

# Identifying the successful drivers

Mansoureh Aghabeig





# Main steps

- Understanding data
- Preparing data for applying machine learning model
- Analysis data
- Applying machine learning model
- Steps for implementing the machine learning model
- Choosing machine learning model
- Suggestion for future work



# Understanding data



# Understanding data

Checking data from different aspect:

- Determining different attributes (columns) in data
- Type of data to be sure it is consistent with model we are going to apply
- Missing values in data

# Statistical information of data

```
|: import pandas as pd
import matplotlib.pyplot as plt

df = pd.read_csv ('TV_CASE.csv', sep=';')
print(df.info())
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 5347 entries, 0 to 5346
Data columns (total 9 columns):
#   Column                Non-Null Count  Dtype
---  -
0   Programme             5347 non-null   int64
1   Station               5347 non-null   int64
2   datetime              5347 non-null   object
3   cost                  5347 non-null   float64
4   profit                5332 non-null   float64
5   order_sum             5332 non-null   float64
6   male in %             5347 non-null   float64
7   new customer in %     5347 non-null   float64
8   vis                   5347 non-null   int64
dtypes: float64(5), int64(3), object(1)
memory usage: 376.1+ KB
None
```

We can see :

- **9** attributes in data
- Total number of data is: **5347** rows
- **Two** attributes contain missing values
- We have **one non-numerical** type in data

# Checking missing values

```
|: df.isnull().any()
```

```
|: Programme      False
   Station        False
   datetime       False
   cost           False
   profit         True
   order_sum      True
   male in %      False
   new customer in % False
   vis            False
   dtype: bool
```

```
|: df.isnull().sum()
```

```
|: Programme      0
   Station        0
   datetime       0
   cost           0
   profit         15
   order_sum      15
   male in %      0
   new customer in % 0
   vis            0
   dtype: int64
```

- Checking better data to understand where is the missing values
- The col **profit** and **order\_sum** contain missing values.
- The total number of missing values are **15 rows**



# **Preparing data for applying machine learning model**





# Handling missing values

It is **impossible** to apply a machine learning method on dataset contains missing values.

There are two main solutions for handling missing values:

- Deleting rows contain missing values
- Imputing (replacing) missing values with a value

As the number of rows contain missing values in our data is so low in compared to whole data set, **exactly 0.2%**, we decide to simply remove rows containing missing values





**Analysis data**



# Success metric

We should define **which attribute** in data is the target one. In other word we should decide if one driver is a successful based on this attribute.

In our date, there are four attributes which may show the trend of success in one driver:

- **Profit:** Predicted profit provided by the data science team
- **Order\_sum:** Predicted orders provided by the data science team
- **New customer in:** Predicted new customers in % who saw the spot provided by the data science team
- **Vis:** Predicted visits on the Mister Spex website driven by the spot (provided by the data science team)



# What is the main success metric?

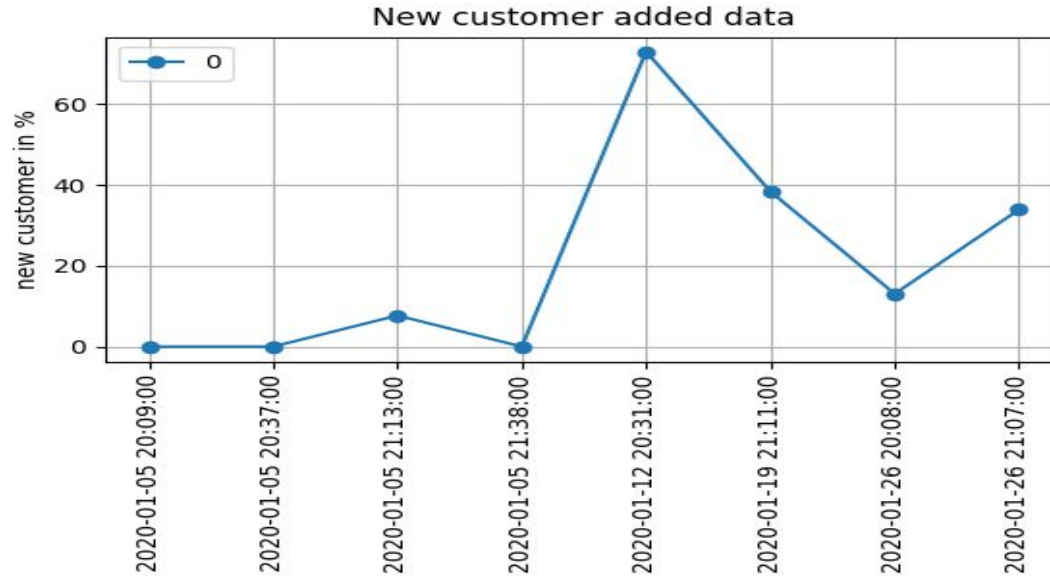
- We know that if each of these four metrics have an increasing trend during time, can expect the successfulness in a driver.
- But we need to choose one or two of these metrics for building our model.
- But which of them?



# Solution for finding the best metric

- Hypothesis 1 :
  - One solution is **plotting time series data** for each driver based on these four metrics.
  - We did it for example for “Programme 0”
  - We can see time series plots with different shapes for each metrics
  - Also we can see different pick during time for each metric
- Interpretation:
  - We cannot simply say that which metric is the best one by seeing their time\_series data
  - Because for example, for “programme 0”, in first days of month, the order sum is quite high while the profit value is low. Or you may see in the mid of the month, the new customer added is increased significantly while the profit stay static.

# Time series plot for Program 0- New customer

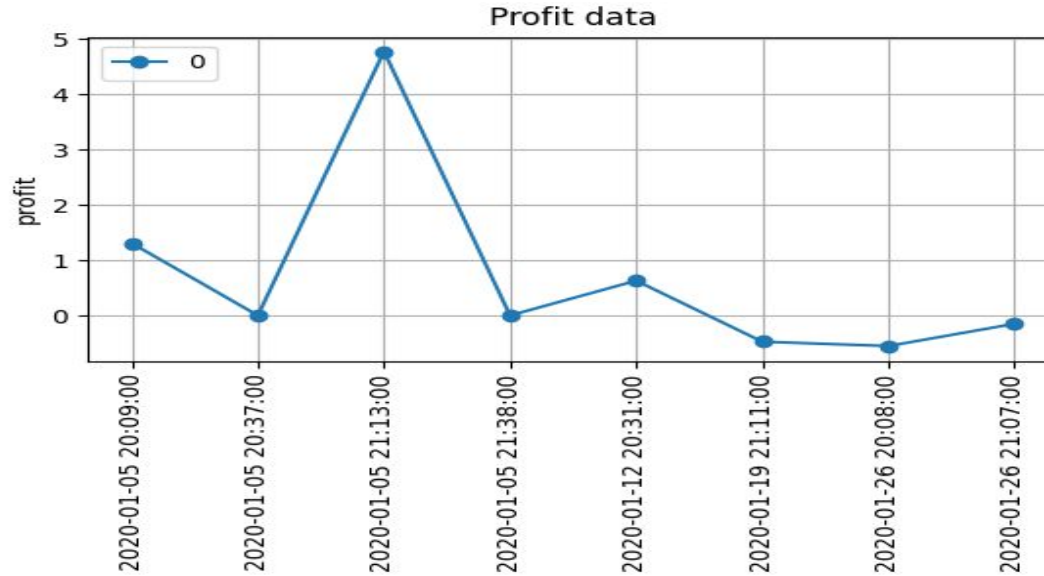




## Time series plot for Program 0- Order\_sum



# Time series plot for Program 0- Profit





# Time series plots for Program 0-Visit







# Solution for finding the best metric

- Hypothesis 2:
  - We should find what is **the relationship between different metric**
  - For this, we need to find the **correlation values between metrics**
  - Please note, a correlation value above zero shows a positive relation between two metrics, in other word by increasing the one, the other increase also and a correlation below zero shows an inverse relation between two metrics



## Correlation values between metrics

Metric	Profit	Visit	Order_sum
Profit	1	-----	-----
Visit	0.051	1	-----
Order_sum	0.23	0.3	1
New_customer	-0.02	-0.02	-0.007



# Solution for finding the best metric

Interpretation for Hypothesis 2:

- We can see there is a **positive relation among all metrics except new\_customer**.
  - This means as profit, order\_sum and number of visit are increased we have decrease in new customer
  - But comparing the correlation value of new customer with other metrics shows very minor relationship



# Solution for finding the best metric

Final decision about successful metric:

- Among three metrics: “**Profit, Sum\_order & Visit**” as they have all the positive correlation with each other, we can consider each of them as success metric
- In this study, we consider **profit** as it has a more meaning for customer.
- About metric new customer as both it has negative and also very low relation with other metrics we can ignore this metric as a successful one



# Applying machine learning model





# Choosing a data science model

## Goal:

- In this project, we want to identify each driver (program) is successful or not.
- There are several approaches for reaching this goal.

## Our approach:

- Predicting profit values for the other month for each driver

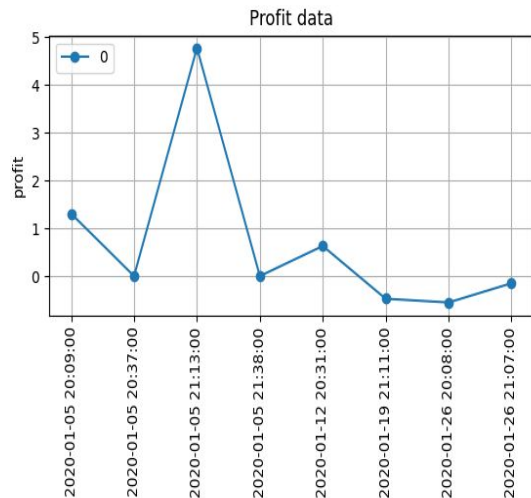


## Our approach

If we can find the **pattern of time series of profit for this month**, maybe we can predict the next time series of profit for the other month too.

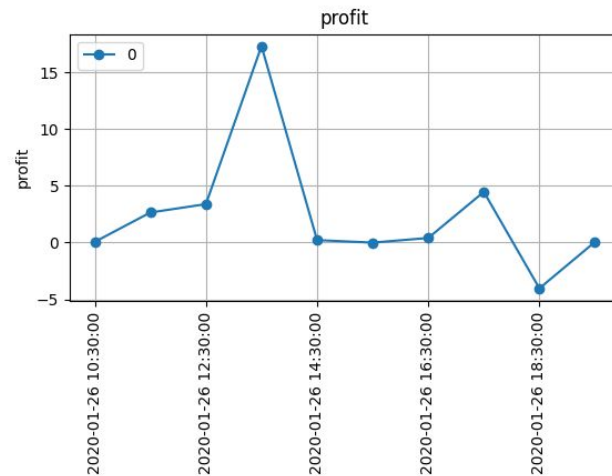
In this way, we can **compare the total profit of each driver** from one month to other month, and if we see an increase in total profit for the other month, we may conclude that this driver has a successful trend and we consider it as a good one.

# Our approach



Machine learning Model

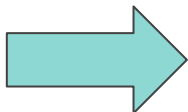
Predict for other month



Total Profit



$N$



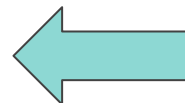
Compute  $M - N$

- If it is positive = Success
- If it is negative = Unsuccess

Total Profit



$M$





# Steps for implementing the machine learning model





# Generate a test data set

We need to generate a test data set which has the same statistical properties with the original dataset

In this way, we are able to compare profit result of this month and the other month

# Generate a test data set

df

	Programme	Station	datetime	cost	profit	order_sum	male
0	1226	25	2020-01-02 08:31:00	8.00000	-0.57495	0.02318	11.74
1	722	26	2020-01-02 08:47:00	8.00000	0.54784	0.00744	5.908
2	32	46	2020-01-02 09:00:00	5.00000	0.00000	0.00000	8.173
3	1296	31	2020-01-02 09:01:00	5.00000	-0.25717	0.00128	21.48
4	836	0	2020-01-02 09:02:00	21.00000	-12.35411	0.06575	30.34
5	1423	23	2020-01-02 09:02:00	17.00000	-10.01115	0.05328	27.44
6	891	4	2020-01-02 09:08:00	27.00000	0.46805	0.01001	4.967
7	890	11	2020-01-02 09:09:00	9.00000	0.24622	0.00527	92.11
8	18	34	2020-01-02 09:10:00	5.00000	0.35770	0.00955	87.27
9	499	18	2020-01-02 09:13:00	9.00000	-0.50961	0.00333	58.01
10	81	26	2020-01-02 09:13:00	8.00000	-0.43094	0.00282	85.85
11	1462	30	2020-01-02 09:25:00	11.00000	1.12289	0.01152	27.59
12	61	25	2020-01-02 09:27:00	8.00000	0.13976	0.00557	6.979
13	639	35	2020-01-02 09:27:00	10.00000	0.15962	0.00636	14.15
14	795	38	2020-01-02 09:29:00	12.00000	3.27619	0.00614	95.64
15	890	11	2020-01-02 09:31:00	9.00000	0.00000	0.00000	46.15

df Format: %s

Original dataset

Keeping average and  
standard deviation for  
each column



Generate Test dataset

test\_df

	Programme	Station	datetime	cost	profit	order_sum	male
0	1226	45	2020-01-02 08:31:00	120	-1.93848	0.01418	55.74
1	722	14	2020-01-02 08:47:00	194	7.49053	0.06572	61.62
2	32	18	2020-01-02 09:00:00	257	-7.59288	0.01968	51.52
3	1296	9	2020-01-02 09:01:00	415	6.29050	0.04780	97.29
4	836	24	2020-01-02 09:02:00	636	9.65812	0.06303	63.06
5	1423	8	2020-01-02 09:02:00	908	-8.32796	0.03482	64.46
6	891	10	2020-01-02 09:08:00	165	1.24195	0.02640	28.45
7	890	7	2020-01-02 09:09:00	294	-10.75535	0.06373	1.932
8	18	1	2020-01-02 09:10:00	1330	-3.31202	0.00788	30.63
9	499	0	2020-01-02 09:13:00	355	5.61490	0.00307	43.84
10	81	23	2020-01-02 09:13:00	675	4.25814	0.01827	18.93
11	1462	24	2020-01-02 09:25:00	913	-6.55704	0.03041	56.01
12	61	28	2020-01-02 09:27:00	506	-17.42243	0.01423	39.32
13	639	44	2020-01-02 09:27:00	848	2.71190	0.05073	49.41
14	795	9	2020-01-02 09:29:00	277	16.88276	0.08257	39.01
15	890	40	2020-01-02 09:31:00	75	4.61911	0.00391	50.77

test\_df Format: %s

Test dataset



# Feature selection

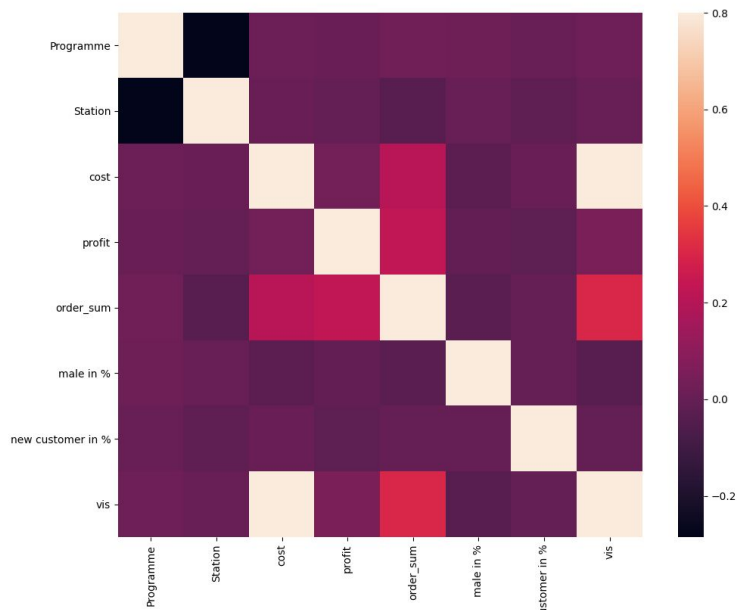
We need to identify and remove unneeded, irrelevant and redundant attributes from data that do not contribute to the accuracy of a predictive model or may in fact decrease the accuracy of the model.

We use two methods for feature (attribute) selection:

- Correlation matrix generation for Profit attribute
- Determining feature importance based on the used machine learning model



# Confusion matrix for Profit

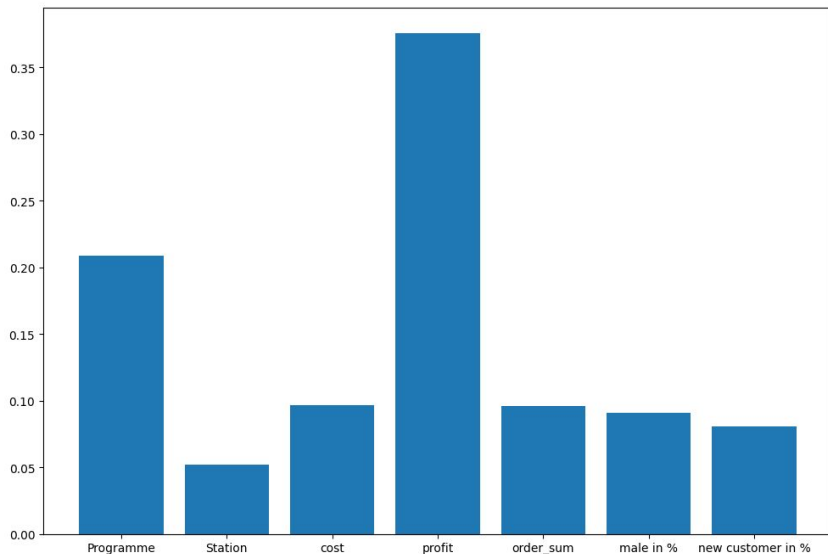


We can observe that:

- There is a high relation between profit and order\_sum, cost & vis
- In other hand, there is a low relation between profit and station & new customer



# Feature importance based on model



We determine the importance of each attribute based on the **machine learning model** we are using:

We can observe that for our model, **program, profit, cost, order\_sum and male** have important influence

While, **station and new customer** do not have a significant influence on our model



# Feature selection

Based on confusion matrix for Profile and also feature importance for our model, we decide to not consider these two attribute in our model

- Station
- New customer



# Choosing machine learning model







# Random forest regressor model

- It is not easy to choose an exact predictive model for the project
- We try to implement the “Random forest regressor” which applied successfully in the project “Sales prediction over time”. As there are some common points between two project



# Steps for implementing model

- Removing the Station and New customer attribute from our data set
- Dividing our whole dataset to two train and test part
- Training model based on train part of data
- We may use test part for determining the accuracy of our model (*in this version, as we just applied one model, we did not evaluate the accuracy of it, but in next version, we need to check accuracy for different models, to find the best one*)
- Applying the trained model to test data set.
- Generate the result data set

# The result file



df\_result

	Programme	Profit_first_month	Profit_second_month	successful
300	0	0.53294	7.08388	0.00000
845	1	-12.45779	5.23311	0.00000
1058	2	0.00000	0.85595	0.00000
449	3	2.94444	4.39006	0.00000
293	4	0.00621	5.61554	0.00000
960	5	0.20064	0.26175	0.00000
665	6	26.21667	5.18383	1.00000
555	7	65.43882	72.02362	0.00000
1281	8	-0.10729	0.30714	0.00000
1244	9	0.00000	3.31906	0.00000
1422	10	1.93976	0.42812	1.00000
180	11	0.00000	3.42333	0.00000
501	12	1.86760	2.88102	0.00000
1290	13	3.91910	6.05723	0.00000
866	14	-16.92593	-1.78538	0.00000
89	15	-6.36162	0.25532	0.00000
995	16	0.30880	0.11412	1.00000

df\_result

Format: %s



# **Suggestion for future work**





# Suggestion

In this project, we just apply “Random forest regressor” for train our model, but in future we may apply more models used for forecasting time series data such as

- Linear regression
- K nearest neighbor
- Stochastic gradient descent
- Decision tree
- Neural network
- Autoregressive–moving-average model

It is important to know, we have a time series forecasting model here.