



# Exploratory Data Analysis

16<sup>th</sup> Juillet 2025

**Project: Bank Marketing Campaign**

**Name: Najma Abdi, Leena Ganta, Adama Sall**

**Batch Code: LISUM45**

**Specialization: Data Science**

# Agenda

Problem Statement

Dataset

Model

Evaluation

Recommendations

# Problem Statement

- ABC Bank is planning to launch a new term deposit product and wants to predict which customers are likely to subscribe based on their past interactions with previous marketing campaigns. By developing a machine learning model, the bank aims to identify high-probability prospects in advance, allowing the marketing team to focus its efforts where they are most likely to succeed. This targeted approach will help reduce overall campaign costs while improving conversion rates.
- The primary business objective is to predict whether a client will subscribe to a term deposit. This binary classification problem helps optimize marketing efforts and reduce unnecessary contact with clients

# Dataset

## [Data Set UCI Link](#)

After downloading the data and uploaded it in Jupyter Notebook along with `import pandas as pd` module, we have found the following:

Dataset: bank-additional-full.csv, Date Analyzed: 2025-06-19

Tool Used: Python (pandas)

Data Overview:

- Shape: 41188 rows × 21 columns
- Target Variable: 'y' (binary: yes/no)
- Class Distribution (Target)
  - No: 36548 (88.73%)
  - Yes: 4640 (11.27%)

- housing: 990
- loan: 990

Duplicate rows: 12

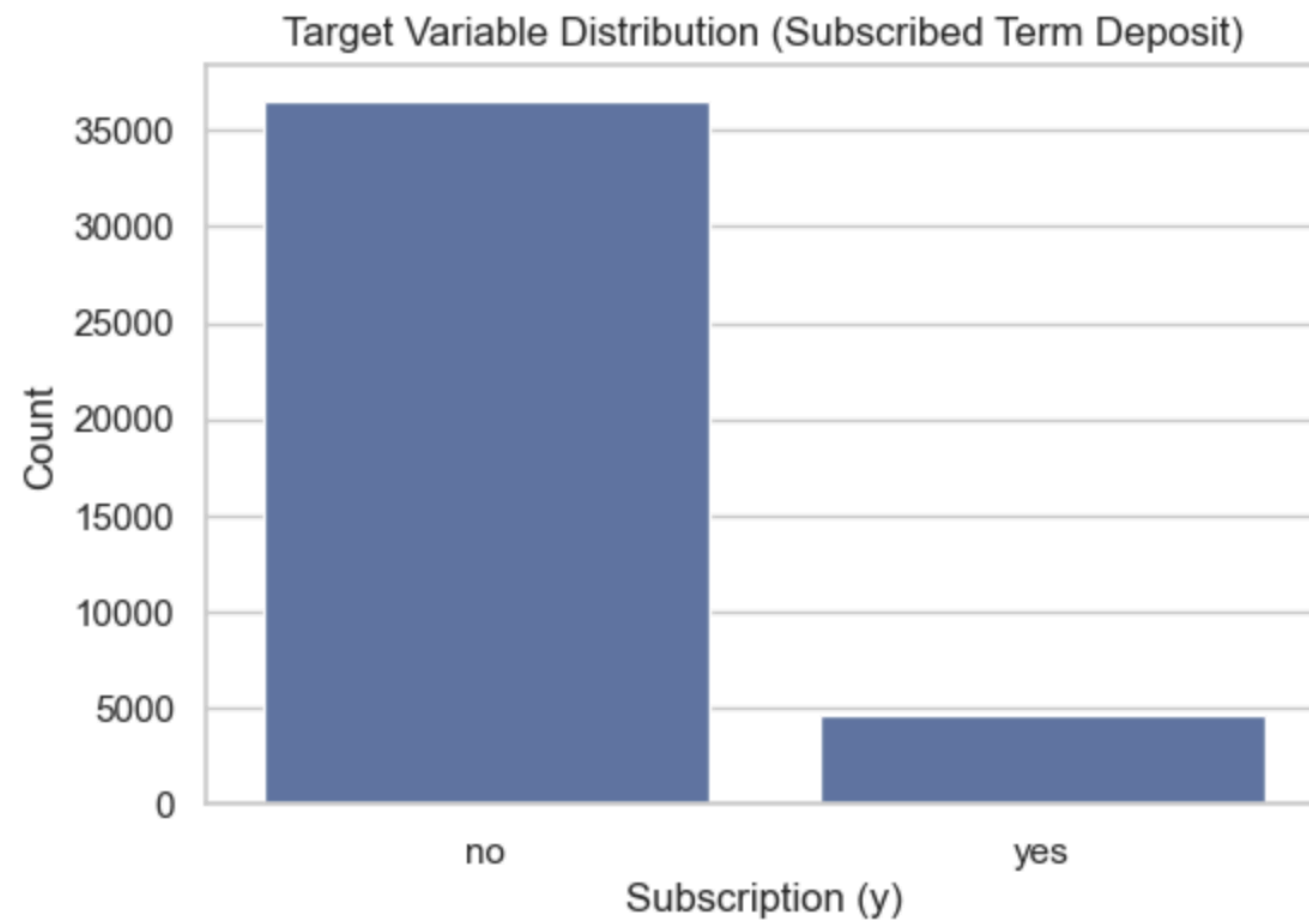
# EDA

# EDA

# 1. Univariate Analysis



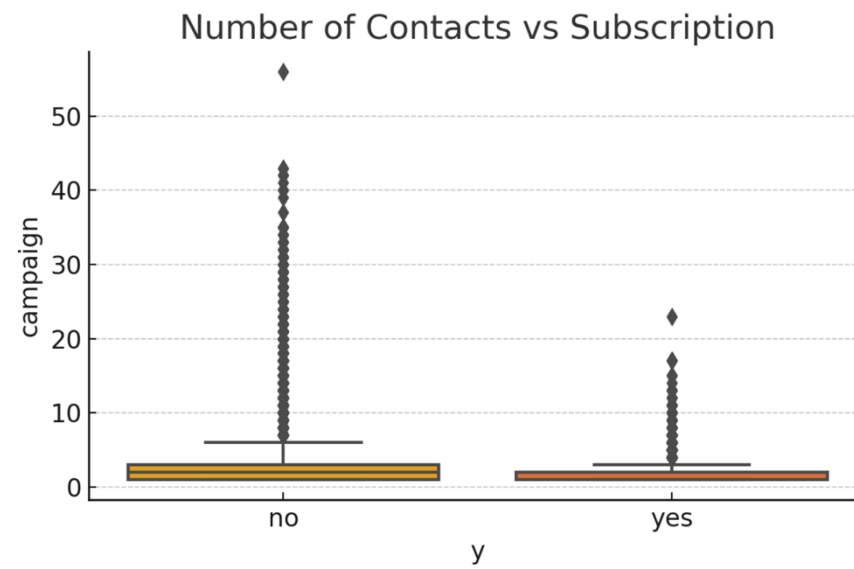
- Majority ( $\sim 89\%$ ) of clients do not subscribe.
- Success rate is  $\sim 11\%$ .
- This indicates class imbalance and campaign challenges.



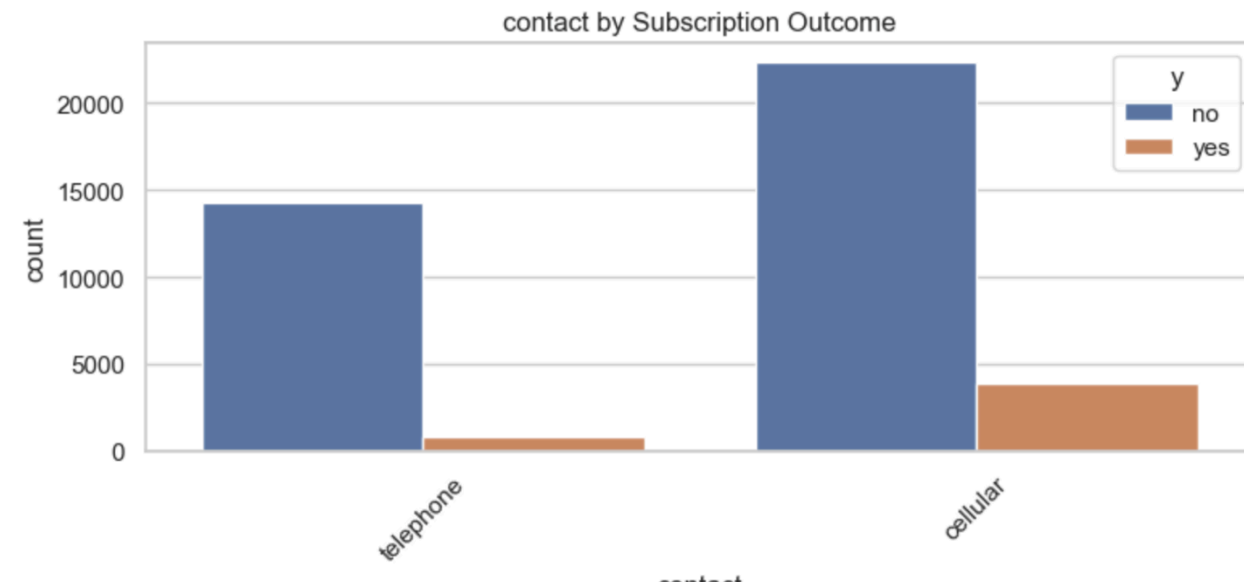


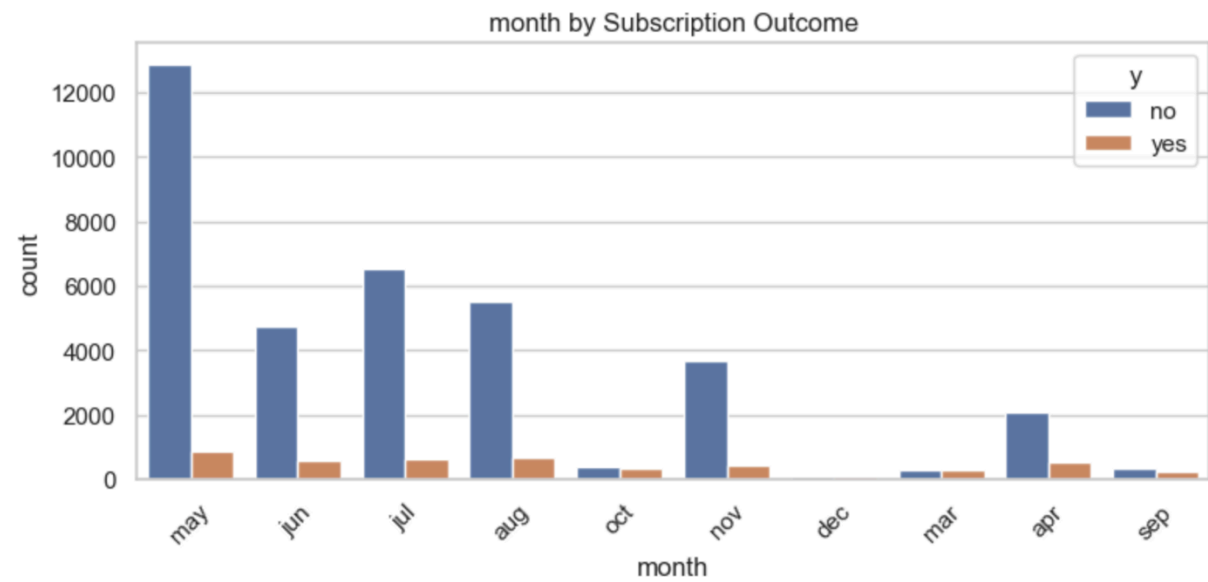
# Profile of Likely Subscribers

- Subscribers are often students or retired individuals.
- More likely with no personal/housing loans.
- Longer call durations correlate with success.

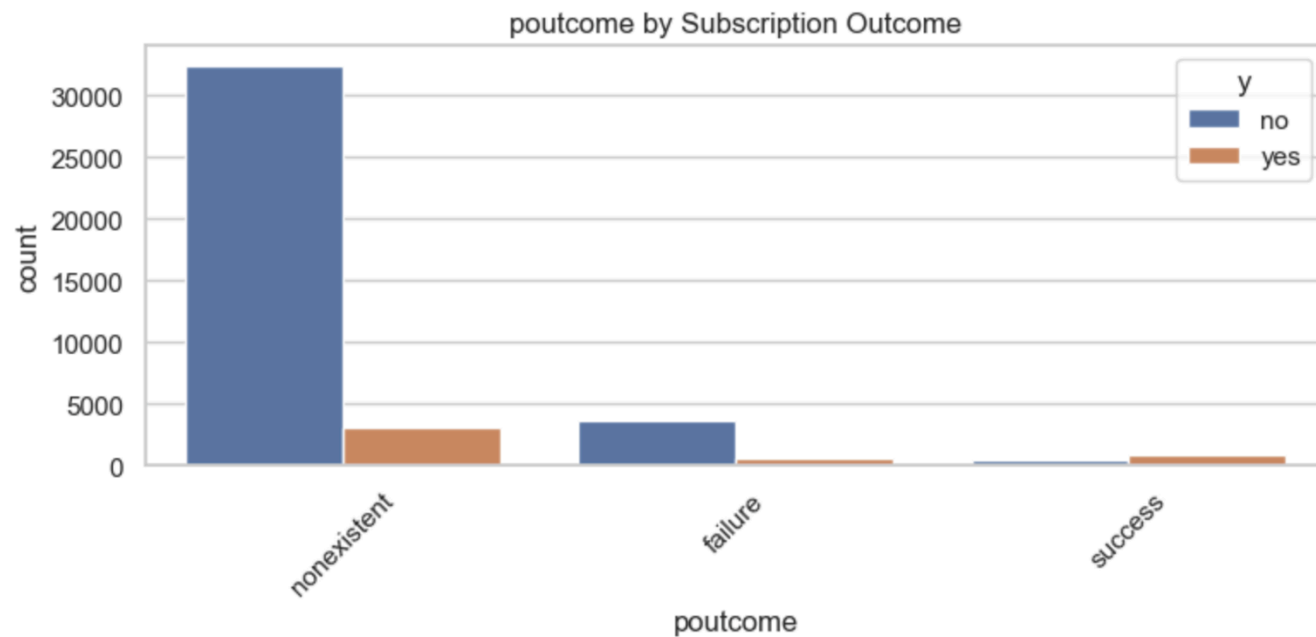


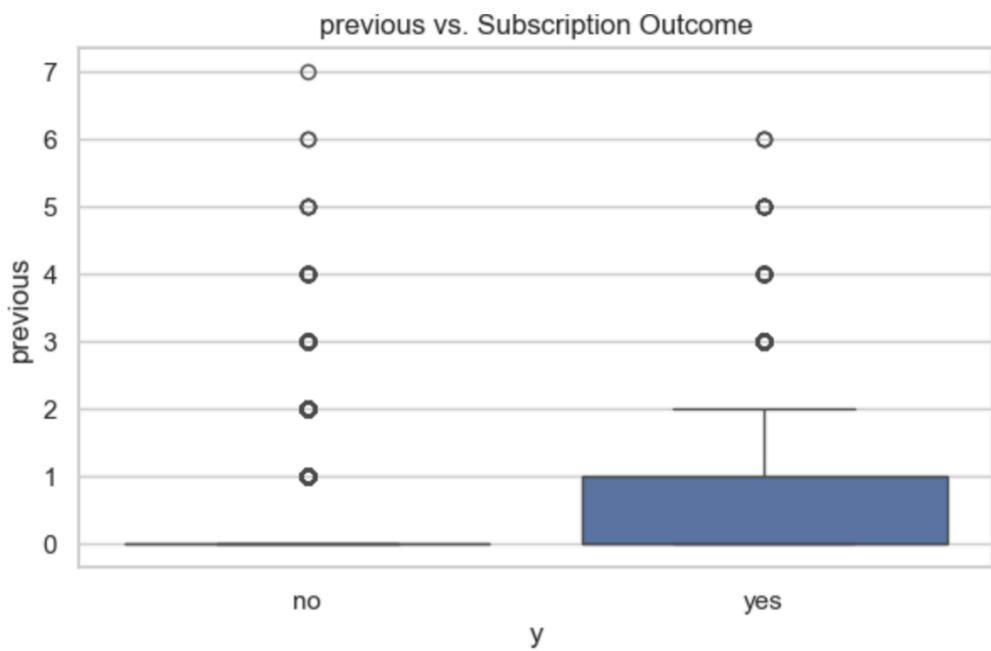
- The type of communication used to contact customers is mostly cellular compared to telephone.
- May seems to be the month with most contacts made.





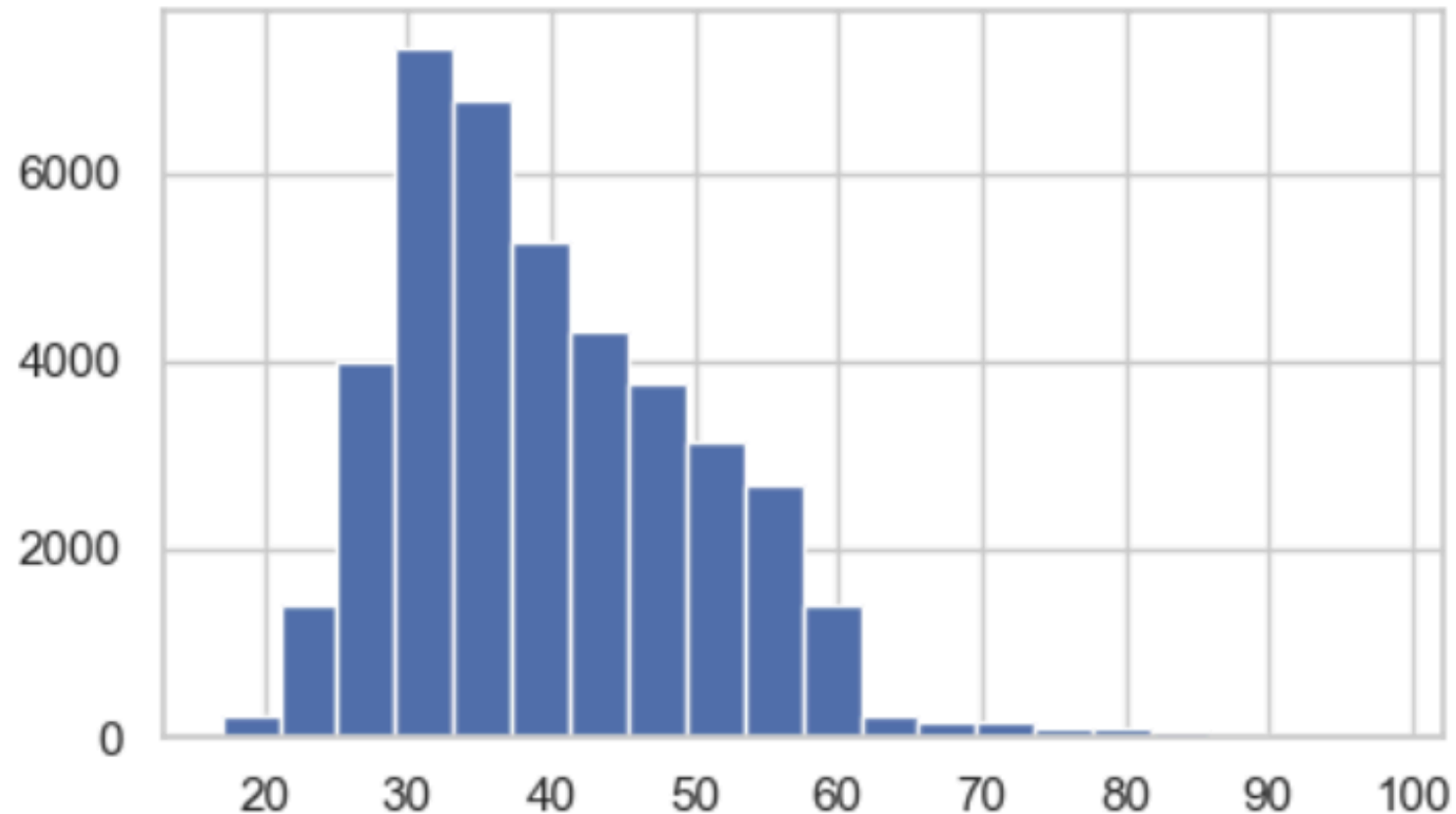
- Most of the previous campaign results whether succeeded or failed is nonexistent.
- In the previous campaign, the percentage of people who subscribed is less to those who did not.





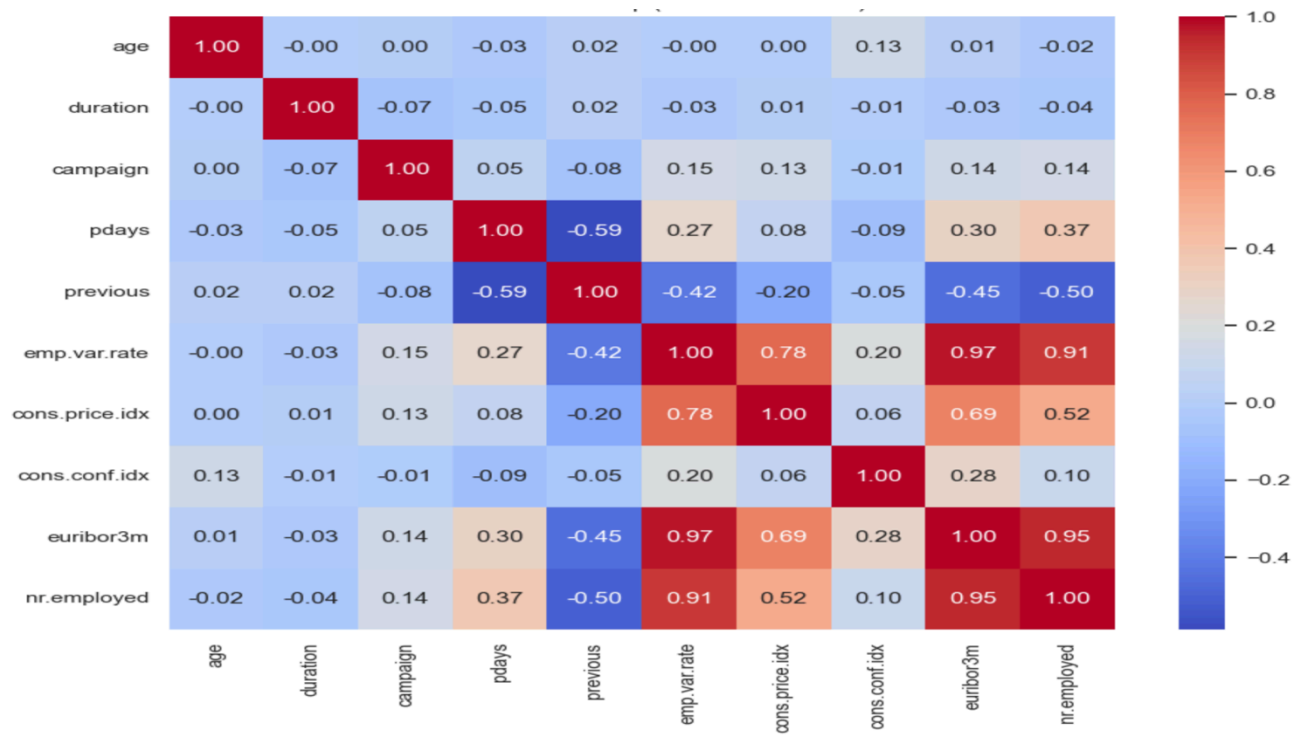


- The age group of customers contacted mostly fall between 20 to 60.
- Number of follow ups made to a customer is less in the previous campaign.



## Correlation map

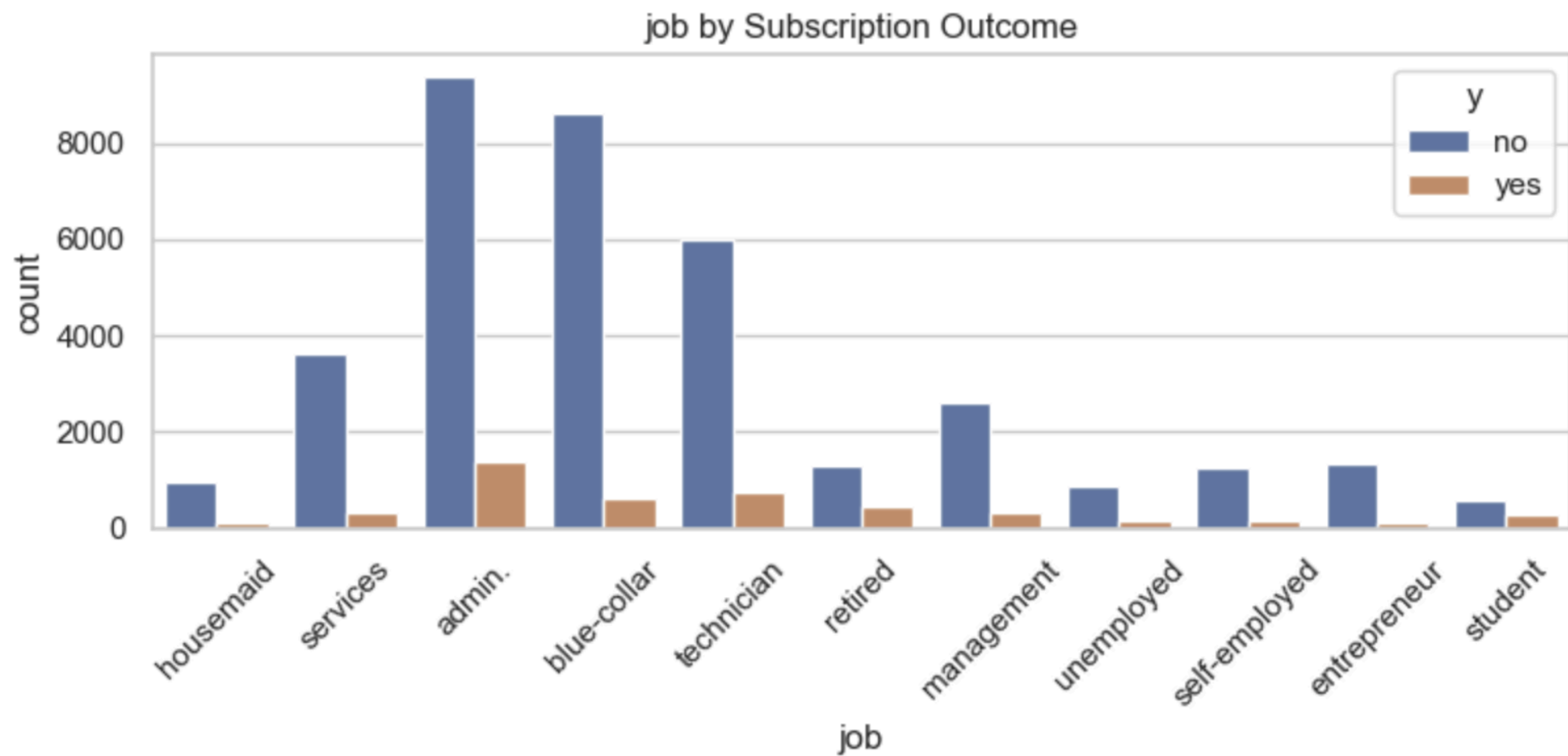
- Most features are not highly correlated, reducing multicollinearity concerns.



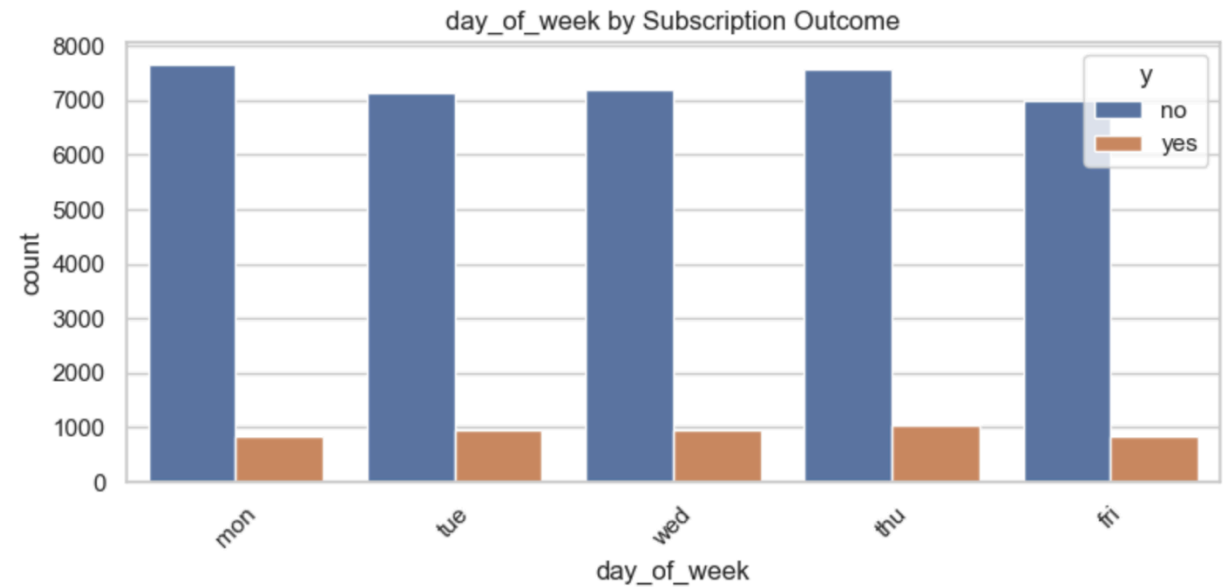
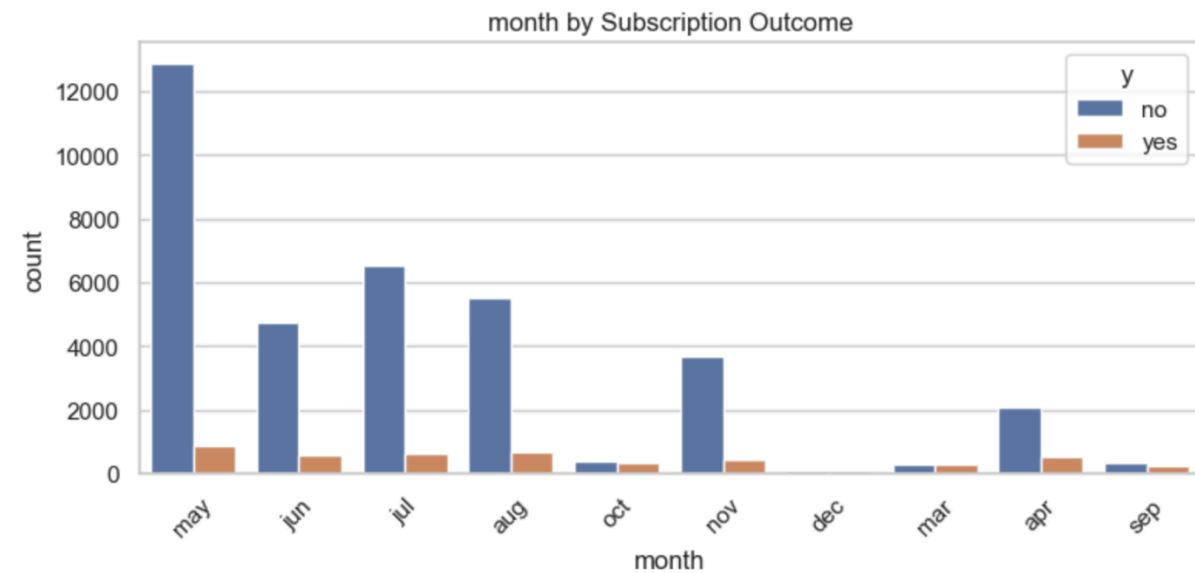
# EDA

## 2. Bivariate Analysis

- Classification of the Job Subscription Outcome.

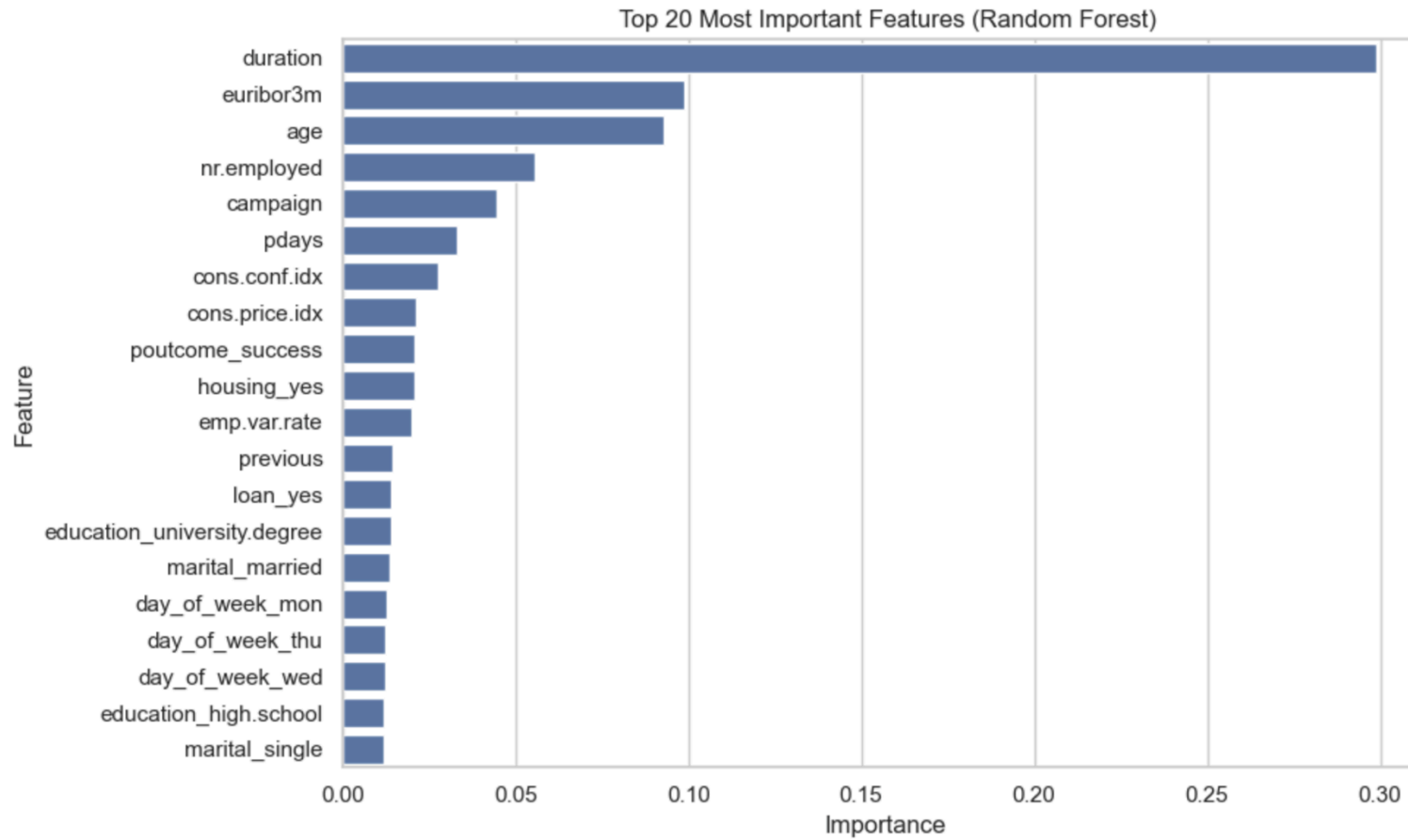


- Looking at the correlation between month by Subscription outcome and day-of-week.

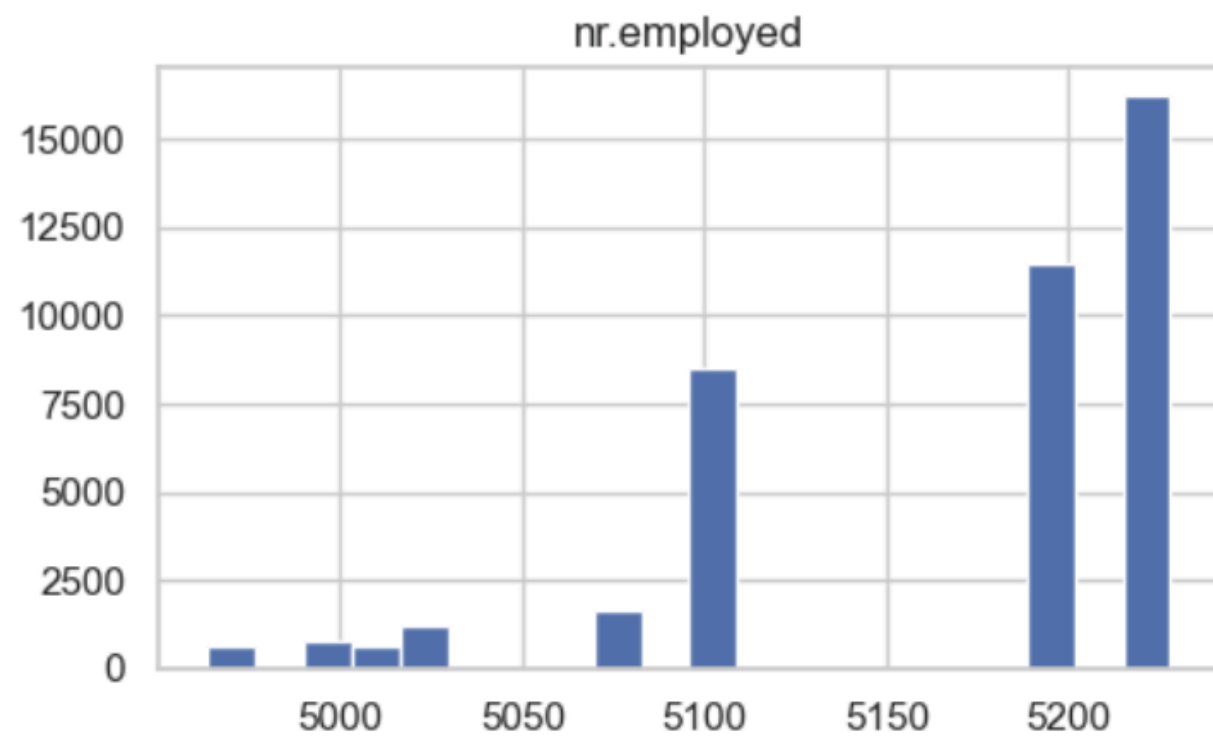
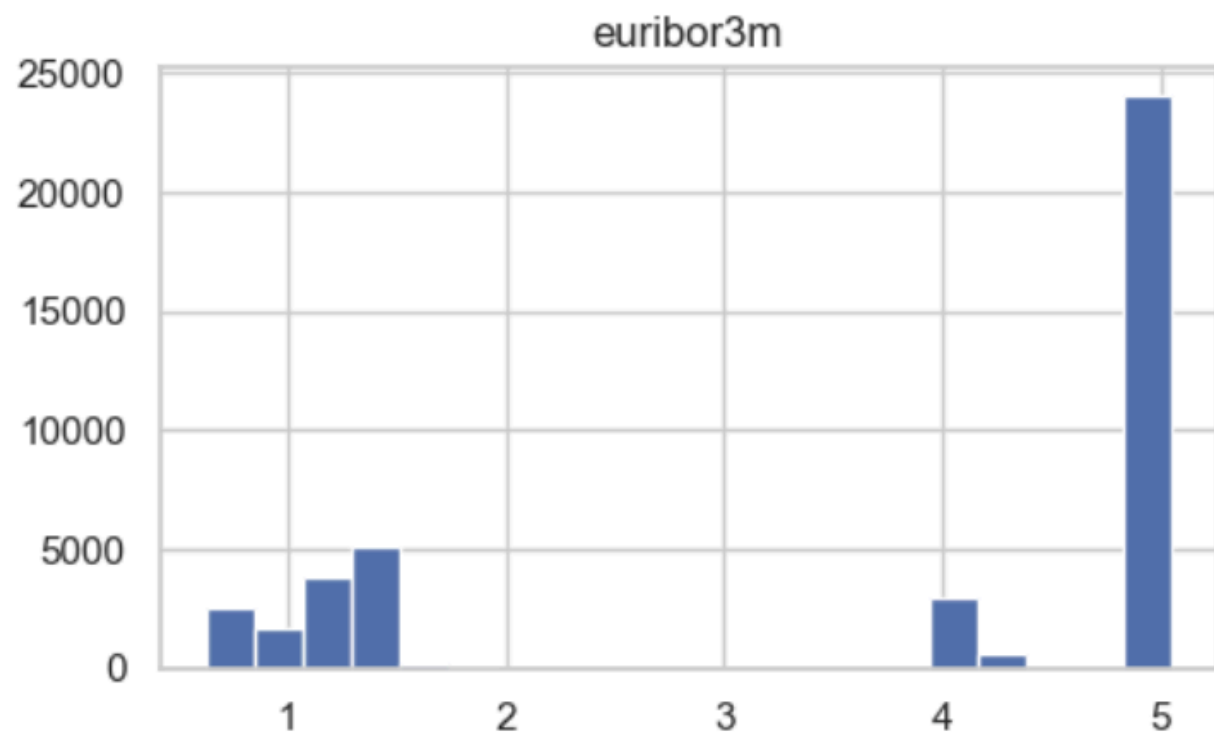




- Looking at relation The Top 20 Most Important Features



- Economic variables (e.g., euribor3m, nr.employed) have strong correlations and predictive potential.



## Recommendation to improve campaign

1. **May is the most effective month to contact customers**
2. **Increase the time of contacts made per customer**
3. **Give more focus on university graduate students and high school degree students**
4. **age-groups of 26-40 and 41-60 have a higher proportion among customers, therefore these groups present a**

**profitable target for the marketing team.**

- 5. Target the admins and technicians for more subscriptions**

## Recommended models for this dataset

**Problem Type: Binary classification with imbalance.**

1. Random Forest: Robust, interpretable
2. XGBoost: High accuracy, handles imbalance
3. Logistic Regression: Good, baselined Boost Decision Tree

**Techniques:**

4. - SMOTE or class weighting
5. - Exclude 'duration' for real-time prediction
6. Deploy with: sklearn, xgboost, joblib, Fast API



# Thank You