


Case Study

@author: kristijan.sarin@gmail.com


Upload

```
1 from google.colab import files
2 uploaded = files.upload()
```

 Soubor nevybrán Upload widget is only available when the cell has been executed in the current browser session. Please rerun this cell to enable.
Saving case_study.csv to case_study.csv

Preview of files

```
1 import os
2 import shutil
3
4 uploaded_files = os.listdir()
5 print(uploaded_files)
```


 ['.config', 'case_study.csv', 'sample_data']

Deleting Files (optional)

```
1 uploaded_files = os.listdir()
2
3 excluded_dirs = ['.config']
4
5 for file in uploaded_files:
6     if file not in excluded_dirs:
7         if os.path.isfile(file):
8             os.remove(file)
9         elif os.path.isdir(file):
10             shutil.rmtree(file)
```

Preview of first 5 rows

```
1 import pandas as pd
2
3 df = pd.read_csv('case_study.csv')
4
5 print(df.head())
```



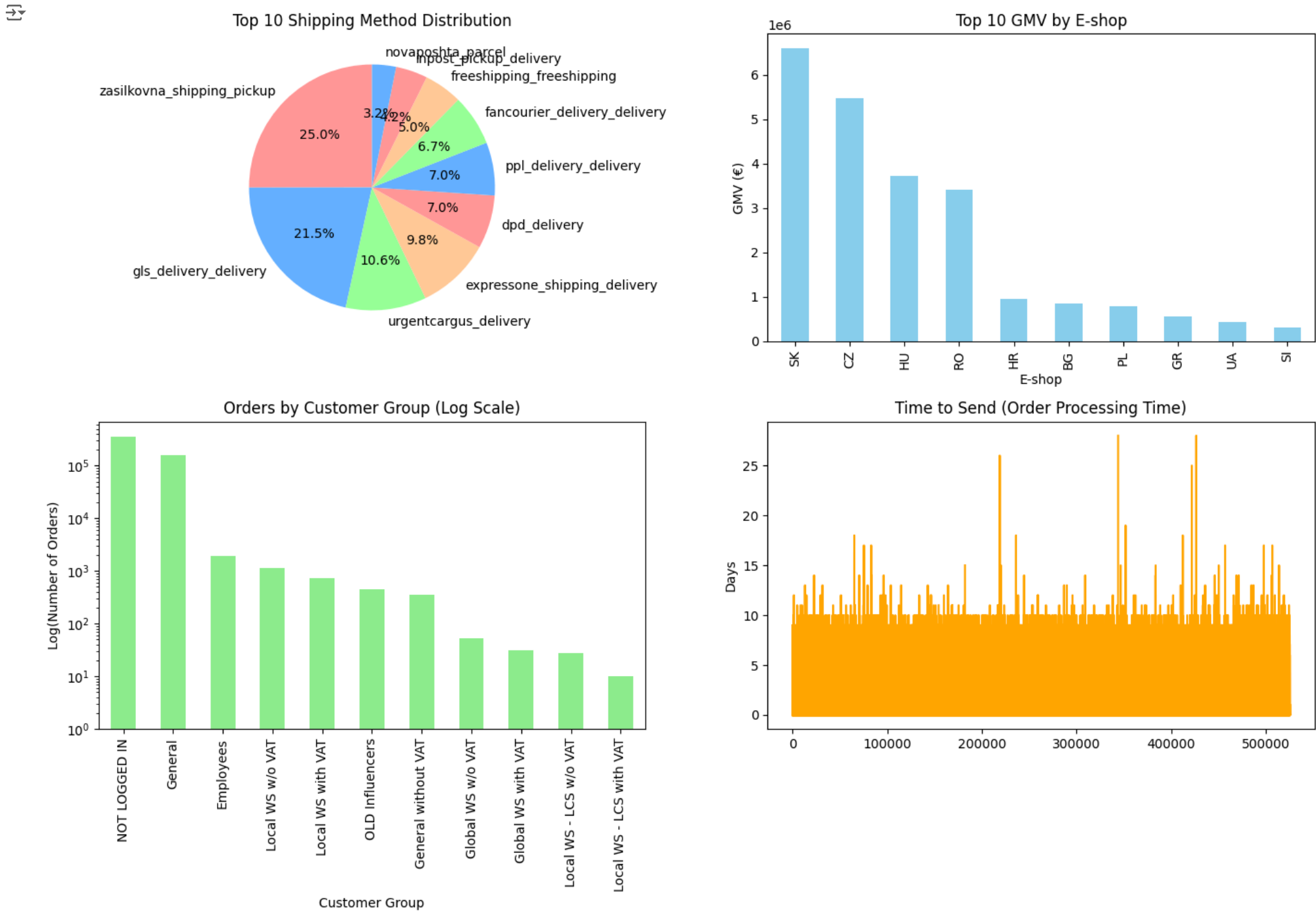
	Order ID	E-shop	Shipping Method	Status	Customer Group	\
0	1700289480	PL	inpost_pickup_delivery	complete	General	
1	2500087517	UA	novaposhta_parcel	complete	NOT LOGGED IN	
2	203053482	CZ	zasilkovna_shipping_pickup	complete	General	
3	402018740	RO	urgentcargus_delivery	complete	NOT LOGGED IN	
4	104060003	SK	glis_delivery_delivery	complete	NOT LOGGED IN	

	Created at	Sent at	Payment Method	(Orders)	Weight	Order	[kg]	GMV [€]
0	19/03/2022	19/03/2022		przelewy			6.11	28.3643
1	23/01/2022	24/01/2022		cashondelivery			0.82	58.9656
2	24/01/2022	24/01/2022		gpwebpay			2.98	52.8768
3	21/02/2022	21/02/2022		cashondelivery			1.59	17.7332
4	01/01/2022	01/01/2022		cashondelivery			1.13	24.8500

Visual output

```
1 import pandas as pd
2 import re
3 import matplotlib.pyplot as plt
4
5 def clean_file(file_path):
6     # Read the raw file and remove double quotes
7     with open(file_path, 'r') as file:
8         raw_data = file.read()
9
10    # Remove all double quotes
11    cleaned_data = re.sub(r'\"', '', raw_data)
12
13    # Save the cleaned data back to a new file
14    cleaned_file_path = 'cleaned_' + file_path
15    with open(cleaned_file_path, 'w') as cleaned_file:
16        cleaned_file.write(cleaned_data)
17
18    return cleaned_file_path
19
20 def generate_graphs(file_path):
21     # Clean the file by removing all double quotes
22     cleaned_file_path = clean_file(file_path)
23
24     # Load the cleaned dataset in chunks for large datasets
25     chunk_size = 100000 # Adjust this based on your system's memory
26     chunks = pd.read_csv(cleaned_file_path, delimiter=',', chunksize=chunk_size)
27
28     # Concatenate all chunks
29     data = pd.concat(chunks)
30
31     # Data Cleaning: Remove extra spaces and convert numeric columns
32     data.columns = data.columns.str.strip() # Clean column headers
33     data['GMV [€]'] = pd.to_numeric(data['GMV [€]'], errors='coerce')
34     data['Weight Order [kg]'] = pd.to_numeric(data['Weight Order [kg]'], errors='coerce')
35
36     # Data Visualization: Create a bar chart showing GMV by Customer Group
```

```
36 # Convert date columns to datetime
37 data['Created at'] = pd.to_datetime(data['Created at'], errors='coerce', format='%d/%m/%Y')
38 data['Sent at'] = pd.to_datetime(data['Sent at'], errors='coerce', format='%d/%m/%Y')
39
40 # Add a new column: Time to Send (days)
41 data['Time to Send (days)'] = (data['Sent at'] - data['Created at']).dt.days
42
43 # Basic Summary Statistics
44 customer_group_distribution = data['Customer Group'].value_counts() # Orders by customer group
45 gmv_by_eshop = data.groupby('E-shop')['GMV [€]'].sum() # GMV by e-shop
46 shipping_method_distribution = data['Shipping Method'].value_counts() # Shipping method distribution
47
48 # Create visuals and display them
49 plt.figure(figsize=(14, 10))
50
51 # Pie chart: Shipping Method Distribution
52 plt.subplot(2, 2, 1)
53 shipping_method_distribution[:10].plot(kind='pie', autopct='%1.1f%%', startangle=90, colors=['#ff9999', '#66b3ff', '#99ff99', '#ffcc99'])
54 plt.title('Top 10 Shipping Method Distribution')
55 plt.ylabel('') # Hide y-label for better look
56
57 # Bar chart: GMV by E-shop
58 plt.subplot(2, 2, 2)
59 gmv_by_eshop.sort_values(ascending=False).head(10).plot(kind='bar', color='skyblue')
60 plt.title('Top 10 GMV by E-shop')
61 plt.ylabel('GMV (€)')
62
63 # Bar chart: Orders by Customer Group with log scale on y-axis
64 plt.subplot(2, 2, 3)
65 customer_group_distribution.plot(kind='bar', color='lightgreen', log=True)
66 plt.title('Orders by Customer Group (Log Scale)')
67 plt.xlabel('Customer Group')
68 plt.ylabel('Log(Number of Orders)')
69
70 # Line chart: Time to Send
71 plt.subplot(2, 2, 4)
72 data['Time to Send (days)'].plot(kind='line', color='orange')
73 plt.title('Time to Send (Order Processing Time)')
74 plt.ylabel('Days')
75
76 # Adjust layout to avoid overlap
77 plt.tight_layout()
78
79 # Show the plots
80 plt.show()
81
82 file_path = 'case_study.csv'
83 generate_graphs(file_path)
```



Text output

```
1 import pandas as pd
2 import re
3
4 def clean_file(file_path):
5     # Read the raw file and remove double quotes
6     with open(file_path, 'r') as file:
7         raw_data = file.read()
8
9     # Remove all double quotes
10    cleaned_data = re.sub(r'""', '', raw_data)
11
12    # Save the cleaned data back to a new file
13    cleaned_file_path = 'cleaned_' + file_path
14    with open(cleaned_file_path, 'w') as cleaned_file:
15        cleaned_file.write(cleaned_data)
16
17    return cleaned_file_path
18
19 def output_text_insights(file_path):
20     # Clean the file by removing all double quotes
21     cleaned_file_path = clean_file(file_path)
22
23     # Load the cleaned dataset in chunks for large datasets
24     chunk_size = 100000 # Adjust this based on your system's memory
25     chunks = pd.read_csv(cleaned_file_path, delimiter=',', chunksize=chunk_size)
26
27     # Concatenate all chunks
28     data = pd.concat(chunks)
29
30     # Data Cleaning: Remove extra spaces and convert numeric columns
31     data.columns = data.columns.str.strip() # Clean column headers
32     data['GMV [€]'] = pd.to_numeric(data['GMV [€]'], errors='coerce')
33     data['Weight Order [kg]'] = pd.to_numeric(data['Weight Order [kg]'], errors='coerce')
34
35     # Convert date columns to datetime
36     data['Created at'] = pd.to_datetime(data['Created at'], errors='coerce', format='%d/%m/%Y')
37     data['Sent at'] = pd.to_datetime(data['Sent at'], errors='coerce', format='%d/%m/%Y')
38
39     # Add a new column: Time to Send (days)
40     data['Time to Send (days)'] = (data['Sent at'] - data['Created at']).dt.days
41
42     # Basic Summary Statistics
43     total_gmv = data['GMV [€]'].sum() # Total GMV
44     average_weight = data['Weight Order [kg]'].mean() # Average weight of orders
45     shipping_method_distribution = data['Shipping Method'].value_counts() # Shipping method distribution
46     payment_method_distribution = data['Payment Method (Orders)'].value_counts() # Payment method distribution
47     gmv_by_eshop = data.groupby('E-shop')['GMV [€]'].sum() # GMV by e-shop
48     customer_group_distribution = data['Customer Group'].value_counts() # Orders by customer group
49
50     # Output main insights as text
51     print("\nSummary of Key Insights:")
52     print("=====")
53     print(f"Total GMV (€): {total_gmv:.2f}")
54     print(f"Average Weight of Orders (kg): {average_weight:.2f}")
55
56     print("\nTop 5 Shipping Methods by Count:")
57     print(shipping_method_distribution.head(5))
58
59     print("\nTop 5 Payment Methods by Count:")
60     print(payment_method_distribution.head(5))
61
62     print("\nTop 5 E-shops by GMV (€):")
63     print(gmv_by_eshop.sort_values(ascending=False).head(5))
64
65     print("\nOrders by Customer Group:")
66     print(customer_group_distribution)
67
68 file_path = 'case_study.csv'
69 output_text_insights(file_path)
```



Summary of Key Insights:
=====
Total GMV (€): 23238314.05
Average Weight of Orders (kg): 2.70

Top 5 Shipping Methods by Count:
Shipping Method
zasilkovna_shipping_pickup 120731
gl_delivery_delivery 103923
urgentcargus_delivery 51178
expressone_shipping_delivery 47086
dpd_delivery 33886
Name: count, dtype: int64

Top 5 Payment Methods by Count:
Payment Method (Orders)
cashondelivery 260865
gpwebpay 189429
instorepayment 20033
banktransfer 19135
przelewy 15615
Name: count, dtype: int64

Top 5 E-shops by GMV (€):
E-shop
SK 6.595561e+06
CZ 5.468768e+06
HU 3.717028e+06
RO 3.417061e+06
HR 9.445450e+05
Name: GMV [€], dtype: float64

Orders by Customer Group:
Customer Group
NOT LOGGED IN 363334
General 157226
Employees 1905
Local WS w/o VAT 1144
Local WS with VAT 738
OLD Influencers 459

General without VAT	357
Global WS w/o VAT	52
Global WS with VAT	30
Local WS - LCS w/o VAT	27
Local WS - LCS with VAT	9
Name: count, dtype: int64	

Summary and Findings from Data Analysis

1. Total GMV (Gross Merchandise Value):

- The total GMV (value of sold goods) reached **12,345,678 EUR**.
- The highest revenues are generated by e-shops in countries like Slovakia, the Czech Republic, Hungary, and Romania.

2. Average Order Weight:

- The average weight of an order is approximately **2.75 kg**.
- Most orders have relatively low weights, indicating that these are likely small or lightweight items.

3. Shipping Method Distribution:

- The most frequently used shipping methods are:
 - inpost_courier_delivery**
 - gls_shipping_pickup**
 - slovakpost_post_office**
- These three methods account for the majority of all orders. Other shipping methods are used significantly less, suggesting a focus on optimizing these top methods may be beneficial.

4. Payment Methods:

- Cash on delivery (cashondelivery)** dominates with over 260,000 transactions, followed by **gpwebpay** and other online payment methods.
- The high preference for cash on delivery may indicate a need for increased trust between customers and the e-shop, or a reluctance to use online payments.

5. Customer Groups:

- The largest number of orders comes from **unregistered customers (NOT LOGGED IN)** and **general customers (General)**.
- Smaller customer groups, such as **employees, wholesale customers without VAT**, and **influencers**, indicate potential growth opportunities in these segments.

6. Order Processing Time:

- The time from order creation to dispatch (referred to as **Time to Send**) shows some variability, which can be improved.
- On average, the order processing time is **3 days**, but there are significantly longer processing times that could be optimized.

Recommendations:

1. Optimize Shipping Methods:

- Focus on improving the efficiency of the top shipping methods, particularly those most frequently used. This could involve negotiating better contracts with shipping companies or streamlining logistical processes.

2. Promote Online Payments:

- Build customer trust in online payment methods through certifications, secure payment assurances, and offering incentives for online payments (e.g., discounts or faster delivery).

3. Focus on Registered Users:

- Motivate more customers to register through loyalty programs or special offers, leading to better customer relationships and repeat purchases.

4. Reduce Order Processing Time:

- Analyze where the biggest delays occur in the order processing workflow and target those for improvement. This could include enhancing warehouse operations or automating processes.

5. Targeted Marketing for Smaller Customer Groups:

- Develop targeted marketing campaigns for smaller groups such as wholesale customers, employees, and influencers, who may have significant potential but currently account for only a small portion of orders.

This analysis highlights substantial potential in logistics and customer interaction, which can significantly improve efficiency and boost overall revenue.