

Group FTR

Week 8: Deliverables

Team Details

Name	Email	Country	College/Company	Specialization
Fabian Umeh	Fabianumeh335@gmail.com	UK	Teesside University	Data Science
Rukevwe Ovuowo	rukevwe10@gmail.com	Nigeria	GBG Data science Academy	Data Science
Olutayo Oladeinbo	oladeinboolutayo@yahoo.com	UK	Teesside University	Data Science

Problem statement

Churn rate is a marketing metric that describes the number of customers who leave a business over a specific time. Every user is assigned a prediction value that estimates their state of churn at any given time.

Business Understanding

Browsing behaviour Historical purchase data among other information It factors in our unique and proprietary predictions of how long a user will remain a customer. This score is updated every day for all users who have a minimum of one conversion. The values assigned are between 1 and 5.

Project lifecycle

Two weeks—deadline (1/09/2022)

Data intake report

Name: Customer Churn score prediction

Report date: 18/08/2022

Internship Batch: LISUM11: 30

Version:<1.0>

Data intake by: Fabian Umeh, Rukevwe Ovuowo, and Olutayo Oladeinbo

Data intake reviewer: Group members

Data storage location: [Github](#)

Tabular data details:

Total number of observations: 36992

Total number of files: 1

Total number of features: 25

Base format of the file: .csv

Size of the data: 8.3 MB

PROBLEMS IN DATA

1.1 Missing values:

Some values in columns such as [region_category, preferred_offer_types, points_in_wallet], appears to be missing.

1.2 Approach:

- The region category was encoded as follows:

City – 3

Town – 2

Village – 1

And as such, missing values were filled with the non-extreme value (2) for town.

- The points_in_wallet column had a relatively small percentage of missing values, and these values were filled using the average of the collected samples.

- The Preferred offer type column was encoded as follows:

Offer – 1

Without offer – 0

And as such the few missing values were replaced with no offers.

```
[ ] new= new.fillna({'region_category' : 2, 'points_in_wallet' : new['points_in_wallet'].mean(),  
                    'preferred_offer_types' : 0})
```



data.info()



```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 36992 entries, 0 to 36991
Data columns (total 25 columns):
#   Column                                Non-Null Count  Dtype
---  -
0   customer_id                          36992 non-null  object
1   Name                                  36992 non-null  object
2   age                                   36992 non-null  int64
3   gender                               36992 non-null  object
4   security_no                          36992 non-null  object
5   region_category                      31564 non-null  object
6   membership_category                 36992 non-null  object
7   joining_date                        36992 non-null  object
8   joined_through_referral             36992 non-null  object
9   referral_id                         36992 non-null  object
10  preferred_offer_types               36704 non-null  object
11  medium_of_operation                 36992 non-null  object
12  internet_option                     36992 non-null  object
13  last_visit_time                     36992 non-null  object
14  days_since_last_login               36992 non-null  int64
15  avg_time_spent                      36992 non-null  float64
16  avg_transaction_value               36992 non-null  float64
17  avg_frequency_login_days            36992 non-null  object
18  points_in_wallet                    33549 non-null  float64
19  used_special_discount               36992 non-null  object
20  offer_application_preference        36992 non-null  object
21  past_complaint                      36992 non-null  object
22  complaint_status                   36992 non-null  object
23  feedback                           36992 non-null  object
24  churn_risk_score                    36992 non-null  int64
dtypes: float64(3), int64(3), object(19)
memory usage: 7.1+ MB
```

2.1 String Error:

The data type in the average frequency login days column was supposed to be 'int' but rather contained some string value 'Error', which was further removed.

2.2 Approach:

The average frequency login days column keeps record of the average days the company site is visited by a customer. So ideally 'Error', most likely represent 0 login days was replaced as such.

Inspection of avg_frequency_login_days for possible anomaly:

From findings, the column contains 'int' and 'str' values, and as a result was not considered by the describe function above.

```
[ ] #Detection of str value from avg_frequency_login_days column
print(set([data for data in data['avg_frequency_login_days'] for a in data if a \
    not in ['0','1','2','3','4','5','6','7','8','9','.','-']]))

{'Error'}
```

```
[ ] #replacing error(str) values with 0
data = data.replace({'avg_frequency_login_days': {'Error': 0}})
#converting to integer
data['avg_frequency_login_days'] = data['avg_frequency_login_days'].astype('float64')
```

3.1 Negative values:

Some columns contained negative values, which appeared to be anomalous and were further excluded from the analysis.

3.2 Approach:

The negative (anomalous) values were few and were excluded from the analysis.

Check and removal of negative values (anomalous values):

```
#code to check for the negative values in the Dataframe we noticed from the describe function
anom = data[['avg_time_spent', 'days_since_last_login', 'points_in_wallet', 'churn_risk_score', 'avg_frequency_login_days']].min(axis=0)
anom[anom < 0]
```

avg_time_spent	-2814.109110
days_since_last_login	-999.000000
points_in_wallet	-760.661236
churn_risk_score	-1.000000
avg_frequency_login_days	-43.652702
dtype:	float64

```
[ ] #code to remove negative values in the avg_time_spent column
data = data.drop(data[data['avg_time_spent'] < 0].index).copy()

#code to remove negative values in the avg_time_spent column
data = data.drop(data[data['days_since_last_login'] < 0].index).copy()

#code to remove negative values in the avg_time_spent column
data = data.drop(data[data['points_in_wallet'] < 0].index).copy()

#code to remove negative values in the avg_time_spent column
data = data.drop(data[data['churn_risk_score'] < 0].index).copy()

#code to remove negative values in the avg_time_spent column
new_data = data.drop(data[data['avg_frequency_login_days'] < 0].index).copy()
```

4.1 Duplicate check:

Duplicate check:

```
[ ] print('counting duplicates')
len(new_data) - len(new_data.drop_duplicates())

counting duplicates
0
```

