

---

Prepared by Sama Osama

# *WEEK 2 OF CELLULA INTERNSHIP*

## *Task'2*

### Booking Cancellation Prediction

Hotel Reservation Dataset Analysis &  
Modeling

---

# Dataset Overview

- 36,285 rows × 17 columns
- Covers reservation behavior, pricing, and booking outcomes
- Target variable: booking\_status (Canceled vs. Not\_Canceled)

[illegible]

# *Data Preprocessing Steps*

1- Standardized column names (lowercase, removed spaces)

3- Check for Null values and duplicates

5- Dropped invalid date entries (NAT)



2- Fixed mislabeled columns (e.g., p-c → preservation\_canceled)

4- Converted date of reservation to datetime

# Feature Engineering

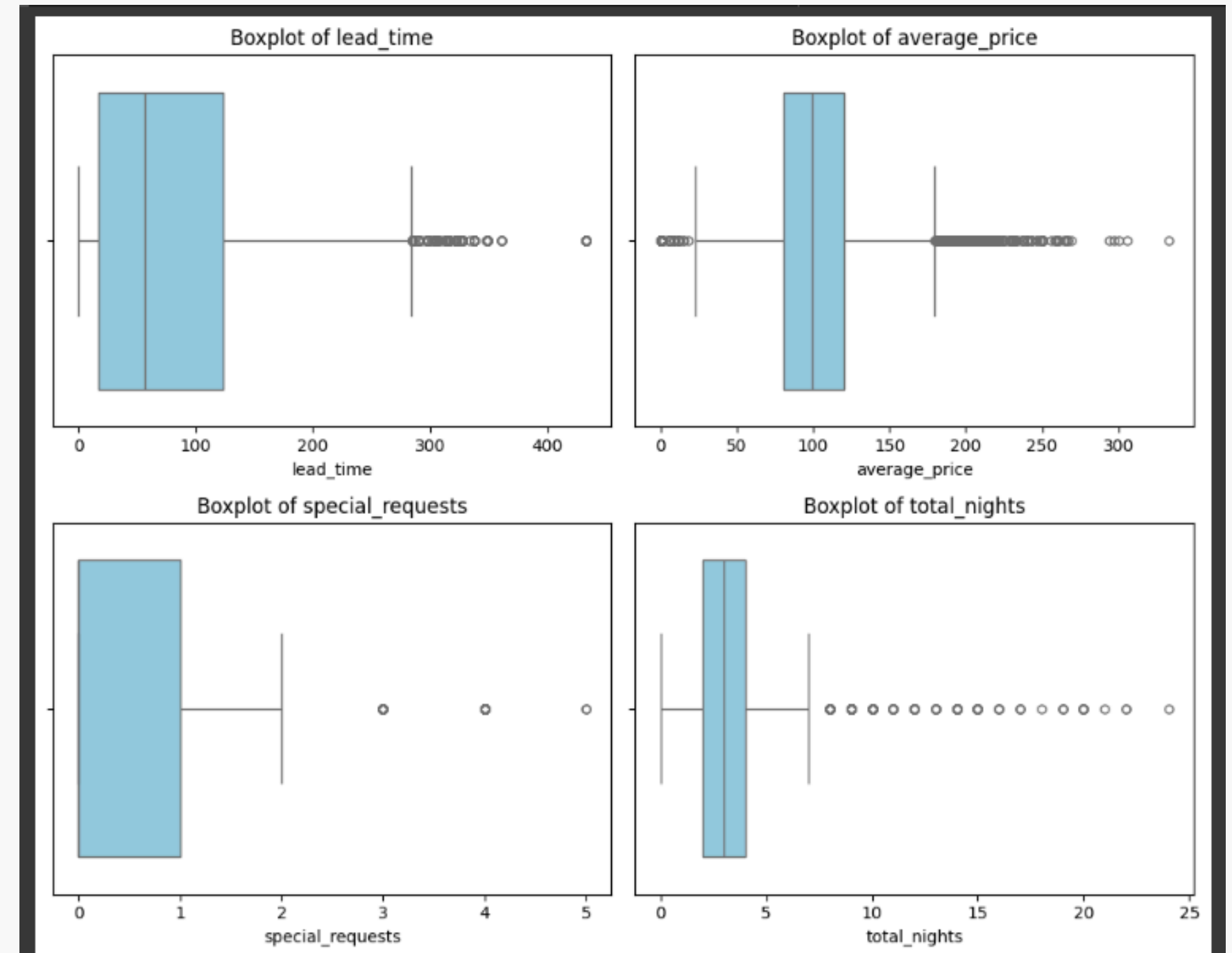
- Extracted: reservation\_year and reservation\_month from date\_of\_reservation
- total\_nights = number\_of\_weekend\_nights + number\_of\_week\_nights
- Dropped irrelevant fields like booking\_id

	booking_id	number_of_adults	number_of_children	number_of_weekend_nights	number_of_week_nights	type_of_meal	car_parking_space	room_type	lead_time	market_segment_type	Visited_Before	Preservation_Canceled	Preservation_not_Canceled	average_price	special_needs
0	INN00001	1	1	2	5	Meal Plan 1	0	Room_Type 1	224	Offline	0	0	0	88.00	
1	INN00002	1	0	1	3	Not Selected	0	Room_Type 1	5	Online	0	0	0	106.68	
4	INN00005	1	0	1	2	Not Selected	0	Room_Type 1	48	Online	0	0	0	77.00	
8	INN00009	1	1	0	4	Meal Plan 1	0	Room_Type 1	121	Offline	0	0	0	96.90	
10	INN00011	1	0	1	0	Not Selected	0	Room_Type 1	0	Online	0	0	0	85.03	
...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...
36272	INN36273	2	0	2	6	Meal Plan 1	0	Room_Type 1	148	Online	0	0	0	98.39	
36275	INN36276	2	0	1	2	Meal Plan 1	0	Room_Type 1	224	Offline	0	0	0	65.00	
36276	INN36277	2	0	2	3	Not Selected	0	Room_Type 1	5	Online	0	0	0	106.68	
36279	INN36281	2	0	1	1	Not Selected	0	Room_Type 1	48	Online	0	0	0	94.50	
36283	INN36285	3	0	0	4	Meal Plan 1	0	Room_Type 1	121	Offline	0	0	0	96.90	

14149 rows x 19 columns

# Outlier Detection and Handling using Z-score

- Applied Z-Score method on:
- lead\_time, average\_price, special\_requests, total\_nights
- Removed records with  $Z > 3$

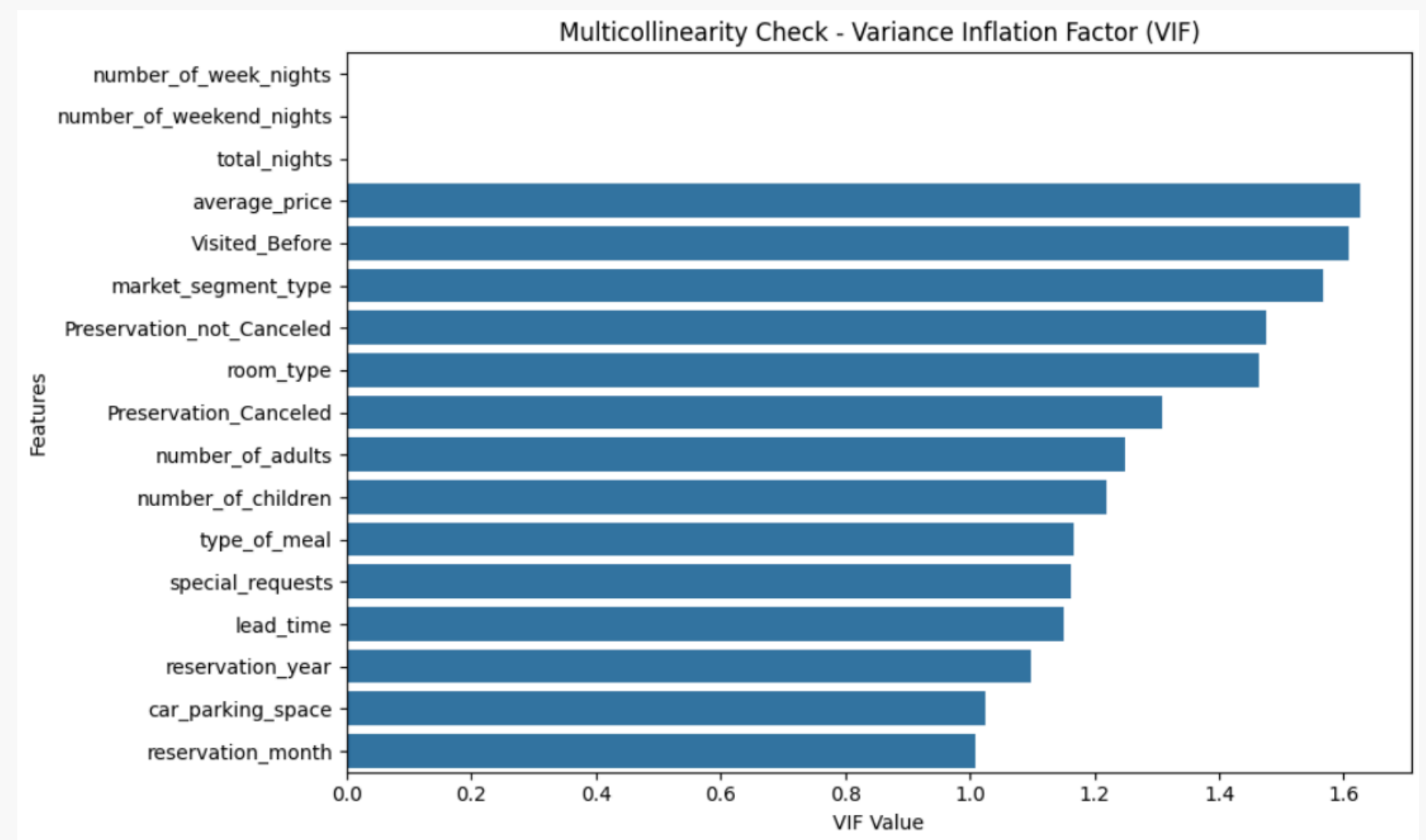


# *Transformation For The Categorical Data*

- Label encoded:
- type\_of\_meal, room\_type, market\_segment\_type
- Converted booking\_status to 0 (Not\_Canceled) and 1 (Canceled)
- Applied StandardScaler to numerical features

# Multicollinearity Check

- Calculate VIF (Variance Inflation Factor)
- Visualize VIF values as a horizontal bar chart
- Use it to detect multicollinearity among features
- number\_of\_week\_nights, number\_of\_weekend\_nights, and total\_nights have  $VIF = \infty$ , indicating perfect multicollinearity.
- Most other features have  $VIF < 2$ , suggesting low multicollinearity.



# *Train-Test Split*

- 80% training / 20% test
- Scaled and resampled training data
- Final sets:
- X\_train\_balanced, X\_test\_scaled, y\_train\_balanced, y\_test

```
Logistic Regression Accuracy: 0.7922
```

```
Classification Report:
```

	precision	recall	f1-score	support
0	0.81	0.89	0.85	1865
1	0.74	0.61	0.67	965
accuracy			0.79	2830
macro avg	0.78	0.75	0.76	2830
weighted avg	0.79	0.79	0.79	2830



# Modeling Approaches

## Results – Logistic Regression

- Accuracy: e.g., 79.22%
- Precision, Recall, F1-score (from classification report)
- Confirms model’s effectiveness in predicting cancellations
- Balanced performance across both classes after SMOTE

Metric	Class 0 (Not Canceled)	Class 1 (Canceled)
Precision	0.81	0.74
Recall	0.89	0.61
F1-Score	0.85	0.67
<ul style="list-style-type: none"><li>• Overall Accuracy: 79.22%</li><li>• Macro F1 Score: 76%</li><li>• Weighted F1 Score: 79%</li></ul>		

---

*Thank you*

---