# Configuration Manual

MSc Research Project
Data Analytics

# Prajwal Keshav Kongi

Student ID: 22205314

School of Computing
National College of Ireland

Supervisor:     Dr. Catherine Mulwa

| Student Name: | Prajwal Keshav Kongi |
|---|---|
| Student ID: | 22205314 |
| Programme: | Data Analytics |
| Year: | 2023-2024 |
| Module: | MSc Research Project |
| Supervisor: | Dr. Catherine Mulwa |
| Submission Due Date: | 12/08/2024 |
| Project Title: | Configuration Manual |
| Word Count: | XXX |
| Page Count: | 10 |

I hereby certify that the information contained in this (my submission) is information pertaining to research I conducted for this project. All information other than my own contribution will be fully referenced and listed in the relevant bibliography section at the rear of the project.

**ALL** internet material must be referenced in the bibliography section. Students are required to use the Referencing Standard specified in the report template. To use other author's written or electronic work is illegal (plagiarism) and may result in disciplinary action.

| Signature: | |
|---|---|
| Date: | 16th September 2024 |

**PLEASE READ THE FOLLOWING INSTRUCTIONS AND CHECKLIST:**

| | |
|---|---|
| Attach a completed copy of this sheet to each project (including multiple copies). | ☐ |
| **Attach a Moodle submission receipt of the online project submission**, to each project (including multiple copies). | ☐ |
| **You must ensure that you retain a HARD COPY of the project**, both for your own reference and in case a project is lost or mislaid. It is not sufficient to keep a copy on computer. | ☐ |

Assignments that are submitted to the Programme Coordinator office must be placed into the assignment box located outside the office.

| **Office Use Only** | |
|---|---|
| Signature: | |
| Date: | |
| Penalty Applied (if applicable): | |

# Configuration Manual

### Prajwal Keshav Kongi
### 22205314

## 1 Introduction

The configuration manual outlines the system specifications, environment setups, and methodologies used in implementing uplift modeling to target potential customers in the fashion e-commerce domain in Russia. The research was carried out on Jupyter Notebook using Python. New packages and libraries for Uplift modeling were installed like scikit-uplift. Data processing and Feature engineering are carried out in this project. Four uplift models are trained, tested, and evaluated when there are multiple treatments in the data. Two-Model with CatBoost, LGBM classifier, Class Transformation, and T-learner with LGBM models are applied. This document provides detailed descriptions and information on the tools and technologies used in developing uplift modeling.

## 2 System Specifications



Figure 1: System Specifications

Figure 1 shows the details about the Device Specifications and Windows specifications used for this research.

Figure 2: Jupyter and Python Versions

Figure 2 shows the details about the Jupyter Notebook and Python versions used in implementing this project.

# 3 Fashion E-Commerce Dataset



Figure 3: Fashion E-Commerce Direct Messaging Data

Figure 3 shows the source of the E-commerce data used in the project which is derived from Kaggle platform.

# 4 Dataset Preparation

In this section, data pre-processing is carried out for both message and campaign datasets.

## Installing New Libraries

```
1  pip install pandas scikit-learn scikit-uplift xgboost
```

```
1  pip install --upgrade scikit-uplift
```

```
1  pip install catboost
```
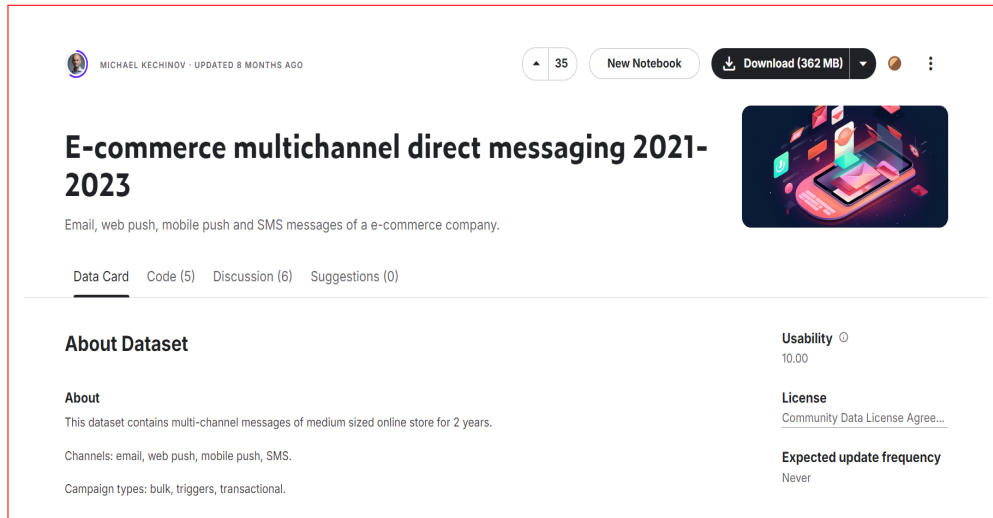
```
1  !{sys.executable} -m pip install scikit-uplift dill lightgbm
```

Figure 4: Installation of New libraries and packages

In Figure 4, new libraries and packages are installed for this project for the implementation of uplift modeling.

## Importing Libraries

```
1   import numpy as np
2   import pandas as pd
3
4   import matplotlib.pyplot as plt
5   import seaborn as sns
6
7   from pymongo import MongoClient
8   import sqlalchemy as db
9
10  import sys
11  from sklearn.model_selection import train_test_split
12  from sklift.models import TwoModels
13  from sklift.models import SoloModel
14  from sklift.models import ClassTransformation
15  from sklearn.ensemble import RandomForestClassifier
16  from catboost import CatBoostClassifier
17  from xgboost import XGBClassifier
18  from lightgbm import LGBMClassifier
19
20  from sklift.metrics import uplift_at_k, uplift_auc_score, qini_auc_score
21  from sklift.viz import plot_qini_curve
```

Figure 5: Libraries used in this project

Figure 5 shows the libraries used in this project.

```
1  # Converting string to int
2
3  df['campaign_id'] = df['campaign_id'].astype('Int32')
```

```
1  # Function to convert 't' or 'f' to boolean values (0 and 1)
2
3  def convert_to_bool(value):
4      if value == 't':
5          return True
6      else:
7          return False
8
9  convert_to_bool_cols = ['is_opened', 'is_clicked', 'is_unsubscribed', 'is_purchased']
10
11
12  df[convert_to_bool_cols] = df[convert_to_bool_cols].applymap(convert_to_bool)
```

```
1  # Converting date columns to datetime format
2
3  date_cols = ['date', 'sent_at', 'opened_first_time_at', 'opened_last_time_at', 'clicked_first_time_at', 'clicked_last_time_a
4              'unsubscribed_at', 'hard_bounced_at', 'soft_bounced_at', 'complained_at', 'blocked_at', 'purchased_at']
5
6  df[date_cols] = df[date_cols].apply(pd.to_datetime)
7
8  df[date_cols].head(2)
9
```

Figure 6: Fixing Datatypes

In Figure 6, the features are converted to appropriate datatypes using apply, and applymap functions.

```
1  # Replacing missing values with Mode of each column
2
3  mode_cols = ['subject_with_personalization', 'subject_with_deadline', 'subject_with_emoji',
4              'subject_with_bonuses', 'subject_with_discount', 'subject_with_saleout']
5
6  for column in mode_cols:
7      mode_value = df3[column].mode()[0]  # Getting the most frequent value
8      df3[column] = df3[column].fillna(mode_value)
```

```
1  # Filling subject_length feature with mean
2
3  df3['subject_length'] = df3['subject_length'].fillna(df3['subject_length'].mean())
4
5  df3['total_count'] = df3['total_count'].fillna(df3['total_count'].mean())
```

Figure 7: Handling Missing Values

The handling of missing values for both Numerical and Categorical variables is shown in Figure 7.

```
1  # Taking only top 7 email_providers and renaming remaining email_providers with 'others'
2
3  df['email_provider'] = df['email_provider'].apply(lambda x: x if x in ['mail.ru', 'gmail.com', 'yandex.ru', 'bk.ru',
4                                                      'list.ru', 'inbox.ru', 'rambler.ru'] else 'others')
```

Figure 8: Grouping Email Providers

Figure 8 shows the apply function used with lambda to keep only the top 7 email providers and group the remaining email providers as 'others'.

**Pre-processing of Campaigns Data**

```python
1  df_camp1 = df_camp.copy()
2
3  # Dropping null values from total_count feature
4  df_camp1 = df_camp1.dropna(subset=['total_count'])
5
6  # Removing duplicates
7  df_camp1 = df_camp1.drop_duplicates(subset=['id'])
8
9  # From ab_test column, we'll get to know which users were in treatment and control groups
10 treatment_campaign_ids = list(df_camp1[df_camp1['ab_test'] == True]['id'].unique())
11
12 # Dropping irrelevant columns
13 df_camp1.drop(columns = ['topic', 'started_at', 'finished_at', 'ab_test', 'hour_limit',
14                          'is_test', 'position'], inplace=True, axis=1)
15
16 # Fixing datatypes
17 # Convert boolean columns to 0 and 1
18
19 camp_binary_cols = ['warmup_mode', 'subject_with_personalization', 'subject_with_deadline',
20                     'subject_with_emoji', 'subject_with_bonuses', 'subject_with_discount', 'subject_with_saleout']
21
22 df_camp1[camp_binary_cols] = df_camp1[camp_binary_cols].apply(lambda x: x.astype(int))
23
24 # Merging Message and Campaigns Data
25 df3 = pd.merge(df3, df_camp1, left_on='campaign_id', right_on='id', how='left')
26
```

Figure 9: Processing of Campaigns Data

Figure 9 outlines the key pre-processing steps undertaken for campaign data.

# 5 Feature Engineering

```python
1  # Applying log transformation for total_count column
2
3  df_camp_log_transformed = df_camp4.copy()
4  df_camp_log_transformed['total_count'] = np.log1p(df_camp_log_transformed['total_count'])
```

```python
1  # Function to clip outliers based on the 5th and 95th percentile for subject_length
2
3  def cap_outliers(df, column):
4      lower_cap = df[column].quantile(0.05)
5      upper_cap = df[column].quantile(0.95)
6      df[column] = df[column].clip(lower=lower_cap, upper=upper_cap)
7      return df
8
9  df_camp_capped = cap_outliers(df_camp4, 'subject_length')
```

Figure 10: Treating Outliers

In Figure 10, functions to detect and treat outliers are shown. The methods used to treat outliers of variables are Log Transformation and Clipping outliers using percentiles.

```
1  # Creating new column with the difference in seconds between message 'sent_at' and 'opened_at', 'clicked_at', 'purchased_at'
2
3  df1['opened_seconds'] = (df1['opened_first_time_at'] - df1['sent_at']).dt.total_seconds()
4  df1['clicked_seconds'] = (df1['clicked_first_time_at'] - df1['sent_at']).dt.total_seconds()
5  df1['purchased_seconds'] = (df1['purchased_at'] - df1['sent_at']).dt.total_seconds()
```

```
1  # Replacing NaN values in these columns with 0
2
3  df1['opened_seconds'].fillna(0, inplace=True)
4  df1['clicked_seconds'].fillna(0, inplace=True)
5  df1['purchased_seconds'].fillna(0, inplace=True)
```

```
1   # Segregating the seconds column's data into 10 segments
2
3   num_bins = 10
4
5   df1['opened_recency'] = pd.cut(df1['opened_seconds'],
6                                  bins=num_bins, labels=list(range(num_bins-1, -1, -1)), include_lowest=True)
7
8   df1['clicked_recency'] = pd.cut(df1['clicked_seconds'],
9                                  bins=num_bins, labels=list(range(num_bins-1, -1, -1)), include_lowest=True)
10
11  df1['purchased_recency'] = pd.cut(df1['purchased_seconds'],
12                                  bins=num_bins, labels=list(range(num_bins-1, -1, -1)), include_lowest=True)
```

```
1  # Manually filling 0 where message opened_at, clicked_at and purchased_at are null values
2
3  df1.loc[df['opened_first_time_at'].isna(), 'opened_recency'] = 0
4  df1.loc[df['clicked_first_time_at'].isna(), 'clicked_recency'] = 0
5  df1.loc[df['purchased_at'].isna(), 'purchased_recency'] = 0
```

Figure 11: Creating Recency Score

Figure 11 shows the steps taken in the creation of recency scores for 3 features ie, message opened, clicked, and purchased.

```
1  # One-hot Encoding
2
3  df3.drop(columns = ['warmup_mode', 'platform', 'id', 'campaign_type', 'channel_y'], inplace=True, axis=1)
4
5  df4 = df3.copy()
6  df4 = pd.get_dummies(df4, columns = ['channel_x', 'message_type', 'email_provider'], drop_first=True, dtype = int)
```

Figure 12: One-hot Encoding

One-hot encoding is implemented in Figure 12, to transform categorical features into numerical features using the 'get_dummies' function.

# 6    Modeling

In this research, four uplift models are trained, tested, and fine-tuned. The following contains the codes of four uplift models implemented on the fashion e-commerce dataset.

## 6.1 Two Model with CatBoost



Figure 13: Two Model approach with CatBoost

Figure 13 shows the implementation of the Two-model approach using the CatBoost classifier.

## 6.2 Class Transformation



Figure 14: Class Transformation Model

In Figure 14, Class Transformation using the CatBoost Classifier is trained on the dataset Zhao and Harinen (2019). The pipeline is used in this code for smooth implementation. The data is split using stratify split method where treatment and control groups are equally distributed among the train and test datasets.

```
3. Two Model Approach with LGBM Classifier

1  X = df4.drop(columns=['is_clicked'])
2  y = df4['is_clicked']
3  treatment = df4['treatment']
4
5  # Stratify split
6  stratify_cols = pd.concat([treatment, y], axis=1)
7  X_train, X_val, trmnt_train, trmnt_val, y_train, y_val = train_test_split(X,
8                                                                              treatment,
9                                                                              y,
10                                                                             stratify=stratify_cols,
11                                                                             test_size=0.3,
12                                                                             random_state=31)
13 print(f"Train shape: {X_train.shape}")
14 print(f"Validation shape: {X_val.shape}")

Train shape: (700000, 28)
Validation shape: (300000, 28)

1  # Initialize the models with regularization
2  treatment_model = LGBMClassifier(
3                      random_state=31,
4                      n_estimators=100,
5                      learning_rate=0.05,
6                      lambda_l1=1.0,  # L1 regularization
7                      lambda_l2=1.0,  # L2 regularization
8                      min_split_gain=0.01,
9                      min_child_weight=1,
10                     subsample=0.8,
11                     colsample_bytree=0.8,
12                     objective='binary')
13
14 control_model = LGBMClassifier(
15                     random_state=31,
16                     n_estimators=100,
17                     learning_rate=0.05,
18                     lambda_l1=1.0,
19                     lambda_l2=1.0,
20                     min_split_gain=0.01,
21                     min_child_weight=1,
22                     subsample=0.8,
23                     colsample_bytree=0.8,
24                     objective='binary')
25
26 # Initialize the TwoModels approach
27 tm = TwoModels(estimator_trmnt=treatment_model, estimator_ctrl=control_model, method='vanilla')
28
29 # Fit and Predict
30 tm = tm.fit(X_train, y_train, trmnt_train)
31 uplift_tm = tm.predict(X_val)
```

Figure 15: Two Model with LGBM

Figure 15 shows the code of Two model approach using LGBM Classifier. Regularization techniques and fine-tuning of the models for both treatment and control groups is carried out in this code.

**4. T-Learner (Meta Learner) with LGBM Regressor**

**Treatment**

```
1   # Treated Units
2   df_treated = df4[df4['treatment'] == 1]
3
4   # Features
5   features_treated = df_treated.drop(columns=['campaign_id', 'is_clicked', 'treatment'], axis = 1)
6
7   # Target action
8   y_treated = df_treated.loc[:, ['is_clicked']]
```

**Control**

```
1   # Control Units
2   df_control = df4[df4['treatment'] == 0]
3
4   # Features
5   features_control = df_control.drop(columns=['campaign_id', 'is_clicked', 'treatment'], axis = 1)
6
7   # Target action
8   y_control = df_control.loc[:, ['is_clicked']]
```

```
1   # features for all the samples
2
3   features = df4.drop(columns=['campaign_id', 'is_clicked', 'treatment'], axis = 1)
```

```
1   # LGBM Regressior
2
3   optimized_lgbm = LGBMRegressor(
4       random_state=42,
5       n_estimators=500,
6       learning_rate=0.05,
7       lambda_l1=1.0,            # L1 regularization
8       lambda_l2=1.0,            # L2 regularization
9       subsample=0.8,            # Subsample to prevent overfitting
10      colsample_bytree=0.8,     # Feature sampling
11      early_stopping_round=10   # Early stopping
12  )
```

```
1   # Fit and Predict
2   t_treated = optimized_lgbm.fit(features_treated, y_treated, eval_set=[(features, df4['is_clicked'])])
3   t_control = optimized_lgbm.fit(features_control, y_control, eval_set=[(features, df4['is_clicked'])])
```
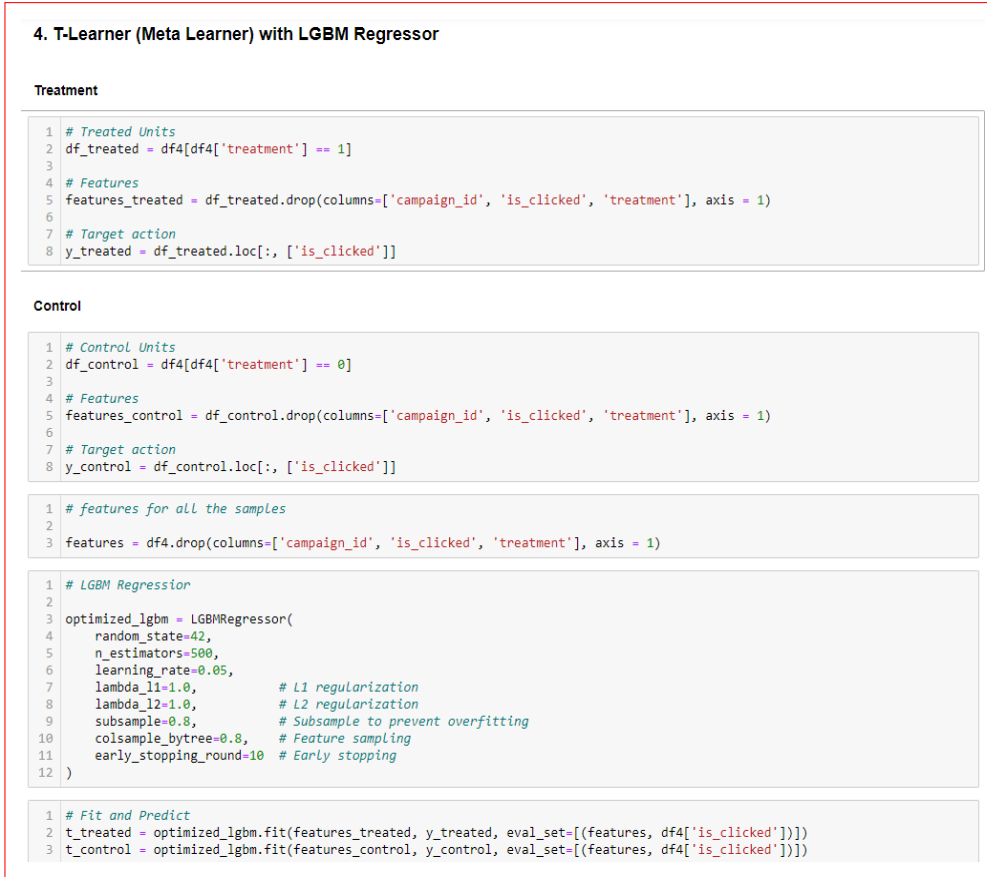
Figure 16: T-Learner with LGBM

In Figure 16, the code outlines the implementation of T-learner with LGBM on both treatment and control groups. The model is fine-tuned with various methods to overcome over-fitting.

# 7    Evaluation

In this research, three evaluation metrics are used to evaluate the uplift models. Uplift Score at 30%, AUUC and AUQC scores.

```
1   # Evaluation
2
3   uplift_score = uplift_at_k(y_test, uplift_two_model, t_test, strategy='overall', k=0.3)
4   auuc = uplift_auc_score(y_test, uplift_two_model, t_test)
5   auqc = qini_auc_score(y_test, uplift_two_model, t_test)
6
7   print(f'Two-Model Approach - Uplift score at 30%: {uplift_score}')
8   print(f'Two-Model Approach - Area Under Uplift Curve (AUUC): {auuc}')
9   print(f'Two-Model Approach - Area Under Qini Curve (AUQC): {auqc}')
10
11
12  # Plotting Uplift Score at 30%
13  plt.figure(figsize=(6, 4))
14  plt.bar('Uplift Score at 30%', uplift_score, color='skyblue')
15  plt.xlabel('Evaluation Metric')
16  plt.ylabel('Score')
17  plt.title('Two Models: Uplift Score at 30%')
18  plt.ylim(0, uplift_score + 0.05)  # Adjust y-axis limit
19  plt.grid(axis='y', linestyle='--', alpha=0.7)
20  plt.show()
21
```

Figure 17: Evaluation Metrics

Figure 17 shows the code for calculating the uplift score at 30%, AUUC, and AUQC scores. Additionally, the uplift score is plotted using the matplotlib library.

```python
1  # Plot the Qini curve
2
3  plt.figure(figsize=(8, 6))
4  plot_qini_curve(y_val, uplift_tm, trmnt_val, perfect=True)
5  plt.title('Qini Curve')
6  plt.xlabel('Number of Targets (Percentile)')
7  plt.ylabel('Cumulative Gain')
8  plt.grid(True)
9  plt.show()
```

Figure 18: Qini Curve

Figure 18 shows the code to plot the Qini curve using 'sklift.viz' library and from 'plot_qini_curve' package.

# References

Zhao, Z. and Harinen, T. (2019). Uplift modeling for multiple treatments with cost optimization, *2019 IEEE International Conference on Data Science and Advanced Analytics (DSAA)*, pp. 422–431.