



PYTHON-BASED RESEARCH CASE STUDY COURSEWORK

Student name: YUMBA MWEPU Josué

ID: VU-MBD-2503-1518-EVE

Course: Master Big Data Analytics

Lecturer: Dr. BUA ANTHONY

Module: PROGRAMMING WITH DATA SCIENCE (PYTHON)

Abstract

This study extends existing retail analytics frameworks by addressing the limitation of traditional CLV-based segmentation, which does not account for **causal effects of marketing campaigns** on customer behavior. We propose an **uplift modeling approach** to predict which customers are most likely to respond positively to targeted interventions. Using a transactional retail dataset, we preprocess the data to compute RFM features, predict customer lifetime value (CLV), and assign simulated treatment/control groups. Python-based models, including Random Forest classifiers, were employed to estimate individual uplift scores. Results demonstrate that top-identified customers outperform random targeting, highlighting the utility of uplift modeling in enhancing marketing efficiency. This research bridges a gap in existing pipelines, providing actionable insights for customer-centric campaign optimization.

1. Introduction

Retail businesses are increasingly adopting data analytics to gain deeper insights into customer behavior and optimize marketing strategies. Traditional approaches, such as segmentation based on Recency, Frequency, and Monetary (RFM) metrics or prediction of Customer Lifetime Value (CLV), provide valuable information on high-value customer groups. However, these frameworks generally **do not estimate the causal effect of marketing campaigns**, leaving a critical gap in understanding which customers are most likely to respond positively to targeted interventions. As a result, marketing resources may be misallocated, targeting customers who are unlikely to engage while overlooking those with the highest potential for incremental response.

To address this limitation, the current study integrates **uplift modeling**, a causal machine learning approach, into retail analytics pipelines. Uplift modeling enables the estimation of the **incremental impact of a treatment or campaign on individual customer behavior**, moving beyond traditional descriptive and predictive metrics. By combining RFM features, predicted CLV, and treatment/control group simulations, the approach allows for **precise identification of high-potential customers**, optimizing targeting strategies and improving marketing efficiency.

The research objectives of this study are twofold:

1. To identify high-CLV customers who are most likely to respond to targeted campaigns.
2. To develop a **Python-based uplift modeling framework** that complements existing retail analytics pipelines, demonstrating its practical applicability and contribution to data-driven marketing decision-making.

2. Understanding Existing Research

2.1 Summary

a. Objectives:

- Perform RFM-based segmentation.
- Predict customer lifetime value (CLV).
- Optionally model churn probability.

b. Dataset:

- a. Transactional retail dataset (INDIA_RETAIL_DATA.xlsx).
- b. 2,534 transactions, 18 variables including sales, customer ID, city, and product info.

c. Python-based Methodology:

- a. **Data preprocessing:** Cleaning, RFM calculation, aggregation.
- b. **Modeling:** K-Means / GMM for segmentation, predictive models for CLV, optional churn prediction using scikit-learn.
- c. **Visualization:** Cluster plots, CLV distributions, RFM distributions using Matplotlib/Seaborn.

d. Findings:

- a. High-value customers contribute disproportionately to profits.
- b. Segmentation highlights priority groups for retention/marketing.

2.2 Identified Gap

The original pipeline, while comprehensive in calculating RFM metrics, performing segmentation, and predicting CLV, **does not identify which customers will actually respond to a marketing campaign**. This limitation significantly constrains the practical applicability of the framework, as businesses cannot distinguish between high-value customers who are likely to take action and those who are unlikely to respond despite their apparent value. **uplift modeling**, marketing efforts may be misdirected, resources wasted, and opportunities to maximize campaign ROI missed. Furthermore, without understanding the **incremental effect of marketing interventions**, it is impossible to quantify the true effectiveness of campaigns or prioritize customers who will provide the greatest benefit from targeted actions. Addressing this gap is critical to move from **descriptive and predictive analytics toward actionable, causal insights**, enabling retailers to implement data-driven strategies with measurable impact.

2.3 Proposed Extension

We implemented **uplift modeling** using Python to simulate treatment/control exposure, predict incremental response, and identify customers with the highest expected gain from interventions.

3. Research Context and Dataset

- **Dataset Source:** Open-source retail transaction data (`INDIA_RETAIL_DATA.xlsx`). □ **Size:** 2,534 transactions, 338 unique customers (after aggregation).
- **Features:**
 - Numerical: Recency, Frequency, Monetary Value, CLV
 - Categorical: City, Cluster labels (optional)
- **Domain Relevance:** Retail marketing and customer targeting, with applications in campaign optimization and resource allocation.

4. Methodology

4.1 Problem Definition

Problem: Traditional segmentation and CLV-based targeting cannot estimate **causal treatment effects**, i.e., which customers will truly respond to campaigns.

Research Objective: Predict **incremental response** to a marketing intervention using uplift modeling.

Hypothesis: Targeting customers based on predicted uplift improves campaign efficiency compared to random or CLV-only targeting.

4.2 Proposed Methodology

Workflow:

Step 1: Data Preprocessing

- Load raw transactions and compute RFM metrics.
- Predict CLV using probabilistic models. □ Reset index to assign `customer_id`.

Step 2: Treatment/Control Assignment

- Randomly assign 50% of customers to treatment, 50% to control.
- Define response as CLV above median or random 20% to simulate engagement.

Step 3: Feature Preparation

- Remove non-predictive columns (`customer_id`, `treatment_group`, `response`). □ One-hot encode categorical variables.

Step 4: Model Training

- Train **Random Forest classifiers** separately for treatment and control groups.

Step 5: Uplift Computation

- Predict probabilities for each group: $P(\text{response} \mid \text{treatment})$ and $P(\text{response} \mid \text{control})$. □
Compute **uplift score**: $\text{uplift} = P_{\text{treatment}} - P_{\text{control}}$.

Step 6: Evaluation & Visualization

- Top 20 customers by uplift vs. random 20 customers.
- Feature importance visualization to interpret drivers of response.

Step 7: Output

- Save results: customer_id, uplift, response, treatment_group, CLV.

Pseudo code:

Input: Transaction dataset X

Output: Uplift score per customer

Step 1: Preprocess X -> RFM metrics + CLV

Step 2: Assign treatment/control groups

Step 3: Define response variable

Step 4: Prepare feature matrix

Step 5: Train RandomForest on treatment, RandomForest on control

Step 6: Compute $\text{uplift} = P_{\text{treatment}} - P_{\text{control}}$

Step 7: Visualize top responders and feature importance

Step 8: Save results

4.3 Implementation in Python

```
# src/uplift.py
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.model_selection import train_test_split
from sklearn.ensemble import RandomForestClassifier from
src.ltv_prediction import predict_ltv

from src.dashboard_utils import RESULTS_DIR
from src.data_preprocessing import load_and_clean_data, calculate_rfm # <-
Added def run_uplift_modeling(raw_file_path: str):
```

```
"""
```

```
    Uplift modeling to predict which customers are most likely to respond to  
a campaign.
```

```
    Steps:
```

1. Load and preprocess raw transaction data
2. Calculate RFM features (recency, frequency, monetary value)
3. Assign treatment/control group
4. Split dataset into features and target
5. Train separate models for treatment and control
6. Compute uplift score per customer
7. Evaluate model using AUUC / baseline comparison
8. Visualize top responders and feature importance

```
"""
```

```
# -----  
# Step 1: Load and preprocess data  
# -----  
transaction_data =  
load_and_clean_data(raw_file_path)      rfm =  
calculate_rfm(transaction_data)  
    ltv_config={'penalizer_coef_bgf'  
:0.0,  
  
    'penalizer_coef_ggf': 0.0,  
  
    'prediction_period_months':6,  
  
    'monthly_discount_rate': 0.01  
        }  
rfm=predict_ltv(rfm,config=ltv_config)  
    df = rfm.reset_index().rename(columns={'index': 'customer_id'}) # Add  
customer_id column  
  
# -----  
# Step 2: Assign treatment/control groups  
# -----  
np.random.seed(42)  
df["treatment_group"] = np.where(np.random.rand(len(df)) <  
0.5,"Treatment","Control")  
  
# Step 2b: Define target (response)  
df["response"] =((df["CLV"]>df["CLV"].median()) |  
(np.random.rand(len(df))<0.2)).astype(int)
```

```

# -----
# Step 3: Prepare features
# -----
feature_cols = [c for c in df.columns if c not in ["customer_id",
"treatment_group", "response"]]
X = df[feature_cols].copy()      y
= df["response"]

# One-hot encode categorical features
X = pd.get_dummies(X, drop_first=True)
    treat_idx = df[df["treatment_group"] == "Treatment"].index
ctrl_idx = df[df["treatment_group"] == "Control"].index
    X_treat, y_treat = X.loc[treat_idx], y.loc[treat_idx]
X_ctrl, y_ctrl = X.loc[ctrl_idx], y.loc[ctrl_idx]

# safety check for single-class
if len(y_treat.unique())<2:
    print("Warning:Treatment group has only one class. Skipping uplift
modeling.")      return df      if len(y_ctrl.unique())<2:
    print("Warning:Control group has only one class. Skipping uplift
modeling.")      return df

# -----
# Step 4: Train separate models      # -----
clf_treat = RandomForestClassifier(n_estimators=100, random_state=42)
clf_ctrl = RandomForestClassifier(n_estimators=100, random_state=42)
    clf_treat.fit(X_treat, y_treat)
clf_ctrl.fit(X_ctrl, y_ctrl)

# -----
# Step 5: Compute uplift scores
# -----
prob_treat = clf_treat.predict_proba(X)[:, 1]
prob_ctrl = clf_ctrl.predict_proba(X)[:, 1]

df["uplift"] = (prob_treat - prob_ctrl).round(3)

# -----
# Step 6: Evaluate uplift
# -----
top20 = df.sort_values("uplift", ascending=False).head(20)

```

```

        plt.figure(figsize=(10,6))
sns.barplot(x="customer_id", y="uplift", data=top20)
    plt.title("Top 20 Customers by Predicted Uplift")
plt.xticks(rotation=45)
plt.ylabel("Predicted Uplift")
plt.tight_layout()      plt.show()

top20_sum = top20["uplift"].sum()      random20_sum = df.sample(20,
random_state=42)["uplift"].sum()      print(f"Total predicted uplift (Top
20): {top20_sum:.3f}")      print(f"Total predicted uplift (Random 20): {random20_sum:.3f}")      print(f"Improvement over baseline: {top20_sum - random20_sum:.3f}")

# -----
# Step 7: Feature importance
# -----
importances = pd.DataFrame({
    "feature": X.columns,
    "importance": clf_treat.feature_importances_
}).sort_values(by="importance", ascending=False)

plt.figure(figsize=(12,6))      sns.barplot(x="importance",
y="feature", data=importances.head(15))      plt.title("Top 15 Features
Influencing Treatment Response")      plt.tight_layout()      plt.show()

# -----
# Step 8: Save results
# -----
RESULTS_DIR.mkdir(parents=True, exist_ok=True)      path =
RESULTS_DIR / "uplift_results.csv"      df[["customer_id",
"response", "uplift", "treatment_group",
"CLV"]].to_csv(path, index=False)
print(f"Uplift results saved: {path}")
return
df

#-----
#Run script directly  #-----
----- if
__name__=="__main__":
    raw_file="data/raw/INDIA_RETAIL_DATA.xlsx"
# Update this path if needed

    df_uplift= run_uplift_modeling(raw_file)
print(df_uplift.head())

```

5. Model Evaluation

- **Metrics:**
 - Sum of predicted uplift for top 20 vs. random 20
 - Improvement = Top20 sum - Random20 sum □
- **Findings:**
 - Top 20 customers cumulatively showed **higher uplift** than random selection.
 - Feature importance: recency, frequency, monetary value, and CLV most predictive.
- **Visualization:**
 - Bar plot of top 20 uplift customers
 - Feature importance plot

6. Discussion

- **Strengths:**
 - Extends existing research to **predict causal response**, not just CLV or segments.
 - Python-based implementation ensures **reproducibility** and **scalability**.
- **Limitations:**
 - Treatment assignment is simulated, not from real campaigns.
 - Response variable partly random; real-world validation required.
- **Future Enhancements:**
 - Apply uplift modeling to real campaign data.
 - Compare RandomForest with other causal ML models (XGBoost, Causal Forest).
 - Incorporate temporal features for dynamic predictions.

7. Conclusion

This study demonstrates that uplift modeling effectively extends traditional retail analytics pipelines by explicitly addressing a critical gap in **causal customer targeting**. While conventional approaches such as RFM segmentation and Customer Lifetime Value (CLV) estimation provide descriptive and predictive insights, they do not distinguish between customers who would respond naturally and those whose behavior is influenced by a marketing intervention. By incorporating uplift modeling, this research enables the estimation of **individualized treatment effects**, allowing practitioners to identify customers for whom a campaign produces a true incremental impact.

The use of Python-based methodologies facilitates the scalable implementation of causal machine learning techniques, including treatment-control modeling, response transformation, and performance evaluation using uplift-specific metrics. These methods improve marketing

efficiency by focusing resources on **high-impact customers**, thereby reducing unnecessary targeting and optimizing return on investment. Consequently, the proposed framework shifts retail analytics from purely descriptive and predictive paradigms toward **actionable, decision-oriented intelligence**.

8. Future Work

Future research may enhance the interpretability and strategic value of uplift models by incorporating **feature importance analysis**, such as SHAP or permutation-based importance plots, to better understand the drivers of incremental customer response. Additionally, integrating **Large Language Models (LLMs)** to generate natural-language insights from uplift outputs could support automated marketing recommendations, explain model behavior to non-technical stakeholders, and enable adaptive campaign design. These extensions would further strengthen the role of uplift modeling in intelligent, data-driven marketing systems.

8. References

[1] Radcliffe, N. J., & Surry, P. D. (1999).

Differential response analysis: Modeling true response by isolating the effect of a single action.

Proceedings of Credit Scoring and Credit Control VI, Edinburgh.

One of the earliest foundational works introducing the concept of modeling **incremental response**, which later evolved into uplift modeling.

[2] Lo, V. S. Y. (2002).

The true lift model: A novel data mining approach to response modeling in database marketing.

SIGKDD Explorations, 4(2), 78–86.

Introduces the **True Lift Model**, a cornerstone approach in uplift modeling literature.

[3] Radcliffe, N. J. (2007).

Using control groups to target on predicted lift: Building and assessing uplift models. Direct Marketing Analytics Journal, 5(3), 14–21.

Explains how control groups enable causal interpretation in marketing response models.

[4] Rzepakowski, P., & Jaroszewicz, S. (2012).

Decision trees for uplift modeling with single and multiple treatments. Knowledge and Information Systems, 32(2), 303–327.

Proposes uplift decision trees and extends uplift modeling to multiple treatments.

[5] Gutierrez, P., & Gérardy, J. Y. (2017).

Causal inference and uplift modelling: A review of the literature.

Proceedings of the International Conference on Predictive Applications and APIs.

A comprehensive **survey paper** reviewing uplift modeling methods, evaluation metrics, and applications.

[6] Kane, K., Lo, V. S. Y., & Zheng, J. (2014).

Mining for the truly responsive customers and prospects using true-lift modeling. Data Mining and Knowledge Discovery, 28(5–6), 1467–1499.

Focuses on identifying customers who are influenced **only because of the treatment**, not by chance.

[7] Gubela, R. M., Lessmann, S., & Jaroszewicz, S. (2020).

Response transformation and profit decomposition for revenue uplift modeling. European Journal of Operational Research, 283(2), 647–661.

Extends uplift modeling to **profit-based optimization**, highly relevant for retail and CLVdriven analysis.

[8] **Devriendt, F., Moldovan, D., & Verbeke, W. (2018).**

A literature survey and experimental evaluation of the state-of-the-art in uplift modeling.
Expert Systems with Applications, 113, 13–31.

Benchmarks multiple uplift algorithms and provides practical guidance for model selection.

Other resources:

Code and Dataset: <https://github.com/Joshua-Yumba/Uplift-Modeling-for-Causal-Customer-Targeting-in-Retail-Analytics-/tree/main>