

```
[1]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import datascist as ds

In [2]: train_df=pd.read_csv(r"C:\Users\ATN-PC\Desktop\kvc zeph\it_spend_files\Train.csv")

In [3]: test_df = pd.read_csv(r"C:\Users\ATN-PC\Desktop\kvc zeph\it_spend_files\Test.csv")

In [4]: #ds.structdata.describe(train_df)

checking columns to drop

In [5]: col_drop = []
for col in train_df.columns:
    if col not in test_df.columns:
        col_drop.append(col)
        print(col)
print(col_drop)

total_cost
['total_cost']

In [6]: train_df.columns

Out[6]: Index(['ID', 'country', 'age_group', 'travel_with', 'total_female',
'total_male', 'purpose', 'main_activity', 'info_source',
'tour_arrangement', 'package_transport_int', 'package_accommodation',
'tour_arrangement', 'package_transport_int', 'package_accommodation',
'package_food', 'package_transport_tz', 'package_sightseeing',
'package_guided_tour', 'package_insurance', 'night_mainland',
'night_zanzibar', 'payment_mode', 'first_trip_tz', 'most_impressing',
'total_cost'],
      dtype='object')

In [7]: train_df.drop(columns="ID",inplace=True)
test_df.drop(columns="ID",inplace=True)

In [8]: train_df.shape,test_df.shape

Out[8]: ((4809, 22), (1601, 21))

In [9]: train_df.isnull().sum()

Out[9]: country          0
age_group            0
travel_with         1114
total_female         3
total_male           5
purpose              0
main_activity        0
info_source          0
tour_arrangement     0
package_transport_int 0
package_accommodation 0
package_food         0
package_transport_tz  0
package_sightseeing  0
package_guided_tour  0
package_insurance    0
night_mainland       0
night_zanzibar       0
payment_mode         0
first_trip_tz        0
most_impressing      313
total_cost           0
dtype: int64

filling missing values

In [10]: #train_df=ds.feature_engineering.fill_missing_num(train_df,method ="mean")

In [11]: train_df[["most_impressing"].fillna(method ="bfill",inplace =True)
train_df[["total_female"].fillna(train_df[["total_female"].mode()[0],inplace =True)
train_df[["total_male"].fillna(train_df[["total_male"].mode()[0],inplace =True)
train_df[["travel_with"].fillna(method ="bfill",inplace =True)

test_df[["most_impressing"].fillna(method ="bfill",inplace =True)
test_df[["total_female"].fillna(test_df[["total_female"].mode()[0],inplace =True)
test_df[["total_male"].fillna(test_df[["total_male"].mode()[0],inplace =True)
test_df[["travel_with"].fillna(method ="bfill",inplace =True)

In [12]: #train_df.info()

processing of categorical data

In [13]: # join train_df and test_df together so as to perform the OHE or getdummies or labe encoder once
all_data,ntrain,ntest = ds.structdata.join_train_and_test(train_df,test_df)

In [14]: ds.structdata.get_cat_feats(train_df)

Out[14]: ['country',
'age_group',
'travel_with',
'purpose',
'main_activity',
'info_source',
'tour_arrangement',
'package_transport_int',
'package_accommodation',
'package_food',
'package_transport_tz',
'package_sightseeing',
'package_guided_tour',
'package_insurance',
'payment_mode',
'first_trip_tz',
'most_impressing']
```

those with two options(like yes and No) as obtained from class count

```
'tour_arrangement','package_transport_int','package_accommodation','package_food','package_transport_tz','package_sightseeing','package_guided_tour','package_insurance',
'first_trip_tz'

In [15]: #ds.structdata.class_count(train_df)

In [16]: #label encode those categories that are too large and use get dummies for ones that are like two features
from sklearn.preprocessing import LabelEncoder
lb =LabelEncoder()
cols_2_encode = ['country',
'age_group',
'travel_with',
'purpose',
'main_activity',
'info_source',
'payment_mode',
'most_impressing']
for col in cols_2_encode:
    lb.fit(all_data[col])
    all_data[col] = lb.transform(all_data[col])

In [17]: all_data.head().T

Out[17]:
```

	0	1	2	3	4
country	101	111	111	111	17
age_group	3	2	2	2	0
travel_with	2	0	0	3	0
total_female	1.0	1.0	0.0	1.0	1.0
total_male	1.0	0.0	1.0	1.0	0.0
purpose	1	1	5	1	1
main_activity	7	3	3	7	7
info_source	0	7	0	5	5
tour_arrangement	Independent	Independent	Independent	Package Tour	Independent
package_transport_int	No	No	No	No	No
package_accommodation	No	No	No	Yes	No
package_food	No	No	No	Yes	No
package_transport_tz	No	No	No	Yes	No
package_sightseeing	No	No	No	Yes	No
package_guided_tour	No	No	No	Yes	No
package_insurance	No	No	No	No	No
night_mainland	13.0	14.0	1.0	11.0	7.0
night_zanzibar	0.0	7.0	31.0	0.0	4.0
payment_mode	0	0	0	0	0
first_trip_tz	No	Yes	No	Yes	Yes
most_impressing	2	6	1	2	4
total_cost	674602.5	3214906.5	3315000.0	7790250.0	1657500.0

```


In [18]: # using getdummies from OHE to perform other at once
all_data = pd.get_dummies(all_data)

In [19]: all_data.head().T

Out[19]:
```

	0	1	2	3	4
country	101.0	111.0	111.0	111.0	17.0
age_group	3.0	2.0	2.0	2.0	0.0
travel_with	2.0	0.0	0.0	3.0	0.0
total_female	1.0	1.0	0.0	1.0	1.0
total_male	1.0	0.0	1.0	1.0	0.0
purpose	1.0	1.0	5.0	1.0	1.0
main_activity	7.0	3.0	3.0	7.0	7.0
info_source	0.0	7.0	0.0	5.0	5.0
night_mainland	13.0	14.0	1.0	11.0	7.0
night_zanzibar	0.0	7.0	31.0	0.0	4.0
payment_mode	0.0	0.0	0.0	0.0	0.0
most_impressing	2.0	6.0	1.0	2.0	4.0
total_cost	674602.5	3214906.5	3315000.0	7790250.0	1657500.0
tour_arrangement_Independent	1.0	1.0	1.0	0.0	1.0
tour_arrangement_Package Tour	0.0	0.0	0.0	1.0	0.0
package_transport_int_No	1.0	1.0	1.0	1.0	1.0
package_transport_int_Yes	0.0	0.0	0.0	0.0	0.0
package_accommodation_No	1.0	1.0	1.0	0.0	1.0
package_accommodation_Yes	0.0	0.0	0.0	1.0	0.0
package_food_No	1.0	1.0	1.0	0.0	1.0
package_food_Yes	0.0	0.0	0.0	1.0	0.0
package_transport_tz_No	1.0	1.0	1.0	0.0	1.0
package_transport_tz_Yes	0.0	0.0	0.0	1.0	0.0
package_sightseeing_No	1.0	1.0	1.0	0.0	1.0
package_sightseeing_Yes	0.0	0.0	0.0	1.0	0.0
package_guided_tour_No	1.0	1.0	1.0	0.0	1.0
package_guided_tour_Yes	0.0	0.0	0.0	1.0	0.0
package_insurance_No	1.0	1.0	1.0	1.0	1.0
package_insurance_Yes	0.0	0.0	0.0	0.0	0.0
first_trip_tz_No	1.0	0.0	1.0	0.0	0.0
first_trip_tz_Yes	0.0	1.0	0.0	1.0	1.0

```


In [20]: # checking all_dataset to know if ready for ML
#all_data.info()

In [21]: train =all_data.iloc[0:ntrain]
test =all_data.iloc[ntrain:]

In [22]: #checking shape accuracy with respect to original shape
train.shape,test.shape

Out[22]: ((4809, 31), (1601, 31))

In [23]: train.to_csv("trainuni1ag2.csv",index=False)
test.to_csv("testunil1ag2.csv",index=False)

In [24]: y= train["total_cost"]
testsz= test.drop(columns="total_cost",inplace =True)
train.drop(columns="total_cost",inplace =True)

C:\Users\ATN-PC\AppData\Local\Temp\ipykernel_8252\2855698061.py:2: SettingWithCopyWarning:
A value is trying to be set on a copy of a DataFrame

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy
testsz= test.drop(columns="total_cost",inplace =True)
C:\Users\ATN-PC\AppData\Local\Temp\ipykernel_8252\2855698061.py:3: SettingWithCopyWarning:
A value is trying to be set on a copy of a slice from a DataFrame

See the caveats in the documentation: https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy
train.drop(columns="total_cost",inplace =True)

In [25]: #test.info()

In [26]: #train.info()

In [27]: from sklearn.model_selection import train_test_split
X_train,X_test,y_train,y_test= train_test_split(train,y,test_size =0.1,random_state =42,shuffle =True)
from sklearn.preprocessing import StandardScaler
sc = StandardScaler()
sc.fit(X_train)
X_train_std = sc.transform(X_train)
X_test_std = sc.transform(X_test)

In [28]: from sklearn.linear_model import LinearRegression

model = LinearRegression()

gbm = model.fit(X_train, y_train)

from sklearn.metrics import mean_squared_error

In [ ]:

In [ ]:

In [29]: from sklearn.metrics import mean_absolute_error

# After fitting linear regression model:
y_pred = model.predict(X_test) # Predicted values on test data
mae = mean_absolute_error(y_test, y_pred) # Calculate MAE

print("Mean Absolute Error:", mae)

Mean Absolute Error: 5698891.961778337

In [ ]:

In [30]: from sklearn.tree import DecisionTreeRegressor,ExtraTreeRegressor
from sklearn.ensemble import RandomForestRegressor,ExtraTreesRegressor,GradientBoostingRegressor
from sklearn.metrics import mean_absolute_error
def mae(y_test,y_pred):
    mae = mean_absolute_error(y_test, y_pred)
    return mae

In [31]: dt = DecisionTreeRegressor(max_depth =1, random_state =1)
dt.fit(X_train_std,y_train)

pred =dt.predict(X_test)
print(mae(y_test,pred))

6261231.509697027
C:\Users\ATN-PC\anaconda3\Lib\site-packages\sklearn\base.py:457: UserWarning: X has feature names, but DecisionTreeRegressor was fitted without feature names
warnings.warn(

In [32]: test.columns

Out[32]: Index(['country', 'age_group', 'travel_with', 'total_female', 'total_male',
'purpose', 'main_activity', 'info_source', 'night_mainland',
'night_zanzibar', 'payment_mode', 'most_impressing',
'tour_arrangement_Independent', 'tour_arrangement_Package Tour',
'package_transport_int_No', 'package_transport_int_Yes',
'package_accommodation_No', 'package_accommodation_Yes',
'package_food_No', 'package_food_Yes', 'package_transport_tz_No',
'package_transport_tz_Yes', 'package_sightseeing_No',
'package_sightseeing_Yes', 'package_guided_tour_No',
'package_guided_tour_Yes', 'package_insurance_No',
'package_insurance_Yes', 'first_trip_tz_No', 'first_trip_tz_Yes'],
      dtype='object')

After testing the model performance with x_train and y_train and checking it's MAE with y_test and y_pred(mae(y_test,pred)). Then i now train the catboost_train_data with whole training and y data
Pool(train,y)

In [33]: from catboost import CatBoostRegressor,Pool
catboost_train_data = Pool(train,y)

In [34]: cb_model =CatBoostRegressor(
    loss_function="MAE",
    eval_metric="MAE",
    iterations=1000,
    learning_rate=0.07,
    random_seed=42,
    depth=4,
    bagging_temperature=0.7,

    early_stopping_rounds=10,
    verbose=False, # Disable training verbosity for a cleaner output
    reg_lambda=0.01
)
cb_model.fit(catboost_train_data,plot=True)

MetricVisualizer(layout=Layout(aligned_self='stretch', height='500px'))
<catboost.core.CatBoostRegressor at 0x114138f1f0>

In [35]: ## y_pred =cb_model.predict(X_test)
#print(mae(y_test,pred))

In [36]: test_id = pd.read_csv(r"C:\Users\ATN-PC\Desktop\kvc zeph\it_spend_files\Test.csv")

In [37]: test_id.ID

Out[37]: 0      tour_1
1      tour_100
2      tour_1001
3      tour_1006
4      tour_1009
...
1596      tour_988
1597      tour_990
1598      tour_992
1599      tour_996
1600      tour_998
Name: ID, Length: 1601, dtype: object

In [38]: y_pred = cb_model.predict(X_test)
print(mae(y_test,y_pred))

4997105.893549782

In [39]: predcdf =cb_model.predict(test)

This result was uploaded to Zindi competition website on tanzania tourism prediction and my best result was obtained using CatBoostRegressor and i obtained 11th position in the competition

In [40]: output =pd.DataFrame({"test_id":test_id.ID ,"total_cost": predcdf})
output.to_csv("Zinicatfff.csv",index =False)
output

Out[40]:
```

test_id	total_cost
0	tour_1 2.134284e+07
1	tour_100 7.098960e+06
2	tour_1001 9.302364e+06
3	tour_1006 2.678122e+06
4	tour_1009 1.726829e+07
...	...
1596	tour_988 3.668368e+05
1597	tour_990 1.922577e+07
1598	tour_992 1.036344e+06
1599	tour_996 3.175156e+05
1600	tour_998 4.454162e+06

1601 rows × 2 columns

```
In [ ]:

In [41]: import lightgbm as lgb

lgb_train = lgb.Dataset(train, label=y)

In [42]: model1 = lgb.LGBMRegressor(objective="regression",iterations=1000,
learning_rate=0.05,
random_seed=42,

metric="mae") # optional early stopping

model1.fit(train,y)

[LightGBM] [Warning] Unknown parameter: iterations
[LightGBM] [Warning] Found whitespace in feature_names, replace with underlines
[LightGBM] [Warning] Unknown parameter: iterations
[LightGBM] [Info] Auto-choosing row-wise multi-threading, the overhead of testing was 0.001788 seconds.
You can set 'force_row_wisetrue' to remove the overhead.
And if memory is not enough, you can set 'force_col_wisetrue'.
[LightGBM] [Info] Total Bins 258
[LightGBM] [Info] Number of data points in the train set: 4809, number of used features: 30
[LightGBM] [Info] Start training from score 8114388.775496

Out[42]: +-----+
|          LGBMRegressor          |
|LGBMRegressor(iterations=1000, learning_rate=0.05, metric='mae', |
|          objective='regression', random_seed=42)                  |
+-----+

In [43]: y_pred = model1.predict(X_test)
print(mae(y_test,y_pred))

[LightGBM] [Warning] Unknown parameter: iterations
4189565.1779530183

In [44]: pred1 =model1.predict(test)

[LightGBM] [Warning] Unknown parameter: iterations

This result was uploaded to Zindi competition website on tanzania tourism prediction

In [45]: output =pd.DataFrame({"test_id":test_id.ID ,"total_cost": pred1})
output.to_csv("Zinilgbm3.csv",index =False)
output

Out[45]:
```

test_id	total_cost
0	tour_1 2.319674e+07
1	tour_100 1.011334e+07
2	tour_1001 1.559144e+07
3	tour_1006 4.958598e+06
4	tour_1009 2.566547e+07
...	...
1596	tour_988 7.371609e+05
1597	tour_990 2.457233e+07
1598	tour_992 1.355749e+06
1599	tour_996 6.764005e+05
1600	tour_998 9.673842e+06

1601 rows × 2 columns

```
In [ ]:

In [ ]:

In [ ]:
```