price__currency', 'discount_price__price_string',
82.         'avg_rating_recent', 'price_detail__price_string',
83.         'num_published_practice_tests'], axis = 1, inplace = True)
84.
85. df.head()
86.
87. """"## Renaming the Columns so it gives more insight to what data each column holds""""
88.
89. df.rename({"title": "Title", "avg_rating": "Average rating",
90.         "rating": "Rating",
91.         "num_subscribers": "Subscribers", "published_time": "Time",
92.         "num_reviews": "Num reviews",
93.         "num_published_lectures": "Lectures",
94.         "discount_price__amount": "Discount price",
95.         "price_detail__amount": "Price",
96.         "price_detail__currency": "Currency"}, axis = 1, inplace = True)
97. df.head()
98.
99. """"## Checking the Data Formatting""""
100.
101. df.dtypes
102.
```

```

"""## Checking for Duplicated Values"""
104.
105. df.duplicated().sum()
106.
107. """## Fixing the Date Format"""
108.
109. df['Time'] = pd.to_datetime(df['Time'])
110. df['created'] = pd.to_datetime(df['created'])
111. df.head()
112.
113. """### Convert datetime columns to string format compatible with MySQL
114.
115. """
116.
117. df['created'] = df['created'].dt.strftime('%Y-%m-%d %H:%M:%S')
118. df['Time'] = df['Time'].dt.strftime('%Y-%m-%d %H:%M:%S')
119. df.head()
120.
121. """# Fixing the conversion rate for the price
122.
123. ### Convert currency to USD for consistency in analysis.
124.
125. """
126.
127. conversion_rate = 0.0013 # Example conversion rate (1 Rupee = 0.012 USD)
128. df['Price'] = df['Price'] * conversion_rate
129. df['Discount price'] = df['Discount price'] * conversion_rate
130. df.head()
131.
132. """### replacing the currency with USD"""
133.
134. df['Currency'].replace("INR", "USD", inplace = True)
135. df.head()
136.
137. """## Adding a discount percentage column"""
138.
139. df['Discount Percentage'] = ((df['Price'] - df['Discount price']) / df['Price']) * 100
140. df['Discount Percentage'] = df['Discount Percentage'].round(2).astype(str) + "%"
141. df.head()
142.
143. """# Checking for null values"""
144.
145. df.isna().sum()
146.
147. """# Fixing the missing values
148.
149. ## Fixing the Prices
150. """
151.
152. for column in ['Price', 'Discount price']:
153.     df.fillna({column : df[column].median()}, inplace = True)
154. df.isna().sum()
155.
156. """## Fixing the missing currency"""
157.
158. df.fillna({"Currency": "Unknown"}, inplace = True) # Changing all the null values in price currency
column to 'unknown'
159. print(df['Currency'].isna().sum())
160.
161. df['Currency'].unique()
162.
163. """# Typo correction"""
164.
165. df['Title'] = df['Title'].str.lower() # Converting the title column to lower case
166. df['Title'] = df['Title'].str.strip().replace(r'\s+', ' ', regex=True) # Removing extra spaces
167. df['Title'] = df['Title'].str.replace(r'^\w\s', '', regex=True) # Removing special
characters

```



```

170. """# Categorizing the Courses"""
171. def categorize_title(title):
172.     title_lower = title.lower()
173.     if 'sql' in title_lower or 'mysql' in title_lower or 'database' in title_lower:
174.         return 'Database'
175.     elif 'tableau' in title_lower or 'power bi' in title_lower or 'data viz' in title_lower:
176.         return 'Data Visualization'
177.     elif 'excel' in title_lower or 'spreadsheet' in title_lower:
178.         return 'Spreadsheet'
179.     elif any(kw in title_lower for kw in ['agile', 'scrum', 'pmp', 'project management']):
180.         return 'Project Management'
181.     elif any(kw in title_lower for kw in ['financial', 'finance', 'accounting']):
182.         return 'Finance'
183.     elif 'mba' in title_lower or 'business' in title_lower or 'enterprise' in title_lower:
184.         return 'Business'
185.     elif any(kw in title_lower for kw in ['write', 'writing', 'editorial']):
186.         return 'Writing'
187.     elif 'sale' in title_lower or 'marketing' in title_lower:
188.         return 'Sales/Marketing'
189.     elif any(kw in title_lower for kw in ['data science', 'analytics', 'machine learning']):
190.         return 'Data Science'
191.     elif 'management' in title_lower:
192.         return 'Management'
193.     elif 'leadership' in title_lower:
194.         return 'Leadership'
195.     elif 'communication' in title_lower:
196.         return 'Communication'
197.     else:
198.         return 'Other'
199.
200. df['Category'] = df['Title'].apply(categorize_title)
201. df.head()
202.
203. df['Category'].value_counts()
204.
205. """# Outlier detection and deletion
206.
207. ## Visualizing Box Plot Graph for Columns with potential outliers
208. """
209.
210. numerical_columns = ['Subscribers', 'Price', 'Discount price']
211.
212. for column in numerical_columns:
213.     plt.boxplot(df[column], vert=False)
214.     plt.title(f'Box Plot of {column}')
215.     plt.ylabel('column')
216.     plt.grid()
217.     plt.show()
218.
219. """## Calculating the standard deviation and mean for each feature"""
220.
221. for column in numerical_columns:
222.     mu = df[column].mean()
223.     sigma = df[column].std()
224.     plt.hist(df[column])
225.     plt.title(f'Box Plot of {column}')
226.     plt.ylabel('column')
227.     plt.axvline(mu, color="black", linestyle="solid", label=f"Mean = {mu:.2f}")
228.     plt.axvline(mu + 3 * sigma, color='red', linestyle="dashed", label=f"+{3} Std Dev")
229.     plt.axvline(mu - 3 * sigma, color='red', linestyle="dashed", label=f"-{3} Std Dev")
230.     plt.grid()
231.     plt.legend()
232.     plt.show()
233.

```

```

235. """## Calculating each outlier info"""
236. outlier_info = {} # Creating an empty dictionary with key:value pair and each value has multiple
sub-values for the same column
237.
238. # Looping through the numerical columns to calculate mean, standard deviation, and outliers
239. for column in numerical_columns:
240.     mean_value = df[column].mean()
241.     std_dev = df[column].std()
242.
243.     # Finding outliers using 3 standard deviations from the mean
244.     outliers = df[(df[column] < mean_value - 3 * std_dev) | (df[column] > mean_value + 3 * std_dev)]
245.
246.     # Storing the details in the dictionary
247.     outlier_info[column] = {
248.         'mean': round(mean_value, 2),
249.         'std_dev': round(std_dev, 2),
250.         'num_outliers': len(outliers),
251.         'outlier_values': outliers[column].values.tolist()[:5] # Show first 5 outliers
252.     }
253.
254. # Iterating through the outlier_info dictionary and printing each column's details
255. for column in outlier_info:
256.     print(f"Outlier info for {column}:")
257.     print(f"    Mean: {outlier_info[column]['mean']}")
258.     print(f"    Standard Deviation: {outlier_info[column]['std_dev']}")
259.     print(f"    Number of Outliers: {outlier_info[column]['num_outliers']}")
260.     print(f"    Sample Outliers: {outlier_info[column]['outlier_values']}")
261.     print("-" * 50)
262.
263. """## Exporting the Cleaned data to a CSV file"""
264.
265. df.to_csv('cleaned_data.csv', index=False)
266.
267. print("Data exported successfully as 'cleaned_data.csv'")

```

```

269. """# MySQL Implementation
270.
271. ## Creating a connection creation function
272. """
273.
274. def create_connection():
275.     """Create and return a MySQL database connection and cursor."""
276.
277.     # Establish the connection
278.     mydb = mysql.connector.connect(
279.         host="localhost",
280.         user="root",
281.         password="",
282.         database=""          # Need to fill this path for the intended database
283.     )
284.
285.     # Create a cursor
286.     conn = mydb.cursor()
287.
288.     print("Connection to the database was successful.")
289.     print("≡" * 50)
290.     return mydb, conn
291.
292. """## Connecting the Database"""
293.
294. mydb, conn = create_connection()
295.
296. """## Database code"""
297.
298. file_path = "cleaned_data.csv"
299.
300. data = pd.read_csv(file_path)
301. data.columns
302.
303. # Re-connecting to the database
304. mydb, conn = create_connection()
305.
306. # Dropping the existing database if it exists and creating a new one
307. conn.execute("DROP DATABASE IF EXISTS Course_db")
308. conn.execute("CREATE DATABASE IF NOT EXISTS Course_db")
309. conn.execute("USE Course_db")
310.
311. # Creating the table with the specified column names
312. create_sql = """
313. CREATE TABLE IF NOT EXISTS courses_table (
314.     id INT PRIMARY KEY,
315.     Title TEXT,
316.     is_paid BOOLEAN,
317.     Subscribers INT,
318.     `Average rating` FLOAT,
319.     Rating FLOAT,
320.     `Num reviews` INT,
321.     Lectures INT,
322.     created DATETIME,
323.     Time DATETIME,
324.     `Discount price` FLOAT,
325.     Price FLOAT,
326.     Currency TEXT,
327.     `Discount Percentage` FLOAT,
328.     Category TEXT
329. );
330. """
331. conn.execute(create_sql)
332.
333. conn.close()
334.

```

```

336. """## Making the max allowed packet size larger so the process of data population becomes easier"""
337. # Setting the max_allowed_packet size to a larger value
338. mydb, conn = create_connection()
339. conn.execute("SET GLOBAL max_allowed_packet=64*1024*1024") # Set to 64MB
340.
341. """## Populating the Database with the data in our exported CSV file"""
342.
343. # Re-connecting to the database
344. mydb, conn = create_connection()
345. conn.execute("USE Course_db")
346.
347. # Inserting data into the MySQL table
348. insert_sql = """
349. INSERT INTO courses_table (id, Title, is_paid, Subscribers, `Average rating`, Rating, `Num reviews`,
350.                             Lectures, created, Time, `Discount price`, Price, Currency,
351.                             `Discount Percentage`, Category)
352. VALUES (%s, %s, %s, %s, %s, %s, %s, %s, %s, %s, %s, %s, %s, %s, %s)
353. """
354. data_values = [
355.     (
356.         row['id'], row['Title'], row['is_paid'], row['Subscribers'], row['Average rating'],
357.         row['Rating'], row['Num reviews'], row['Lectures'],
358.         row['created'], row['Time'], row['Discount price'], row['Price'],
359.         row['Currency'], row['Discount Percentage'], row['Category']
360.     )
361.     for _, row in df.iterrows()
362. ]
363.
364. # Executing batch insert
365. conn.executemany(insert_sql, data_values)
366. mydb.commit()
367.
368. print("Data successfully saved to the MySQL database.")
369.
370. # Closing the connection
371. conn.close()
372. mydb.close()
373.
374. """## Checking the Columns in order to Ensure consistency"""
375.
376. # Re-connecting to the database
377. mydb, conn = create_connection()
378. conn.execute("USE Course_db")
379.
380. conn.execute("DESCRIBE courses_table;")
381. columns = conn.fetchall()
382. print("Table columns:", columns)
383.
384.

```

```

"""# MySQL Queries
385.
386. ## Sort and identify top courses by:
387.
388. ### Rating
389. """
390.
391. # Re-connecting to the database
392. mydb, conn = create_connection()
393. conn.execute("USE Course_db")
394.
395. # Sorting by Rating (Top 10 Courses)
396. Sort_by_Rating = """
397. SELECT Title, `Average rating`
398. FROM courses_table
399. ORDER BY `Average rating` DESC
400. LIMIT 10;
401. """
402.
403. # Execute the query
404. conn.execute(Sort_by_Rating)
405.
406. # Fetch the results
407. results = conn.fetchall()
408.
409. # Check if any results are found and print them
410. if results:
411.     print("Top 10 Courses by Rating:")
412.     print("-"*30)
413.     for row in results:
414.         print(row)
415.         print("-"*70)
416. else:
417.     print("No results found.")
418.
419. """### Number of subscribers (if available)
420.
421. """
422.
423. # Re-connecting to the database
424. mydb, conn = create_connection()
425. conn.execute("USE Course_db")
426.
427. # Sorting by Subscribers (Top 10 Courses)
428. Sort_by_Subscribers = """
429. SELECT Title, Subscribers
430. FROM courses_table
431. ORDER BY Subscribers DESC
432. LIMIT 10;
433. """
434.
435. # Execute the query
436. conn.execute(Sort_by_Subscribers)
437.
438. # Fetch the results
439. results = conn.fetchall()
440.
441. # Check if any results are found and print them
442. if results:
443.     print("Top 10 Courses by Subscribers:")
444.     print("-"*30)
445.     for row in results:
446.         print(row)
447.         print("-"*70)
448. else:
449.     print("No results found.")

```

```

"""### Discount percentage (created during preprocessing)"""
452.
453. # Re-connecting to the database
454. mydb, conn = create_connection()
455. conn.execute("USE Course_db")
456.
457. # Sorting by Discount Percentage (Top 10 Courses)
458. Sort_by_Discount_Percentage = """
459. SELECT Title, `Discount Percentage`
460. FROM courses_table
461. ORDER BY `Discount Percentage` DESC
462. LIMIT 10;
463. """
464.
465. # Execute the query
466. conn.execute(Sort_by_Discount_Percentage)
467.
468. # Fetch the results
469. results = conn.fetchall()
470.
471. # Check if any results are found and print them
472. if results:
473.     print("Top 10 Courses by Discount Percentage:")
474.     print("-"*40)
475.     for row in results:
476.         print(row)
477.         print("-"*70)
478. else:
479.     print("No results found.")
480.
481. """### Identify the most expensive courses
482.
483. """
484.
485. # Re-connecting to the database
486. mydb, conn = create_connection()
487. conn.execute("USE Course_db")
488.
489. # Sorting by Discount Percentage (Top 10 Courses)
490. Sort_by_most_expensive = """
491. SELECT Title, price
492. FROM courses_table
493. ORDER BY price DESC
494. LIMIT 10;
495. """
496.
497. # Execute the query
498. conn.execute(Sort_by_most_expensive)
499.
500. # Fetch the results
501. results = conn.fetchall()
502.
503. # Check if any results are found and print them
504. if results:
505.     print("Top 10 Most Expensive Courses in USD: ")
506.     print("-"*40)
507.     for row in results:
508.         print(row)
509.         print("-"*70)
510. else:
511.     print("No results found.")
512.

```

```

514. """## Group the data by category and calculate:
515. ### Number of courses in each category
516. """
517.
518. # Re-connecting to the database
519. mydb, conn = create_connection()
520. conn.execute("USE Course_db")
521.
522. # Calculating the number of courses per price category
523. Sort_by_price_category = """
524. SELECT Category, COUNT(*) AS num_courses
525. FROM courses_table
526. GROUP BY Category;
527. """
528.
529. # Execute the query
530. conn.execute(Sort_by_price_category)
531.
532. # Fetch the results
533. results = conn.fetchall()
534.
535. # Check if any results are found and print them
536. if results:
537.     print("Count of courses for each Price Category: ")
538.     print("-"*40)
539.     for row in results:
540.         print(row)
541.         print("-"*70)
542. else:
543.     print("No results found.")
544.
545. """### Average number of reviews per category"""
546.
547. # Re-connecting to the database
548. mydb, conn = create_connection()
549. conn.execute("USE Course_db")
550.
551. # Calculating the average number of reviews per Price category
552. Sort_by_reviews_category = """
553. SELECT Category, AVG(`Num reviews`) AS avg_reviews
554. FROM courses_table
555. GROUP BY Category;
556. """
557.
558. # Execute the query
559. conn.execute(Sort_by_reviews_category)
560.
561. # Fetch the results
562. results = conn.fetchall()
563.
564. # Check if any results are found and print them
565. if results:
566.     print("Average Number of Reviews per Price Category: ")
567.     print("-"*50)
568.     for row in results:
569.         print(row)
570.         print("-"*70)
571. else:
572.     print("No results found.")
573.

```

```

575. """## Analyze discounts:
576. ### Calculate average discount percentage offered per category
577. """
578.
579. # Re-connecting to the database
580. mydb, conn = create_connection()
581. conn.execute("USE Course_db")
582.
583. # Categorizing by Price Range
584. Sort_by_avg_discount_percentage = """
585. SELECT Category, AVG(`Discount Percentage`) AS avg_discount_percentage
586. FROM courses_table
587. GROUP BY Category;
588. """
589.
590. # Execute the query
591. conn.execute(Sort_by_avg_discount_percentage)
592.
593. # Fetch the results
594. results = conn.fetchall()
595.
596. # Check if any results are found and print them
597. if results:
598.     print("Average Discount Percentage Offered per Price Category: ")
599.     print("-"*50)
600.     for row in results:
601.         print(row)
602.         print("-"*70)
603. else:
604.     print("No results found.")
605.
606. """## Analyze price by paid/free status (if available):
607.
608. ### Calculate the median price for courses in each category.
609. """
610.
611. # Re-connecting to the database
612. mydb, conn = create_connection()
613. conn.execute("USE Course_db")
614.
615. # Categorizing paid and free courses based on 'is_paid' column
616. checking_for_paid_or_free = """
617. SELECT
618.     CASE
619.         WHEN is_paid = 0 THEN 'Free'
620.         WHEN is_paid = 1 THEN 'Paid'
621.     END AS Status,
622.     COUNT(*) AS num_courses
623. FROM courses_table
624. GROUP BY Status;
625. """
626.
627. # Execute the query
628. conn.execute(checking_for_paid_or_free)
629.
630. # Fetch the results
631. results = conn.fetchall()
632.
633. # Check if any results are found and print them
634. if results:
635.     print("Paid / Free Category count: ")
636.     print("-"*50)
637.     for row in results:
638.         print(row)
639.         print("-"*70)
640. else:
641.     print("No results found.")
642.
643. # Closing the connection
644. conn.close()
645. mydb.close()

```



```

647. """As the Paid courses are the only ones I can actually calculate their median price value, since
Free courses are free so their median would be zero. I will only calculate for the paid courses"""
648.
649. # Re-connecting to the database
650. mydb, conn = create_connection()
651. conn.execute("USE Course_db")
652.
653. # Categorizing paid and free courses based on 'is_paid' column
654. calculate_median = """
655. SELECT
656.     price
657. FROM courses_table
658. WHERE is_paid = 1 AND price IS NOT NULL
659. ORDER BY price;
660. """
661.
662. # Execute the query
663. conn.execute(calculate_median)
664.
665. # Fetch all the prices
666. prices = conn.fetchall()
667.
668. # Calculate median
669. num_prices = len(prices)
670. if num_prices > 0:
671.     sorted_prices = [price[0] for price in prices]
672.     if num_prices % 2 == 1: # Odd number of prices
673.         median_price = sorted_prices[num_prices // 2]
674.     else: # Even number of prices
675.         median_price = (sorted_prices[num_prices // 2 - 1] + sorted_prices[num_prices // 2]) / 2
676.     print(f"Median Price for Paid Courses: {median_price}")
677. else:
678.     print("No valid price data found for Paid courses.")
679.
680. # Closing the connection
681. conn.close()
682. mydb.close()
683.

```

## References

### 1) Kaggle Dataset

"Finance & Accounting Courses - Udemy (13K+ course) - Kaggle."

<https://www.kaggle.com/datasets/jilkothari/finance-accounting-courses-udemy-13k-course>.

### 2) Libraries and Tools used.

- 1) Visual Studio Code Application
- 2) Python Programming Language
- 3) Jupyter Notebook
- 4) Numpy Library
- 5) Pandas Library
- 6) Matplotlib Library
- 7) Seaborn Library
- 8) mysql.connector Library
- 9) Tableau Dashboard