

T 201: Experimentation and Uplift Testing

This notebook serves as the solution of the second task **Experimentation and Uplift Testing of Quantum Data Analytics Virtual Experience Program**.

In the previous task we conducted analysis on the client's transaction dataset and identified customer purchasing behaviours to generate insights and provide customer recommendations.

We now extend our analysis from Task 1 to help us identify benchmark stores that allow you to test the impact of the trial store layouts on customer sales.

Category Manager for Chips, has asked us to test the impact of the new trial layouts with a data driven recommendation to whether or not the trial layout should be rolled out to all their stores.

The new trial layouts were rolled in three stores, store number **77, 86, and 88** during trial period. The trial period was for three months, i.e. **February, March, and April 2019**.

We will complete this task in two phases.

- **Phase 1:** We will find control store corresponding to the trial stores
- **Phase 2:** Assessment of trial stores during the trial period

What is Control Store?

Store that have similar performance to trial store in terms of measure like monthly sales, number of customers, etc.

Sequence of Analysis

1. We will be defining *monthly sales* and *number of customers* to find the potential control store.
1. Then we will define driving metrics to rank the stores on how similar they are to the trial store.
 - Correlation Score
 - Absolute Distance Score
1. We will select the store having highest score as control store for the corresponding trial store.
 - Only pre-trial period will be considered to select the control store
 - We will also visualise the performance to cross check that we have selected right control store
1. Once we get control store for each trial store, we will assess the trial store performance during the trial period.
 - We will use statistical tests and visualizations to test the impact of the new trial layouts for all these stores
1. Finally, we will make conclusion on whether or not the trial layout should be rolled out to all their stores.

Importing Libraries

```
In [1]: import numpy as np
import pandas as pd

import matplotlib.pyplot as plt
import seaborn as sns
from scipy import stats
sns.set(font_scale=1)

import warnings
warnings.filterwarnings("ignore")
```

Loading Dataset

```
In [2]: #loading dataset
df = pd.read_csv("C:/Users/amit/Jupyter Notebooks/Quantum/Task 2/QVI_data.csv")
```

```
In [3]: #checking first 5 rows
df.head()
```

	LYLTY_CARD_NBR	DATE	STORE_NBR	TXN_ID	PROD_NBR	PROD_NAME	PROD_QTY	TOT_SALES	PACK_SIZE	BRAND
0	1000	10-10-17	1	1	5	Natural Chip Corny SesSalt175g	2	6.0	175	NATURAL
1	1002	10-08-16	1	2	58	Red Rock Deli Chkn&Garlic Aioli 150g	1	2.7	150	RDR
2	1003	2019-03-07	1	3	52	Grain Waves Sour Cream&Chives 210G	1	3.6	210	GRNWVES
3	1003	2019-03-08	1	4	106	Natural Chip&Co Henry Soy Chkn175g	1	3.0	175	NATURAL
4	1004	2018-11-02	1	5	96	WW Original Stacked Chips 160g	1	1.9	160	WOOLWORTHS

Checking data type of variables

```
In [4]: #checking data types
df.dtypes
```

```
Out[4]: LYLTY_CARD_NBR    int64
DATE                  object
STORE_NBR             int64
TXN_ID               int64
PROD_NBR             int64
PROD_NAME            object
PROD_QTY             int64
TOT_SALES            int64
PACK_SIZE            int64
BRAND                object
LIFESTAGE            object
PREMIUM_CUSTOMER     object
dtype: object
```

Creating Year-Month column

Since we are interested in analyzing numbers on monthly basis, therefore, we are extracting year and month from date column and storing it in `yyyymm` format.

```
In [5]: #creating new column MonthID in yyyymm format from the DATE column
df["MonthID"] = df["DATE"].str.slice(0,7)
df["MonthID"] = df["MonthID"].str.replace(pat="-", repl="")
df.head()
```

	LYLTY_CARD_NBR	DATE	STORE_NBR	TXN_ID	PROD_NBR	PROD_NAME	PROD_QTY	TOT_SALES	PACK_SIZE	BRAND
0	1000	10-10-17	1	1	5	Natural Chip Corny SesSalt175g	2	6.0	175	NATURAL
1	1002	10-08-16	1	2	58	Red Rock Deli Chkn&Garlic Aioli 150g	1	2.7	150	RDR
2	1003	2019-03-07	1	3	52	Grain Waves Sour Cream&Chives 210G	1	3.6	210	GRNWVES
3	1003	2019-03-08	1	4	106	Natural Chip&Co Henry Soy Chkn175g	1	3.0	175	NATURAL
4	1004	2018-11-02	1	5	96	WW Original Stacked Chips 160g	1	1.9	160	WOOLWORTHS

```
In [6]: #group data by STR_NBR and MonthID, and taking sum of TOT_SALES and unique count of LYLTY_CARD_NBR.
measureOverTime = df.groupby(["STORE_NBR", "MonthID"]).agg({"TOT_SALES": "sum", "LYLTY_CARD_NBR": "pd.Series.nunique()}).reset_index()
```

```
#renaming columns to reflect the true meaning of the features
measureOverTime.rename(columns={"TOT_SALES": "totSales", "LYLTY_CARD_NBR": "nCustomers"}, inplace=True)
```

```
Out[6]:
```

	STORE_NBR	MonthID	totSales	nCustomers
0	1	201807	206.9	49
1	1	201808	176.1	42
2	1	201809	278.8	59
3	1	201810	168.1	44
4	1	201811	192.6	46
...
3164	272	201902	395.5	45
3165	272	201903	442.3	50
3166	272	201904	445.1	54
3167	272	201905	314.6	34
3168	272	201906	312.1	34

3169 rows x 4 columns

Stores with full observations

We are interested in stores, who are operated or had transactions for all 12 months. We will consider only these stores for analysis and leave other stores.

For this we simply count MonthID for each store and consider those having equal to 12.

```
In [7]: #having store numbers having count of MonthIDs = 12
storesWithFullObs = df.groupby("STORE_NBR")["MonthID"].count().reset_index()
storesWithFullObs = storesWithFullObs[storesWithFullObs["MonthID"] == 12][["STORE_NBR"]].unique()
```

```
Out[7]: array([ 1,  2,  3,  4,  5,  6,  7,  8,  9, 10, 12, 13, 14,
        15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27,
        28, 29, 30, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41,
        42, 43, 45, 46, 47, 48, 49, 50, 51, 52, 53, 54, 55,
        56, 57, 58, 59, 60, 61, 62, 63, 64, 65, 66, 67, 68,
        69, 70, 71, 72, 73, 74, 75, 77, 78, 79, 80, 81, 82,
        83, 84, 86, 87, 88, 89, 90, 91, 93, 94, 95, 96, 97,
        98, 99, 100, 101, 102, 103, 104, 105, 106, 107, 108, 109, 110,
        111, 112, 113, 114, 115, 116, 118, 119, 120, 121, 122, 123, 124,
        125, 126, 127, 128, 129, 130, 131, 132, 133, 134, 135, 136, 137,
        138, 139, 140, 141, 142, 143, 144, 145, 146, 147, 148, 149, 150,
        151, 152, 153, 154, 155, 156, 157, 158, 159, 160, 161, 162, 163,
        164, 165, 166, 167, 168, 169, 170, 171, 172, 173, 174, 175, 176,
        177, 178, 179, 180, 181, 182, 183, 184, 185, 186, 187, 188, 189,
        190, 191, 192, 194, 195, 196, 197, 198, 199, 200, 201, 202, 203,
        204, 205, 207, 208, 209, 210, 212, 213, 214, 215, 216, 217, 219,
        220, 221, 222, 223, 224, 225, 226, 227, 228, 229, 230, 231, 232,
        233, 234, 235, 236, 237, 238, 239, 240, 241, 242, 243, 244, 245,
        246, 247, 248, 249, 250, 251, 253, 254, 255, 256, 257, 259, 259,
        260, 261, 262, 263, 264, 265, 266, 267, 268, 269, 270, 271, 272],
      dtype=int64)
```

Changing data type of MonthID from object to integer for ease

```
In [8]: measureOverTime.dtypes
```

```
Out[8]: STORE_NBR    int64
MonthID          object
totSales         float64
nCustomers        int64
dtype: object
```

```
In [9]: measureOverTime["MonthID"] = measureOverTime["MonthID"].astype(int)
```

```
Out[9]:
```

	STORE_NBR	MonthID	totSales	nCustomers
0	1	201807	206.9	49
1	1	201808	176.1	42
2	1	201809	278.8	59
3	1	201810	168.1	44
4	1	201811	192.6	46
...
3164	272	201902	395.5	45
3165	272	201903	442.3	50
3166	272	201904	445.1	54
3167	272	201905	314.6	34
3168	272	201906	312.1	34

3169 rows x 4 columns

Stores with full observations

We are interested in stores, who are operated or had transactions for all 12 months. We will consider only these stores for analysis and leave other stores.

For this we simply count MonthID for each store and consider those having equal to 12.

```
In [7]: #having store numbers having count of MonthIDs = 12
storesWithFullObs = df.groupby("STORE_NBR")["MonthID"].count().reset_index()
storesWithFullObs = storesWithFullObs[storesWithFullObs["MonthID"] == 12][["STORE_NBR"]].unique()
```

```
Out[7]: array([ 1,  2,  3,  4,  5,  6,  7,  8,  9, 10, 12, 13, 14,
        15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27,
        28, 29, 30, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41,
        42, 43, 45, 46, 47, 48, 49, 50, 51, 52, 53, 54, 55,
        56, 57, 58, 59, 60, 61, 62, 63, 64, 65, 66, 67, 68,
        69, 70, 71, 72, 73, 74, 75, 77, 78, 79, 80, 81, 82,
        83, 84, 86, 87, 88, 89, 90, 91, 93, 94, 95, 96, 97,
        98, 99, 100, 101, 102, 103, 104, 105, 106, 107, 108, 109, 110,
        111, 112, 113, 114, 115, 116, 118, 119, 120, 121, 122, 123, 124,
        125, 126, 127, 128, 129, 130, 131, 132, 133, 134, 135, 136, 137,
        138, 139, 140, 141, 142, 143, 144, 145, 146, 147, 148, 149, 150,
        151, 152, 153, 154, 155, 156, 157, 158, 159, 160, 161, 162, 163,
        164, 165, 166, 167, 168, 169, 170, 171, 172, 173, 174, 175, 176,
        177, 178, 179, 180, 181, 182, 183, 184, 185, 186, 187, 188, 189,
        190, 191, 192, 194, 195, 196, 197, 198, 199, 200, 201, 202, 203,
        204, 205, 207, 208, 209, 210, 212, 213, 214, 215, 216, 217, 219,
        220, 221, 222, 223, 224, 225, 226, 227, 228, 229, 230, 231, 232,
        233, 234, 235, 236, 237, 238, 239, 240, 241, 242, 243, 244, 245,
        246, 247, 248, 249, 250, 251, 253, 254, 255, 256, 257, 259, 259,
        260, 261, 262, 263, 264, 265, 266, 267, 268, 269, 270, 271, 272],
      dtype=int64)
```

```
In [8]: measureOverTime.dtypes
```

```
Out[8]: STORE_NBR    int64
MonthID          object
totSales         float64
nCustomers        int64
dtype: object
```

```
In [9]: measureOverTime["MonthID"] = measureOverTime["MonthID"].astype(int)
```

```
Out[9]:
```

	STORE_NBR	MonthID	totSales	nCustomers
0	1	201807	206.9	49
1	1	201808	176.1	42
2	1	201809	278.8	59
3	1	201810	168.1	44
4	1	201811	192.6	46
...
3164	272	201902	395.5	45
3165	272	201903	442.3	50
3166	272	201904	445.1	54
3167	272	201905	314.6	34
3168	272	201906	312.1	34

3169 rows x 4 columns

Defining Driving Metrics

We will be using two main metrics for finding the similarity between control stores and trial stores.

- **Correlation Score**
 - We will calculate how correlated the performance of each store is to the trial store
 - Correlation value lies between -1 to 1. Value near to 1 means similarity value near to 1 means dis-similarity
- **Absolute Difference Score**
 - We will calculate the absolute difference between the trial store's performance and each control store's performance
 - We will standardize the distance so that the measure ranges from 0 to 1
 - We will take (1 - standardized distance) so that it has similar effect as correlation score. Value near to 1 means similarity value near to 0 means dis-similarity

We will be using these metrics to calculate score for both `totSales` and `nCustomers` i.e. we will be having total 4 scores for each control store to the trial store. Then we will combine all these scores to get score for each store and we sort them in descending order. **Store with highest score other than trial store itself becomes the control store of the respective trial store.**

Let us define the functions to calculate all 4 scores and to store them into a single dataframe. It will reduce our time of re-write the code again.

```
In [11]: def corrSales(data, trialStore):
    """
    It calculates the correlation score of each store to the trial store based on total sales.

    Parameters:
        data: (dataframe) It takes dataframe containing store number, months, sales and number o
        f customers
        trialStore: (integer) It takes the trial store number from which we need to compare all
        stores

    Return: It returns a dataframe with store number and correlation score corresponding to that store
    and the trial store
    """
    df = data[["STORE_NBR", "MonthID", "totSales"]]
    storeNBR = df["STORE_NBR"].unique()
    corr = list()

    ts = df[df["STORE_NBR"] == trialStore].iloc[:, -1]
    for n in storeNBR:
        cs = df[df["STORE_NBR"] == n].iloc[:, -1]
        corr.append(np.corrcoef(ts, cs)[0, 1])

    corrSalesDf = pd.DataFrame({"STORE_NBR": storeNBR, "corrSales": corr})

    return corrSalesDf
```

```
In [12]: def corrCust(data, trialStore):
    """
    It calculates the correlation score of each store to the trial store based on number of customer
    s.

    Parameters:
        data: (dataframe) It takes dataframe containing store number, months, sales and number o
        f customers
        trialStore: (integer) It takes the trial store number from which we need to compare all
        stores

    Return: It returns a dataframe with store number and correlation score corresponding to that store
    and the trial store
    """
    df = data[["STORE_NBR", "MonthID", "nCustomers"]]
    storeNBR = df["STORE_NBR"].unique()
    corr = list()

    ts = df[df["STORE_NBR"] == trialStore].iloc[:, -1]
    for n in storeNBR:
        cs = df[df["STORE_NBR"] == n].iloc[:, -1]
        corr.append(np.corrcoef(ts, cs)[0, 1])

    corrCustDf = pd.DataFrame({"STORE_NBR": storeNBR, "corrCust": corr})

    return corrCustDf
```

```
In [13]: def magDistSales(data, trialStore):
    """
    It calculates the absolute distance of each store to the trial store based on monthly sales.

    Parameters:
        data: (dataframe) It takes dataframe containing store number, months, sales and number o
        f customers
        trialStore: (integer) It takes the trial store number from which we need to compare all
        stores

    Return: It returns a dataframe with store number and absolute distance corresponding to that store
    and the trial store
    """
    df = data[["STORE_NBR", "MonthID", "totSales"]]
    storeNBR = df["STORE_NBR"].unique()
    dist = list()

    ts = np.array(df[df["STORE_NBR"] == trialStore].iloc[:, -1])
    for n in storeNBR:
        cs = np.array(df[df["STORE_NBR"] == n].iloc[:, -1])
        d = sum(np.abs(ts - cs))
        dist.append(d)

    magDistSalesDf = pd.DataFrame({"STORE_NBR": storeNBR, "distSales": dist})

    return magDistSalesDf
```

```
In [14]: def magDistCust(data, trialStore):
    """
    It calculates the absolute distance of each store to the trial store based on number of customer
    s.

    Parameters:
        data: (dataframe) It takes dataframe containing store number, months, sales and number o
        f customers
        trialStore: (integer) It takes the trial store number from which we need to compare all
        stores

    Return: It returns a dataframe with store number and absolute distance corresponding to that store
    and the trial store
    """
    df = data[["STORE_NBR", "MonthID", "nCustomers"]]
    storeNBR = df["STORE_NBR"].unique()
    dist = list()

    ts = np.array(df[df["STORE_NBR"] == trialStore].iloc[:, -1])
    for n in storeNBR:
        cs = np.array(df[df["STORE_NBR"] == n].iloc[:, -1])
        d = sum(np.abs(ts - cs))
        dist.append(d)

    magDistCustDf = pd.DataFrame({"STORE_NBR": storeNBR, "distCust": dist})

    return magDistCustDf
```

```
In [15]: def stdMagDis(data, colList):
    """
    It standardizes the absolute distance so that the measure ranges from 0 to 1.

    Parameters:
        data: (dataframe) It takes dataframe containing absolute distance
        colList: (list) It takes the list of columns which needs to be standardized.
        To standardize, it uses MinMaxScaler from scikit-learn library

    Return: It returns the dataframe with (1 - standardized distance) and all other given columns
    """
    df = data.copy()

    from sklearn.preprocessing import MinMaxScaler
    minMaxScaler = MinMaxScaler()
    df[colList] = minMaxScaler.fit_transform(df[colList])

    for col in colList:
        df.loc[:, col] = 1 - df.loc[:, col]

    return df
```

```
In [16]: def controlStores(data, trialStore):
    """
    It returns a dataframe with all 4 scores correlation score for sales and number of customers, absol
    ute difference score for sales and number of customers of each store to the given trial store.

    It merges all 4 dataframes returned from different functions into one single dataframe.

    Parameters:
        data: (dataframe) It takes dataframe containing store number, months, sales and number o
        f customers
        trialStore: (integer) It takes the trial store number from which we need to compare all
        stores

    """
    controlStoresDf = corrSales(data, trialStore)
    controlStoresDf = corrCust(data, trialStore)
    controlStoresDf = magDistSales(data, trialStore)
    controlStoresDf = magDistCust(data, trialStore)

    #Standardizing the magnitude distance so that the measure ranges from 0 to 1
    controlStoresDf = stdMagDis(controlStoresDf, ["distSales", "distCust"])

    #creating a column to show the compared trial store number
    controlStoresDf.loc[:, "trialStore"] = trialStore

    #returning final dataframe
    return controlStoresDf
```

```
Out[16]:
```

	STORE_NBR	corrSales	corrCust	distSales	distCust	trialStore
0	1	0.875218	0.322168	0.959093	0.938744	77
1	2	-0.263079	-0.572051	0.933198	0.921899	77
2	3	0.806644	0.834207	0.317676	0.318530	77
3	4	-0.263300	-0.295639	0.130226	0.157734	77
4	5	-0.110652	0.370659	0.528558	0.460949	77
...
255	268	-0.021429	0.672672	0.155314	0.325527	88
256	269	-0.157578	-0.474293	0.418570	0.315467	77
257	270	0.315430	-0.131259	0.414653	0.333844	77
258	271	0.356487	0.019629	0.523884	0.459418	77
259	272	0.176222	0.223217	0.877596	0.946401	77

260 rows x 6 columns

Let's see all four score corresponding to Trial Store (77)

```
In [17]: controlStores77 = controlStores(preTrialMeasures, 77)
controlStores77
```

```
Out[17]:
```

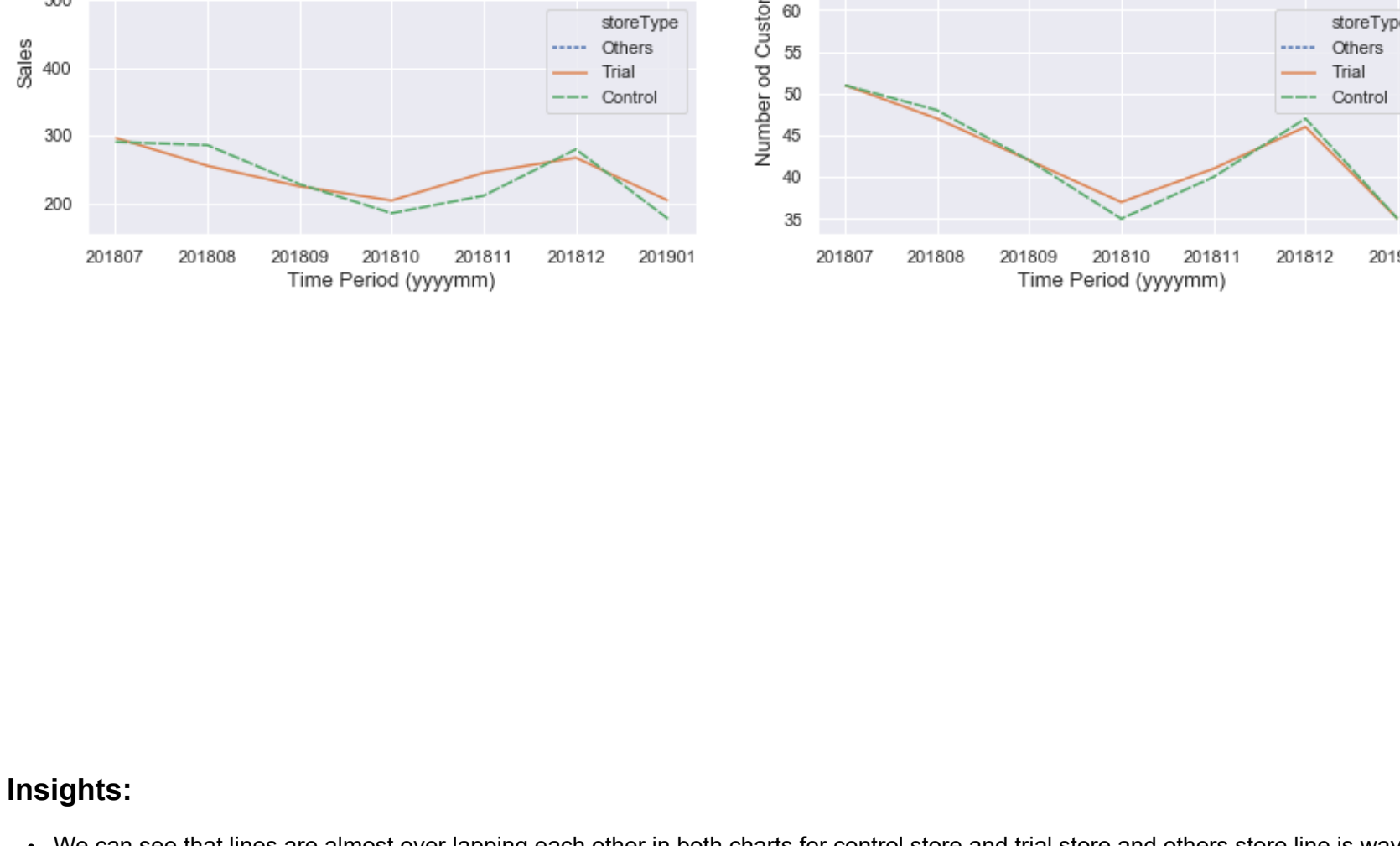
```
controlStoresDf = controlStoresDf.merge(right=corrCust(data, trialStore), how="inner", on=
#)

#merging controlStoresDf and mapDistSales -> controlStoresDf
controlStoresDf = controlStoresDf.merge(right=mapDistSales(data, trialStore), how="inner",
E_NBR)

#merging controlStoresDf and mapDistCust -> controlStoresDf
controlStoresDf = controlStoresDf.merge(right=mapDistCust(data, trialStore), how="inner",
NBR)

#Standardising the magnitude distance so that the measure ranges from 0 to 1
```


In [27]: controlStoresPriorTrialPeriodViz(measureOverTime, controlStore=233, trialStore= 77, savfig=True)



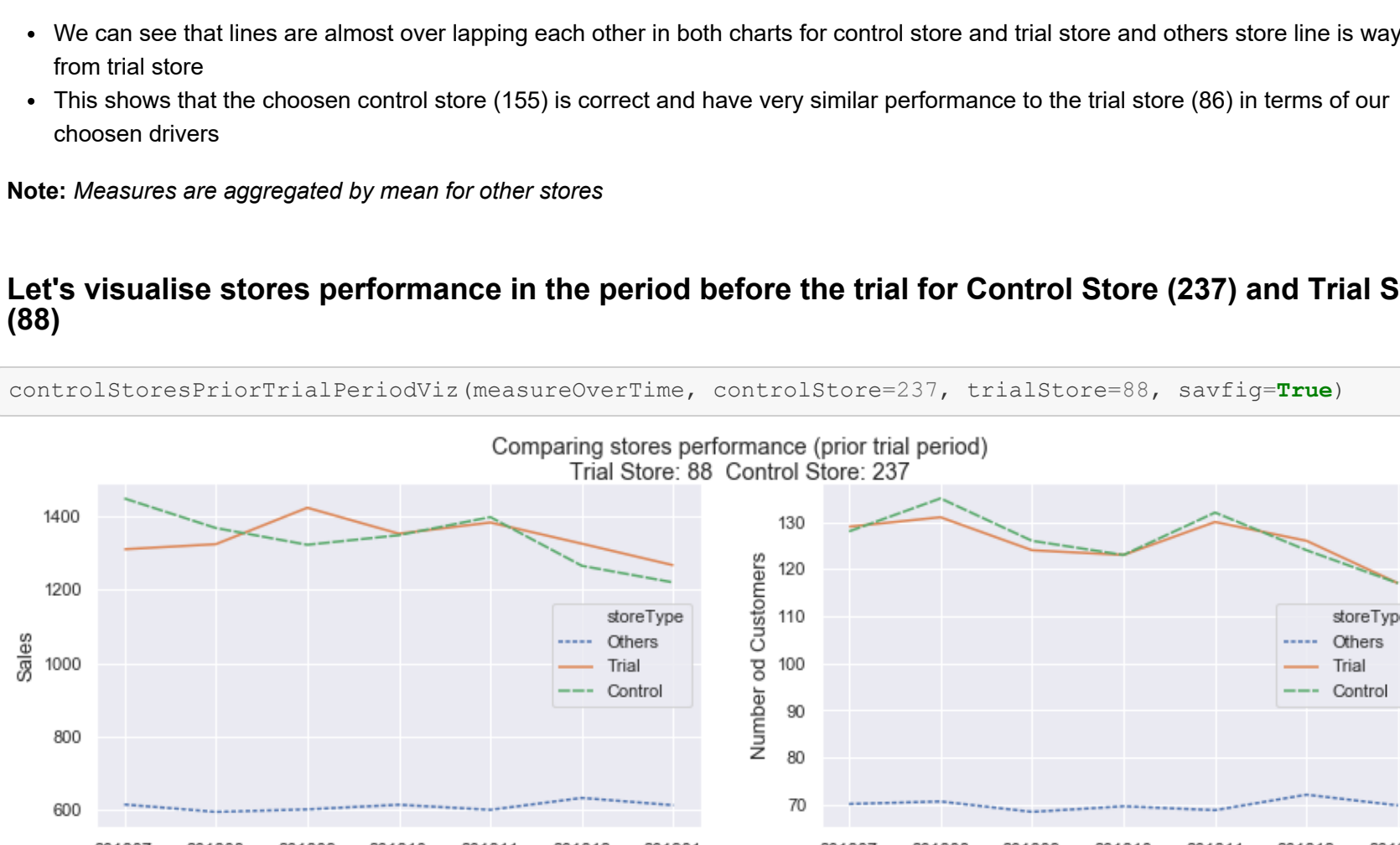
Insights:

- We can see that lines are almost over lapping each other in both charts for control store and trial store and others store line is way far from trial store
- This shows that the chosen control store (233) is correct and have very similar performance to the trial store (77) in terms of our chosen drivers

Note: Measures are aggregated by mean for other stores

Let's visualise stores performance in the period before the trial for Control Store (155) and Trial Store (86)

In [28]: controlStoresPriorTrialPeriodViz(measureOverTime, controlStore=155, trialStore=86, savfig=True)



Insights:

- We can see that lines are almost over lapping each other in both charts for control store and trial store and others store line is way far from trial store
- This shows that the chosen control store (155) is correct and have very similar performance to the trial store (86) in terms of our chosen drivers

Note: Measures are aggregated by mean for other stores

Let's visualise stores performance in the period before the trial for Control Store (237) and Trial Store (88)

In [29]: controlStoresPriorTrialPeriodViz(measureOverTime, controlStore=237, trialStore=88, savfig=True)



Insights:

- We can see that lines are almost over lapping each other in both charts for control store and trial store and others store line is way far from trial store
- This shows that the chosen control store (237) is correct and have very similar performance to the trial store (88) in terms of our chosen drivers

Note: Measures are aggregated by mean for other stores

Phase 2: Assessment of Trial Stores during the Trial Period

(February 2019 - April 2020)

Our goal is to see if there has been an uplift in overall chip sales in trial stores. We now have control stores to compare compare the results and make conclusion.

To do that, *first we need to scale the control store's sales to a level similar to control for any differences between the two stores outside of the trial period.*

Assessment of sales of Trial Store (77)

Calculating Scaling Factor for Control Store's Sale

$$Scaling\ Factor = \frac{\sum Trial\ Store's\ Sales\ (prior\ trial\ period)}{\sum Control\ Store's\ Sales\ (prior\ trial\ period)}$$

In [30]: scalingFactorForControlSales = preTrialMeasures[preTrialMeasures["STORE_NBR"]==77][["totSales"].sum() / preTrialMeasures[preTrialMeasures["STORE_NBR"]==233][["totSales"].sum()]

Out[30]: 1.023617303289533

Creating a dataframe with control store's data only for all months

In [31]: scaledControlSales = measureOverTime[measureOverTime["STORE_NBR"]==233][["STORE_NBR", "MonthID", "totSales"]]

Out[31]:

	STORE_NBR	MonthID	totSales
2699	233	201807	290.7
2700	233	201808	285.9
2701	233	201809	228.6
2702	233	201810	185.7
2703	233	201811	211.6
2704	233	201812	279.8
2705	233	201901	177.5
2706	233	201902	244.0
2707	233	201903	199.1
2708	233	201904	158.6
2709	233	201905	344.4
2710	233	201906	221.0

Scaling control store's sales

$$controlSales = Control\ Store's\ Sales * Scaling\ Factor$$

In [32]: scaledControlSales["controlSales"] = scaledControlSales["totSales"] * scalingFactorForControlSales

Out[32]:

	STORE_NBR	MonthID	totSales	controlSales
2699	233	201807	290.7	297.565550
2700	233	201808	285.9	292.652187
2701	233	201809	228.6	233.998916
2702	233	201810	185.7	190.085733
2703	233	201811	211.6	216.597421
2704	233	201812	279.8	286.408121
2705	233	201901	177.5	181.692071
2706	233	201902	244.0	249.762622
2707	233	201903	199.1	203.802205
2708	233	201904	158.6	162.345704
2709	233	201905	344.4	352.533799
2710	233	201906	221.0	226.219424

Merging trial store's sales in the same dataframe

In [33]: trialSales = measureOverTime[measureOverTime["STORE_NBR"]==77][["MonthID", "totSales"]]

scaledControlSales = scaledControlSales.merge(trialSales, on="MonthID")

scaledControlSales.rename(columns={"totSales_x": "totSales", "totSales_y": "trialSales"}, inplace=True)

Out[33]:

	STORE_NBR	MonthID	totSales	controlSales	trialSales
0	233	201807	290.7	297.565550	296.8
1	233	201808	285.9	292.652187	255.5
2	233	201809	228.6	233.998916	225.2
3	233	201810	185.7	190.085733	204.5
4	233	201811	211.6	216.597421	245.3
5	233	201812	279.8	286.408121	267.3
6	233	201901	177.5	181.692071	204.4
7	233	201902	244.0	249.762622	235.0
8	233	201903	199.1	203.802205	278.5
9	233	201904	158.6	162.345704	263.5
10	233	201905	344.4	352.533799	299.3
11	233	201906	221.0	226.219424	264.7

Calculating Percentage Difference in Sales

Now that we have comparable sales figures for the abs store, we can calculate the percentage difference between the scaled control sales and the trial store's sales during the trial period.

In [34]: scaledControlSales["percentageDiff"] = np.abs(scaledControlSales["controlSales"] - scaledControlSales["trialSales"]) / scaledControlSales["controlSales"]

Out[34]:

	STORE_NBR	MonthID	totSales	controlSales	trialSales	percentageDiff
0	233	201807	290.7	297.565550	296.8	0.002573
1	233	201808	285.9	292.652187	255.5	0.126950
2	233	201809	228.6	233.998916	225.2	0.037602
3	233	201810	185.7	190.085733	204.5	0.075830
4	233	201811	211.6	216.597421	245.3	0.132516
5	233	201812	279.8	286.408121	267.3	0.066716
6	233	201901	177.5	181.692071	204.4	0.124980
7	233	201902	244.0	249.762622	235.0	0.059107
8	233	201903	199.1	203.802205	278.5	0.366521
9	233	201904	158.6	162.345704	263.5	0.625080
10	233	201905	344.4	352.533799	299.3	0.151003
11	233	201906	221.0	226.219424	264.7	0.170103

Check if the difference is significant!

We need to check and confirm this statistically, and we will do hypothesis testing. Our hypothesis will be,

Null Hypothesis: There is no difference in sales in trial and pre-trial period

Alternate Hypothesis: There is a difference in sales in trial and pre-trial period

Let's calculate the standard deviation based on the scaled percentage difference in the pre-trial period

In [35]: stdDev = scaledControlSales[scaledControlSales["MonthID"] < 201902][["percentageDiff"].std()

Out[35]: 0.049940762641425544

Let's calculate the t-value

Our focus will be for trial period only

We are testing with a null hypothesis of there being 0 difference between trial and control stores, that's why we have subtracted zero, $t - value = \frac{percentage\ Difference - 0}{Standard\ Deviation}$

In [36]: scaledControlSales["tValue"] = (scaledControlSales["percentageDiff"] - 0) / stdDev

Out[36]:

	STORE_NBR	MonthID	totSales	controlSales	trialSales	percentageDiff	tValue
0	233	201807	290.7	297.565550	296.8	0.002573	0.051515
1	233	201808	285.9	292.652187	255.5	0.126950	2.540117
2	233	201809	228.6	233.998916	225.2	0.037602	0.752940
3	233	201810	185.7	190.085733	204.5	0.075830	1.518406
4	233	201811	211.6	216.597421	245.3	0.132516	2.653490
5	233	201812	279.8	286.408121	267.3	0.066716	1.339911
6	233	201901	177.5	181.692071	204.4	0.124980	2.502571
7	233	201902	244.0	249.762622	235.0	0.059107	1.183534
8	233	201903	199.1	203.802205	278.5	0.366521	7.339116
9	233	201904	158.6	162.345704	263.5	0.625080	12.476373
10	233	201905	344.4	352.533799	299.3	0.151003	3.023650
11	233	201906	221.0	226.219424	264.7	0.170103	4.406093

The 95th percentile of the t-distribution with the appropriate degrees of freedom to check whether the hypothesis is statistically significant.

Since there are 7 months in the pre-trial period

Hence, degrees of freedom = 6 (7 - 1)

In [37]: #0.95 taken because we are interested to know 95 percentile of t-distribution

degreesOfFreedom = 6

stats.t.ppf(0.95, df=degreesOfFreedom)

Out[37]: 1.9431802803927816

Comparing t-value and 95th percentile of the t-distribution for trial period

In [38]: scaledControlSales[scaledControlSales["MonthID"].isin([201902, 201903, 201904])]

Out[38]:

	STORE_NBR	MonthID	totSales	controlSales	trialSales	percentageDiff	tValue
7	233	201902	244.0	249.762622	235.0	0.059107	1.183534
8	233	201903	199.1	203.802205	278.5	0.366521	7.339116
9	233	201904	158.6	162.345704	263.5	0.625080	12.476373

Observation:

We can observe that the t-value is much larger than the 95th percentile value of the t-distribution for March and April - i.e. the increase in sales in the trial store in March and April is statistically greater than in the control store.

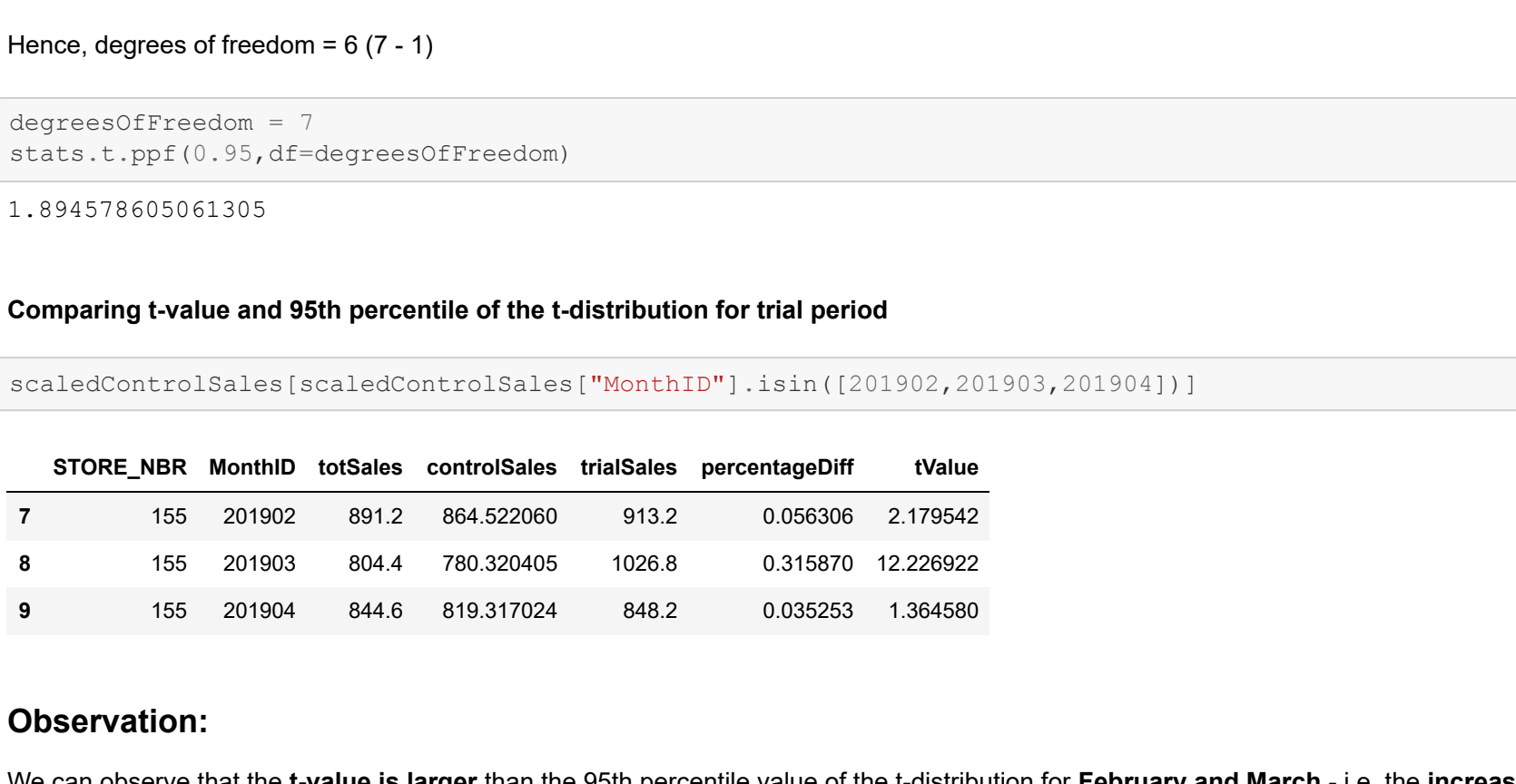
Visualizing the Sales of Trial and Control Store

Let's create a more visual version of this by plotting the sales of the control store, trial store, other stores and the 5% and 95% confidence level of control store's sales.

Let's define a function for this which we can use for other stores also

In [39]: def trialPeriodViz(data, controlStore, trialStore, stdDev, on, savfig=False):
 """
 It plots a line plot sales/ number of customers of Trial Store, Control Store, and Others Store along with 5% and 95% percentile confidence level of control store's sales/ number of customers over the whole time period.
 Parameters:
 data: (dataframe) It takes dataframe with store number, months, sales and number of customers for all months
 controlStore: (integer) It takes control store's number
 trialStore: (integer) It takes trial store's number
 stdDev: (float) It takes the standard deviation which is used in hypothesis testing
 on: (string) It takes only two values "totSales", "nCustomers" to plot using corresponding measure
 savfig: (boolean) It asks to save the chart on the local system or not, by default it is False
 It creates a single line for stores other than control and trial store by aggregating values by mean
 """
 df = data[["STORE_NBR", "MonthID", on]]
 df["storeType"] = df[df["STORE_NBR"]==controlStore].apply(lambda x: "Control" if x==controlStore else "Others")
 df.loc[df["STORE_NBR"]==trialStore, "storeType"] = "Trial"
 Control95 = pd.DataFrame(columns=["STORE_NBR", "MonthID", on, "storeType"])
 Control95[on] = df[df["STORE_NBR"]==controlStore][on] * (1 + 2 * stdDev)
 Control95.loc[:, "storeType"] = "Control 5th % confidence interval"
 df = df.append(Control95, ignore_index=True)
 df = df.append(Control95, ignore_index=True)
 df["MonthID"] = df["MonthID"].astype(str)
 name=""
 if on == "totSales":
 name="Sales"
 else:
 name="Number of customers"
 plt.figure(figsize=(15,6))
 sns.lineplot(data=df, x="MonthID", y=on, hue="storeType", estimator="mean", err_style=None, style="storeType")
 dashes=[(1,0), (1,0), (1,0), (3,1), (3,1)]
 plt.xlabel("Time Period (yyyymm)", size=15)
 plt.ylabel(name, size=15)
 areaMax = df.groupby(["MonthID", "storeType"])[on].mean().max()
 plt.fill_between(["201902", "201904"], areaMax, color="skyblue", alpha=0.5)
 ytext = df[df["storeType"]=="Others"].groupby(["MonthID"])[on].mean().min() - 20
 plt.text("201902", ytext, s="Others")
 plt.title(f"{trialStore} (TrialStore) Control Store (controlStore)", size=15)
 plt.suptitle(name + " over the period", size=20)
 if savfig==True:
 path=f"C:/Users/amit/Jupyter Notebooks/Quantum/Task 2/{name+ ' All ' + str(trialStore)}.jpg"
 plt.savefig(path, dpi=500)

Out[39]: trialPeriodViz(measureOverTime, controlStore=233, trialStore=77, stdDev=stdDev, on="totSales", savfig=True)



Observation:

The results show that the trial in store 77 is significantly different to its control store in the trial period as the trial store performance lies outside the 5% and 95% confidence interval of the control store in two of the three trial months.

Assessment of Number of Customers of Trial Store (77)

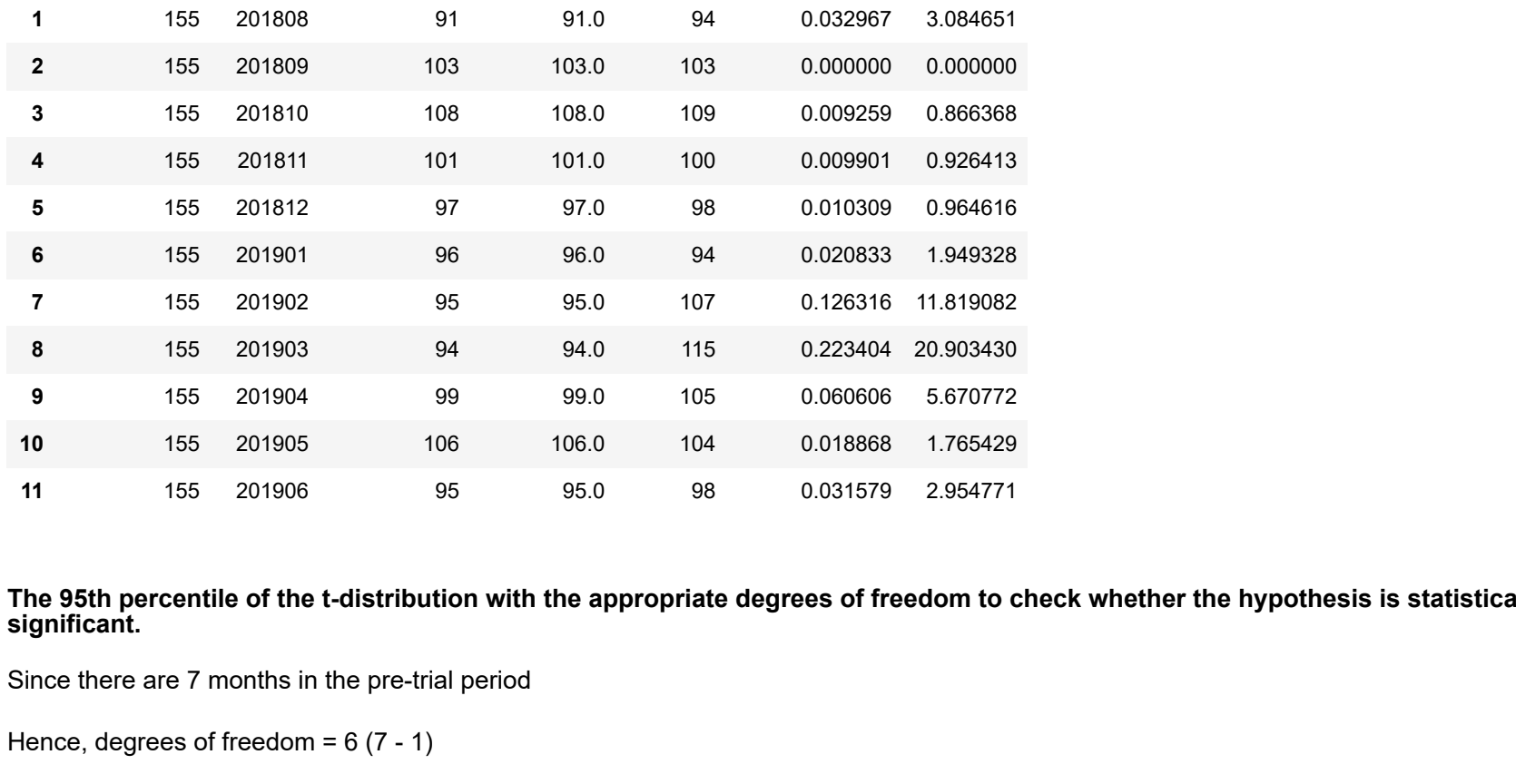
Let's now see if there has been an uplift in number of customers in trial stores during trial period.

To do this, we need perform the same calculations we did above on sales, but now on number of customers.

Let's define a function for all this, it will help in assessing other trial stores also and it will save our lot of time

In [41]: def trialPeriodDiff(data, controlStore, trialStore, scalingFactorForControl, on):
 """
 This function returns the dataframe like we got above for hypothesis testing and it also returns standard deviation which is calculated on percentage difference of pre-trial period
 Parameters:
 data: (dataframe) It takes dataframe with store number, months, sales and number of customers for all months
 controlStore: (integer) It takes control store's number
 trialStore: (integer) It takes trial store's number
 scalingFactorForControl: (float) It takes the scaling factor through which we will get control sales/no. of cust
 on: (string) It takes only two values "totSales", "nCustomers" to plot using corresponding measure
 Return: Dataframe, standard deviation (float)
 """
 contname=""
 trialname=""
 if on=="totSales":
 contname="controlSales"
 trialname="trialSales"
 else:
 contname="controlCust"
 trialname="trialCust"
 scaledControl = data[data["STORE_NBR"]==controlStore][contname]
 scaledControl[contname] = scaledControl[on] * scalingFactorForControl
 trial = data[data["STORE_NBR"]==trialStore][contname]
 trial.rename(columns=(contname), inplace=True)
 scaledControl = scaledControl.merge(trial, on="MonthID")
 scaledControl["percentageDiff"] = np.abs(scaledControl[contname] - scaledControl[trialname]) / scaledControl[contname]
 stdDev = scaledControl[scaledControl["MonthID"] < 201902][contname].std()
 scaledControl["tValue"] = (scaledControl["percentageDiff"] - 0) / stdDev
 return scaledControl, stdDev

Out[41]: trialPeriodDiff(data, controlStore=233, trialStore=77, scalingFactorForControl=on="nCustomers", savfig=True)



Observation:

The results show that the trial in store 86 is not significantly different to its control store in the trial period as the trial store performance lies inside the 5% and 95% confidence interval of the control store in two of the three trial months.

Assessment of Number of Customers of Trial Store (86)

Calculating Scaling Factor for Control Store's Number of Customers

In [52]: scalingFactor= preTrialMeasures[preTrialMeasures["STORE_NBR"]==86][["nCustomers"].sum() / preTrialMeasures[preTrialMeasures["STORE_NBR"]==155][["nCustomers"].sum()]

Out[52]: 1.0

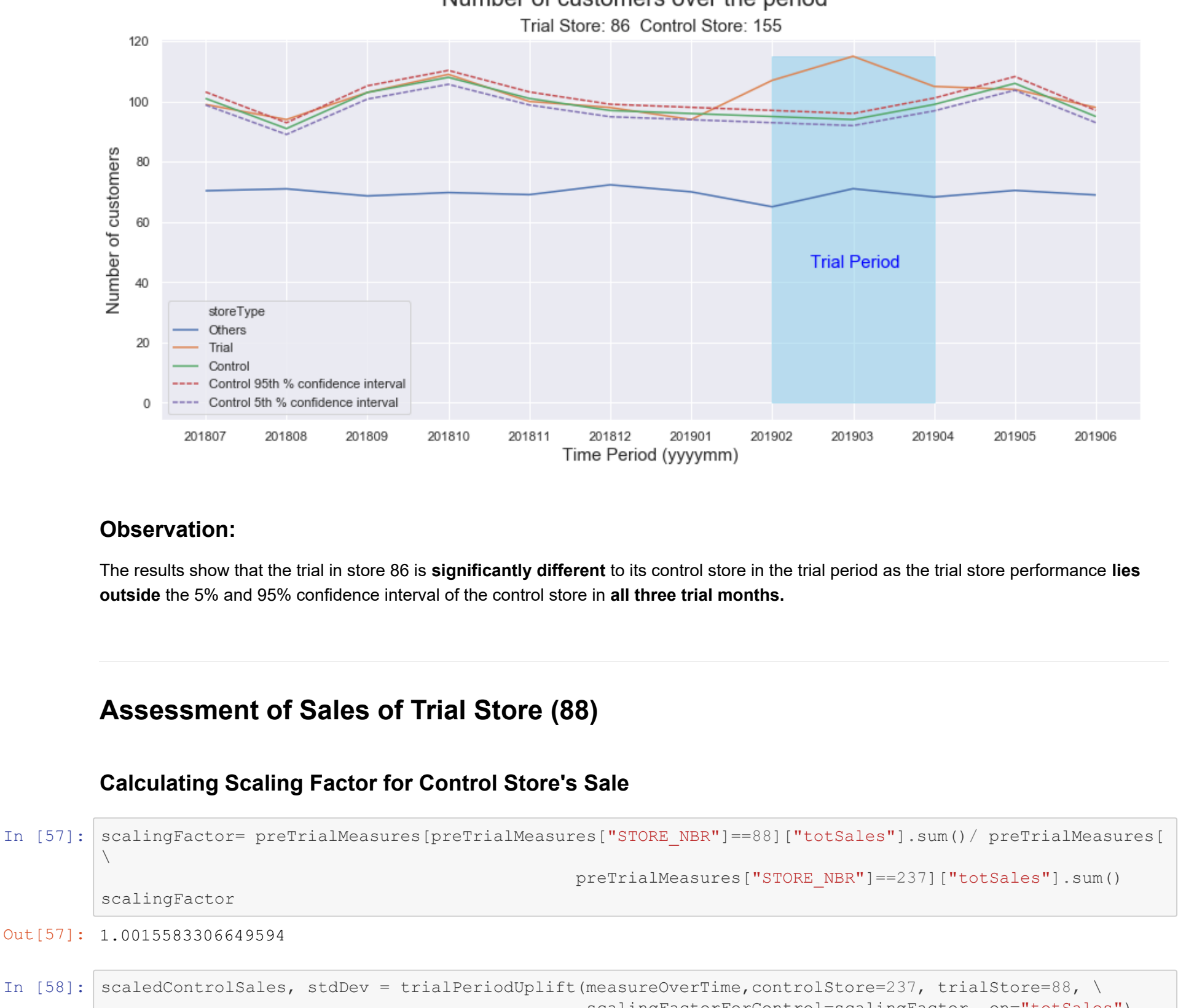
In [53]: scaledControlCust, stdDev = trialPeriodDiff(measureOverTime, controlStore=155, trialStore=86, \ scalingFactorForControl=scalingFactor, on="nCustomers")

Out[53]:

	STORE_NBR	MonthID	nCustomers	controlCust	trialCust	percentageDiff	tValue
0	155	201807	101	101.0	99	0.019802	1.852826
1	155	201808	91	91.0	94	0.032667	3.084651
2	155	201809	103	103.0	103	0.000000	0.000000
3	155	201810	108	108.0	109	0.009259	0.866368
4	155	201811	101	101.0	100	0.009901	0.926413
5	155	201812	97	97.0	98	0.010309	0.964616
6	155	201901	96	96.0	94	0.020833	1.949328
7	155	201902	95	95.0	107	0.126316	11.819082
8	155	201903	94	94.0	115	0.223404	20.903430
9	155	201904	99	99.0	105	0.	

Visualizing the Number of Customers of Trial and Control Store

Let's create a more visual version of this by plotting the number of customers of the control store, trial stores, other stores and the 5% and 95% confidence level of control store's number of customers.



Observation:

The results show that the trial in store 86 is **significantly different** to its control store in the trial period as the trial store performance **lies outside** the 5% and 95% confidence interval of the control store in **all three trial months**.

Assessment of Sales of Trial Store (88)

Calculating Scaling Factor for Control Store's Sale

```
In [57]: scalingFactor= preTrialMeasures[preTrialMeasures["STORE_NBR"]==88][["totSales"].sum()/ preTrialMeasures["STORE_NBR"]==237][["totSales"].sum())

Out[57]: 1.001558330649594
```

```
In [58]: scaledControlSales, stdDev = trialPeriodUplift(measureOverTime,controlStore=237, trialStore=88, \
scaledControlSales

Out[58]:
```

	STORE_NBR	MonthID	totSales	controlSales	trialSales	percentageDiff	tValue
0	237	201807	1448.4	1450.657086	1310.00	0.066961	2.897136
1	237	201808	1367.8	1369.831485	1323.80	0.033674	1.006168
2	237	201809	1322.2	1324.260425	1423.00	0.074562	2.227870
3	237	201810	1348.3	1350.401097	1352.40	0.001480	0.044228
4	237	201811	1397.6	1399.777923	1382.80	0.012129	0.362408
5	237	201812	1265.0	1266.971288	1325.20	0.049559	1.373227
6	237	201901	1219.7	1221.600696	1266.40	0.036673	1.065756
7	237	201902	1404.8	1406.989143	1370.20	0.026147	0.781270
8	237	201903	1208.2	1210.082775	1477.20	0.220743	6.595668
9	237	201904	1204.6	1206.477165	1439.40	0.193060	5.768527
10	237	201905	1199.3	1201.168906	1308.25	0.089147	2.663672
11	237	201906	1153.6	1155.397690	1354.60	0.172410	5.151513

The **95th percentile** of the **t-distribution** with the appropriate degrees of freedom to check whether the hypothesis is statistically significant.

Since there are 7 months in the pre-trial period

Hence, degrees of freedom = 6 (7 - 1)

```
In [59]: degreesOfFreedom = 7
stats.t.ppf(0.95,df=degreesOfFreedom)

Out[59]: 1.894578605061305
```

Comparing t-value and 95th percentile of the t-distribution for trial period

```
In [60]: scaledControlSales[scaledControlSales["MonthID"].isin([201902,201903,201904])]
```

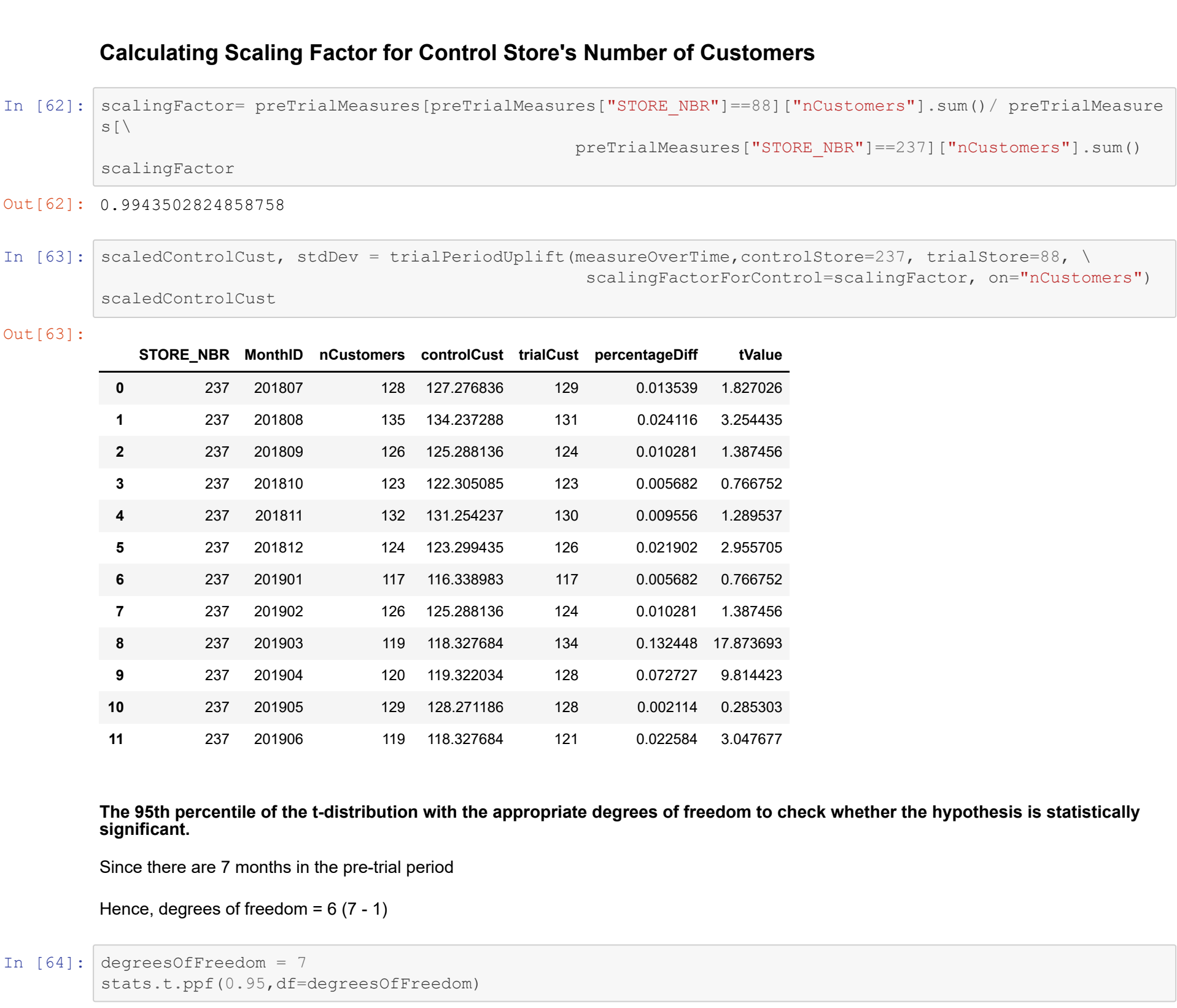
	STORE_NBR	MonthID	totSales	controlSales	trialSales	percentageDiff	tValue
7	237	201902	1404.8	1406.989143	1370.20	0.026147	0.781270
8	237	201903	1208.2	1210.082775	1477.20	0.220743	6.595668
9	237	201904	1204.6	1206.477165	1439.40	0.193060	5.768527

Observation:

We can observe that the **t-value** is **larger** than the 95th percentile value of the t-distribution for **March and April** - i.e. the **increase in sales** in the trial store in March and April is statistically greater than in the control store.

Visualizing the Sales of Trial and Control Store

Let's create a more visual version of this by plotting the sales of the control store, trial stores, other stores and the 5% and 95% confidence level of control store's sales.



Observation:

The results show that the trial in store 88 is **significantly different** to its control store in the trial period as the trial store performance **lies outside** the 5% and 95% confidence interval of the control store in **two of the three trial months**.

Assessment of Number of Customers of Trial Store (88)

Calculating Scaling Factor for Control Store's Number of Customers

```
In [62]: scalingFactor= preTrialMeasures[preTrialMeasures["STORE_NBR"]==88][["nCustomers"].sum()/ preTrialMeasures["STORE_NBR"]==237][["nCustomers"].sum())

Out[62]: 0.9943502824858758
```

```
In [63]: scaledControlCust, stdDev = trialPeriodUplift(measureOverTime,controlStore=237, trialStore=88, \
scaledControlCust

Out[63]:
```

	STORE_NBR	MonthID	nCustomers	controlCust	trialCust	percentageDiff	tValue
0	237	201807	128	127.276836	129	0.013539	1.827026
1	237	201808	135	134.237288	131	0.024116	3.254435
2	237	201809	126	125.289136	124	0.010281	1.387456
3	237	201810	123	122.305085	123	0.005682	0.766752
4	237	201811	132	131.254237	130	0.009556	1.289537
5	237	201812	124	123.299435	126	0.021902	2.955705
6	237	201901	117	116.338983	117	0.005682	0.766752
7	237	201902	126	125.288136	124	0.010281	1.387456
8	237	201903	119	118.327684	134	0.132448	17.873693
9	237	201904	120	119.322034	128	0.072727	9.814423
10	237	201905	129	128.277186	128	0.002114	0.285303
11	237	201906	119	118.327684	121	0.022584	3.047677

The **95th percentile** of the **t-distribution** with the appropriate degrees of freedom to check whether the hypothesis is statistically significant.

Since there are 7 months in the pre-trial period

Hence, degrees of freedom = 6 (7 - 1)

```
In [64]: degreesOfFreedom = 7
stats.t.ppf(0.95,df=degreesOfFreedom)

Out[64]: 1.894578605061305
```

Comparing t-value and 95th percentile of the t-distribution for trial period

```
In [65]: scaledControlCust[scaledControlCust["MonthID"].isin([201902,201903,201904])]
```

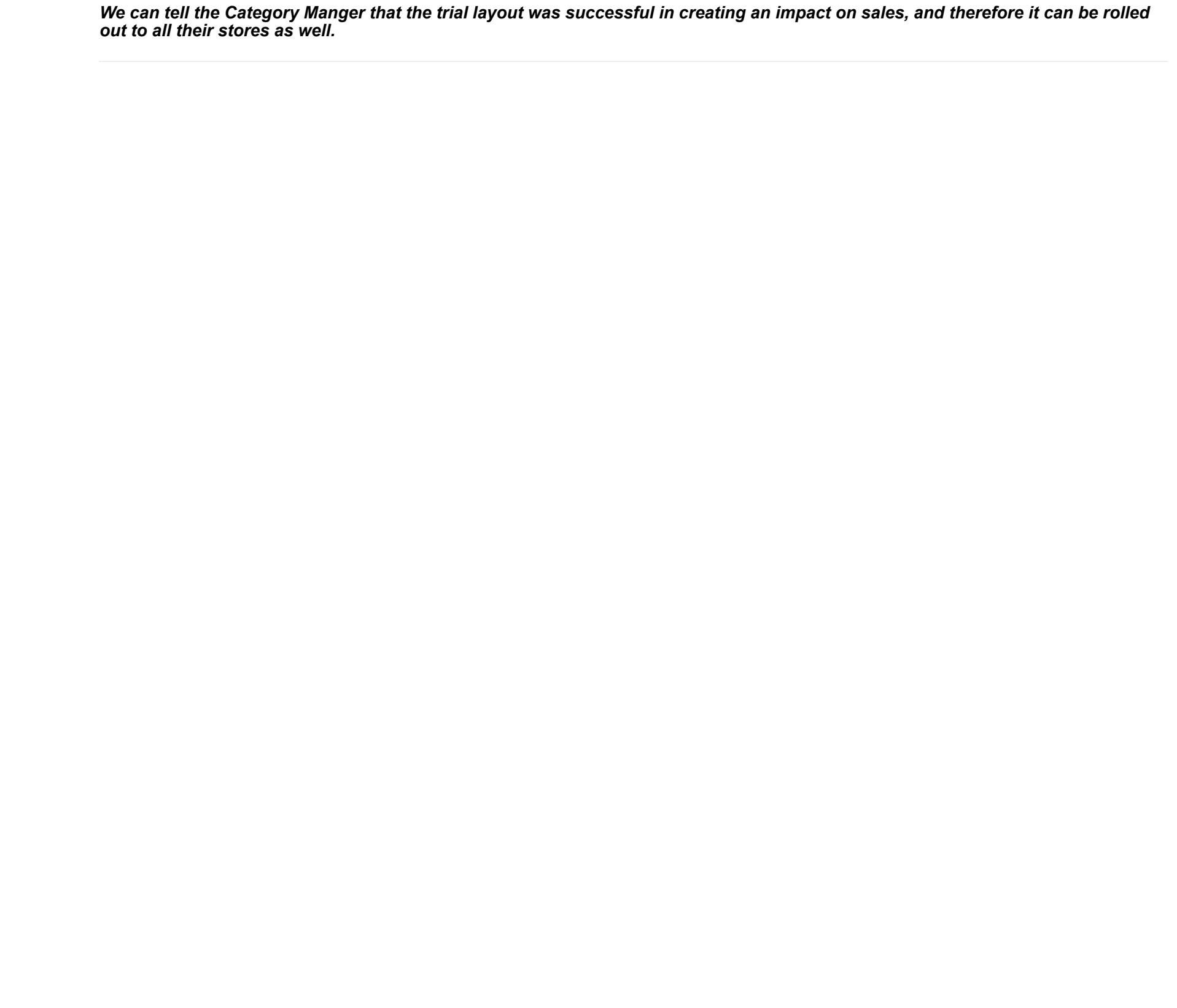
	STORE_NBR	MonthID	nCustomers	controlCust	trialCust	percentageDiff	tValue
7	237	201902	126	125.288136	124	0.010281	1.387456
8	237	201903	119	118.327684	134	0.132448	17.873693
9	237	201904	120	119.322034	128	0.072727	9.814423

Observation:

We can observe that the **t-value** is **much larger** than the 95th percentile value of the t-distribution for **March and April** - i.e. the **increase in number of customers** in the trial store in March and April is statistically greater than in the control store.

Visualizing the Sales of Trial and Control Store

Let's create a more visual version of this by plotting the number of customers of the control store, trial stores, other stores and the 5% and 95% confidence level of control store's number of customers.



Observation:

The results show that the trial in store 88 is **significantly different** to its control store in the trial period as the trial store performance **lies outside** the 5% and 95% confidence interval of the control store in **two of the three trial months**.

Conclusion

First we found control stores for trial stores using pre-trial period data

- Control store 233 for trial store 77
- Control store 155 for trial store 86
- Control store 237 for trial store 88

Then we assessed trial store's performance in trial period

- We found there was a significant difference in the sales and number of customers for two out of three trial months for trial stores 77 and 88
 - Sales and number of customers increased for those months
- However, for trial store 86 there was no significant difference in the sales but there was significant difference in the number of customers
 - Number of customers were increased for all three trial months
 - Sales was increased for one out of three trial months
- This shows that, trial store 86 was able to attract significantly more customers in trial period than pre-trial period however, sales did not increased significantly. We need to check with the category manager, if some different strategy was applied at this store.

Overall, impact of the new trial layout at trial stores shows a significant increase in sales.

We can tell the Category Manager that the trial layout was successful in creating an impact on sales, and therefore it can be rolled out to all their stores as well.