

# Customer Analytics¶

## What is A/B test

- test two or more variants against each other
- to evaluate which one performs best
- in the context of randomized experiment

## A/B testing Process

1. develop a hypothesis about your product or business
2. randomly assign users to two different groups
  - group 1 to current product rules
  - group 2 to a product that test hypothesis
3. pick according to KPI

## Key performance indicators (KPIs)

\*\*How to identify KPIs:

- experience+Domain knowledge + exploratory data analysis
- stability over time
- the importance with other factors

## Question: which paywall has a higher conversion rate

- Current paywall: "I hoped that you enjoyed your free-trial, please consider subscribing"(control)
- Proposed paywall: "Your free-trial has ended, do not miss out, subscribe today"(experiment)

## Process

- Randomly subset the users and show one set the control and one the treatment
- Monitor the conversion rates of each group to see which is better

## The importance of randomness

- isolate the impact of the change made
- reduce the potential impact of confounding variables
- use a assignment criteria may introduce confounders

## Pros and Cons

- Pros: users are impacted individually; testing changes that directly impact their behavior
- Cons: challenging to segment the users into groups; difficult to untangle the impact of the test

# Initial A/B test design

## Responsabe variable

- The quantity used to measure the impact of your change
- Should either be a KPI or directly related to a KPI
- The easier to measure, the better

## Factors & variants

- Factors: the type of variable you are changing
  - The paywall color
- Variants: Particular changes you are testing
  - A red vs a blue paywall

## Experimental unit of our test

- The smallest unit you are measuring the change over
- Individual users make a convenient experimental unit

# Preparing to run an A/B test

## test sensitivity

- First question: what size of impact is meaningful to test
- Smaller changes=more difficult to dect
- Sensitivity: the minimum level of change we want to be able to detect in our test
  - Evaluate different sensitivity values

## Data variability

- Standard Error
  - the variance and standard deviation of the estimate
- Hypothesis
  - Null Hypothesise: The launch of feature does not statistical significant change the KPI
  - Reject Null Hypothesise: The launch of feature does statistical significant change the KPI

## Parameters

- Statistical Power
  - Probability of finding a statistically significant result when the Null Hypothesis is false
- Estimate sample size
  - needed level of sensitivity
  - our desired test power & confidence level
- Effect size

## Significance or Not?

- P-value
  - probability if the Null Hypothesis is true...

- Low p-values: the observation is unlikely to have happened due to randomness
- Test
  - how significant the differences between groups are
  - One-sample:same population
  - Two-sample: different population

## Our data set

We are looking at data from an app. The app is very simple and has just 4 pages:

- The first page is the home page. When you come to the site for the first time, you can only land on the home page as a first page.
- From the home page, the user can perform a search and land on the search page.
- From the search page, if the user clicks on a product, she will get to the payment page (paywall), where she is asked to provide payment information in order to subscribe.
- If she does decide to buy, she ends up on the confirmation page

**Data set overview** We have 5 files, 4 of them contains page\_visit information and 1 of them contains user information.

- page\_visit\_information
  - home\_page\_table.csv
  - search\_page\_table.csv
  - payment\_page\_table.csv
  - payment\_confirmation\_table.csv
- user\_information page
  - user\_table.csv

## Import

In [1]:

```
1 import warnings
2 warnings.filterwarnings('ignore')
3 import pandas as pd
4 import numpy as np
5 import matplotlib.pyplot as plt
6 import seaborn as sns
7 from functools import reduce
8 from sklearn import preprocessing
9 from scipy import stats
```

In [2]:

```
1 user_table = pd.read_csv('user_table.csv')
```

In [3]:

```
1 user_table
```

Out[3]:

|       | user_id | date       | device  | sex    |
|-------|---------|------------|---------|--------|
| 0     | 450007  | 2015-02-28 | Desktop | Female |
| 1     | 756838  | 2015-01-13 | Desktop | Male   |
| 2     | 568983  | 2015-04-09 | Desktop | Male   |
| 3     | 190794  | 2015-02-18 | Desktop | Female |
| 4     | 537909  | 2015-01-15 | Desktop | Male   |
| ...   | ...     | ...        | ...     | ...    |
| 90395 | 307667  | 2015-03-30 | Desktop | Female |
| 90396 | 642989  | 2015-02-08 | Desktop | Female |
| 90397 | 659645  | 2015-04-13 | Desktop | Male   |
| 90398 | 359779  | 2015-03-23 | Desktop | Male   |

## Create control and test groups

randomized two groups

In [4]:

```
1 length=len(user_table['user_id'])
```

In [5]:

```
1 k=np.random.binomial(1,0.495,length)
```

In [6]:

```
1 user_table['group']=k
2 user_table.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 90400 entries, 0 to 90399
Data columns (total 5 columns):
#   Column      Non-Null Count  Dtype
---  -
0   user_id     90400 non-null  int64
1   date        90400 non-null  object
2   device      90400 non-null  object
3   sex         90400 non-null  object
4   group       90400 non-null  int64
dtypes: int64(2), object(3)
memory usage: 3.4+ MB
```

In [7]:

```
1 user_table['group']=user_table['group'].replace(1,'test')
2 user_table['group']=user_table['group'].replace(0,'control')
```

## import datasets

In [8]:

```
1 home_page = pd.read_csv('home_page_table.csv')
2 payment_confirmation= pd.read_csv('payment_confirmation_table.csv')
3 payment_page= pd.read_csv('payment_page_table.csv')
4 search_page= pd.read_csv('search_page_table.csv')
```

In [9]:

```
1 home_page.columns
2 payment_confirmation.columns
```

Out[9]:

```
Index(['user_id', 'page'], dtype='object')
```

In [171]:

```
1 # Compile the list of dataframes you want to merge
2 data_frames = [user_table, home_page, search_page, payment_page, payment_confirm
3 df_merged = reduce(lambda left,right: pd.merge(left,right,on=['user_id'], how='
4 df_merged.columns = ['user_id', 'date', 'device', 'sex', 'group', 'home_page', '
5                     'payment_page', 'payment_confirm']
6 df_merged.info()
```

```
<class 'pandas.core.frame.DataFrame'>
```

```
Int64Index: 90400 entries, 0 to 90399
```

```
Data columns (total 9 columns):
```

| # | Column          | Non-Null Count | Dtype  |
|---|-----------------|----------------|--------|
| 0 | user_id         | 90400 non-null | int64  |
| 1 | date            | 90400 non-null | object |
| 2 | device          | 90400 non-null | object |
| 3 | sex             | 90400 non-null | object |
| 4 | group           | 90400 non-null | object |
| 5 | home_page       | 90400 non-null | object |
| 6 | search_page     | 45200 non-null | object |
| 7 | payment_page    | 6030 non-null  | object |
| 8 | payment_confirm | 452 non-null   | object |

```
dtypes: int64(1), object(8)
```

```
memory usage: 6.9+ MB
```

In [10]:

```
1 df=user_table.merge(home_page,on=['user_id'],how='left')
```

In [11]:

```
1 df=df.merge(search_page,on=['user_id'],how='left')
```

In [12]:

```
1 df=df.merge(payment_page,on=['user_id'],how='left')
```

In [13]:

```
1 df=df.merge(payment_confirmation,on=['user_id'],how='left')
```

In [14]:

1 df

Out[14]:

|       | user_id | date       | device  | sex    | group   | page_x    | page_y      | page_x | page_y |
|-------|---------|------------|---------|--------|---------|-----------|-------------|--------|--------|
| 0     | 450007  | 2015-02-28 | Desktop | Female | control | home_page | NaN         | NaN    | NaN    |
| 1     | 756838  | 2015-01-13 | Desktop | Male   | control | home_page | NaN         | NaN    | NaN    |
| 2     | 568983  | 2015-04-09 | Desktop | Male   | test    | home_page | search_page | NaN    | NaN    |
| 3     | 190794  | 2015-02-18 | Desktop | Female | test    | home_page | search_page | NaN    | NaN    |
| 4     | 537909  | 2015-01-15 | Desktop | Male   | control | home_page | NaN         | NaN    | NaN    |
| ...   | ...     | ...        | ...     | ...    | ...     | ...       | ...         | ...    | ...    |
| 90395 | 307667  | 2015-03-30 | Desktop | Female | test    | home_page | NaN         | NaN    | NaN    |
| 90396 | 642989  | 2015-02-08 | Desktop | Female | test    | home_page | search_page | NaN    | NaN    |
| 90397 | 659645  | 2015-04-13 | Desktop | Male   | test    | home_page | search_page | NaN    | NaN    |
| 90398 | 359779  | 2015-03-23 | Desktop | Male   | test    | home_page | NaN         | NaN    | NaN    |
| 90399 | 438929  | 2015-03-26 | Mobile  | Female | control | home_page | NaN         | NaN    | NaN    |

90400 rows × 9 columns

In [15]:

```
1 df.columns=['user_id','date','device','sex','group','home_page','search_page',
2            'payment_page','payment_confirmation']
```

In [16]:

```
1 df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 90400 entries, 0 to 90399
Data columns (total 9 columns):
#   Column                Non-Null Count  Dtype
---  -
0   user_id               90400 non-null  int64
1   date                  90400 non-null  object
2   device                90400 non-null  object
3   sex                   90400 non-null  object
4   group                 90400 non-null  object
5   home_page             90400 non-null  object
6   search_page           45200 non-null  object
7   payment_page          6030 non-null   object
8   payment_confirmation  452 non-null    object
dtypes: int64(1), object(8)
memory usage: 6.9+ MB
```

In [17]:

```
1 df['date']=pd.to_datetime(df['date'],utc=True)
```

In [18]:

```
1 df['home_page']=df['home_page'].fillna(0)
2 df['search_page']=df['search_page'].fillna(0)
3 df['payment_page']=df['payment_page'].fillna(0)
4 df['payment_confirmation']=df['payment_confirmation'].fillna(0)
```

In [19]:

```
1 df['home_page']=df['home_page'].replace('home_page',1)
2 df['search_page']=df['search_page'].replace('search_page',1)
3 df['payment_page']=df['payment_page'].replace('payment_page',1)
4 df['payment_confirmation']=df['payment_confirmation'].replace('payment_confirmation',1)
```

In [20]:

```
1 df['search_page'].value_counts()
```

Out[20]:

```
0    45200
1    45200
Name: search_page, dtype: int64
```

In [21]:

```
1 df['payment_page'].value_counts()
```

Out[21]:

```
0    84370
1     6030
Name: payment_page, dtype: int64
```

In [23]:

```
1 df['payment_confirmation'].value_counts()
```

Out[23]:

```
0    89948
```

```
1      452
```

```
Name: payment_confirmation, dtype: int64
```

## EAD

- overall Conversion rate
- group by gender
- group by device

In [24]:

```
1 df_purchase=df.groupby(by=['date'],as_index=False)
```

In [25]:

```
1 df_purchase=df_purchase.agg({'payment_confirmation':['sum','count']})
```

In [26]:

```
1 df_purchase.columns=df_purchase.columns.droplevel(level=0)
```

In [27]:

```
1 df_purchase.columns=['date','sum','count']
```

In [28]:

```
1 df_purchase
```

Out[28]:

|     | date                      | sum | count |
|-----|---------------------------|-----|-------|
| 0   | 2015-01-01 00:00:00+00:00 | 6   | 712   |
| 1   | 2015-01-02 00:00:00+00:00 | 5   | 721   |
| 2   | 2015-01-03 00:00:00+00:00 | 7   | 760   |
| 3   | 2015-01-04 00:00:00+00:00 | 7   | 713   |
| 4   | 2015-01-05 00:00:00+00:00 | 11  | 754   |
| ... | ...                       | ... | ...   |
| 115 | 2015-04-26 00:00:00+00:00 | 1   | 792   |
| 116 | 2015-04-27 00:00:00+00:00 | 1   | 779   |
| 117 | 2015-04-28 00:00:00+00:00 | 1   | 736   |
| 118 | 2015-04-29 00:00:00+00:00 | 0   | 713   |



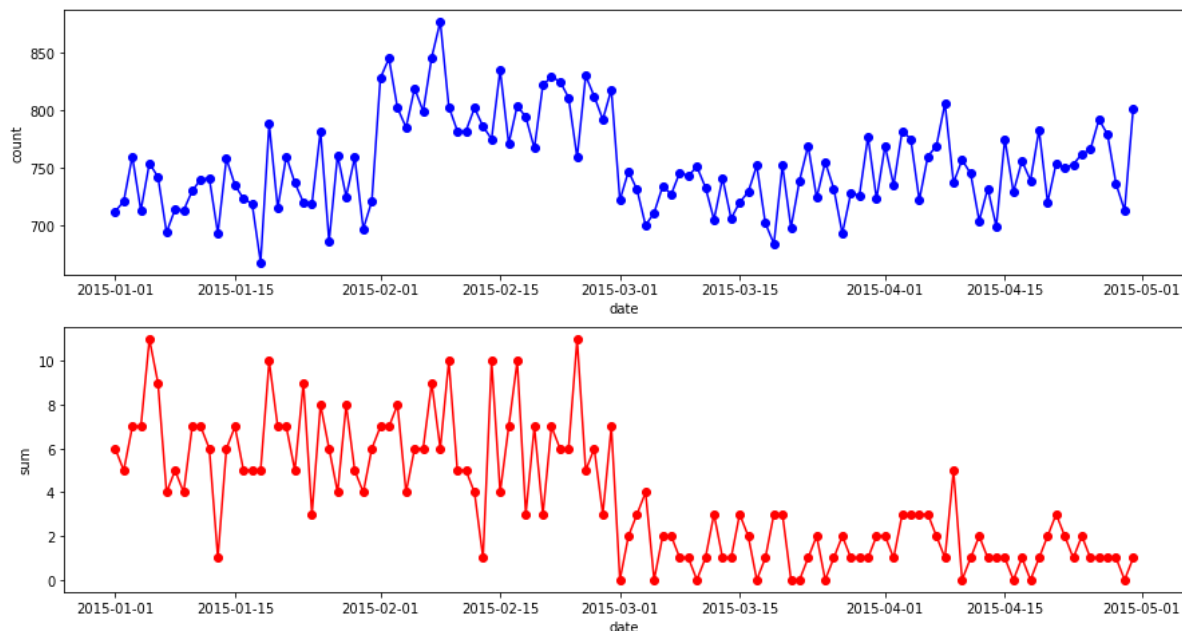
In [29]:

```

1 fig,ax=plt.subplots(2,1,figsize=(15,8))
2 ax[0].plot(df_purchase['date'],df_purchase['count'],color='b',linestyle='-',marker='o')
3 ax[0].set_ylabel('count')
4 ax[0].set_xlabel('date')
5 ax[1].plot(df_purchase['date'],df_purchase['sum'],color='r',linestyle='-',marker='o')
6 ax[1].set_ylabel('sum')
7 ax[1].set_xlabel('date')
8 fig.suptitle(f'Daily Payment Confirmations', fontsize=24)
9 plt.show()

```

## Daily Payment Confirmations



## Sex

In [56]:

```

1 ### group by sex
2 sex_purchases = df.groupby(by=['date','sex']).agg({'payment_confirmation':['sum']})
3 sex_purchases.columns = sex_purchases.columns.droplevel(level=1)
4 sex_purchases.reset_index(inplace=True)
5 sex_purchases.columns=['date','sex','sum','count']

```

In [57]:

```

1 sex_pivot = pd.pivot_table(sex_purchases, values=['sum'], columns=['sex'], index=['date'])

```

In [58]:

```

1 sex_pivot.columns=sex_pivot.columns.droplevel(level=0)

```

In [59]:

```

1 sex_pivot=sex_pivot.reset_index()

```

In [60]:

```
1 sex_pivot
```

Out[60]:

| sex | date                      | Female | Male |
|-----|---------------------------|--------|------|
| 0   | 2015-01-01 00:00:00+00:00 | 3      | 3    |
| 1   | 2015-01-02 00:00:00+00:00 | 2      | 3    |
| 2   | 2015-01-03 00:00:00+00:00 | 3      | 4    |
| 3   | 2015-01-04 00:00:00+00:00 | 2      | 5    |
| 4   | 2015-01-05 00:00:00+00:00 | 6      | 5    |
| ... | ...                       | ...    | ...  |
| 115 | 2015-04-26 00:00:00+00:00 | 1      | 0    |
| 116 | 2015-04-27 00:00:00+00:00 | 1      | 0    |
| 117 | 2015-04-28 00:00:00+00:00 | 1      | 0    |
| 118 | 2015-04-29 00:00:00+00:00 | 0      | 0    |
| 119 | 2015-04-30 00:00:00+00:00 | 0      | 1    |

120 rows × 3 columns

Devices

In [35]:

```

1 ### group by sex
2 device_purchases = df.groupby(by=['date', 'device']).agg({'payment_confirmation':
3 device_purchases

```

Out[35]:

|                           |         | payment_confirmation |       |
|---------------------------|---------|----------------------|-------|
|                           |         | sum                  | count |
| date                      | device  |                      |       |
| 2015-01-01 00:00:00+00:00 | Desktop | 1                    | 493   |
|                           | Mobile  | 5                    | 219   |
| 2015-01-02 00:00:00+00:00 | Desktop | 1                    | 484   |
|                           | Mobile  | 4                    | 237   |
| 2015-01-03 00:00:00+00:00 | Desktop | 3                    | 507   |
|                           | ...     | ...                  | ...   |
| 2015-04-28 00:00:00+00:00 | Mobile  | 1                    | 260   |
| 2015-04-29 00:00:00+00:00 | Desktop | 0                    | 453   |
|                           | Mobile  | 0                    | 260   |
| 2015-04-30 00:00:00+00:00 | Desktop | 0                    | 538   |
|                           | Mobile  | 1                    | 263   |

240 rows × 2 columns

In [36]:

```

1 device_purchases.reset_index(inplace=True)

```

In [37]:

```

1 device_purchases.columns=device_purchases.columns.droplevel(level=0)

```

In [38]:

```

1 device_purchases.columns=['date', 'device', 'sum', 'count']

```

In [39]:

```

1 device_pivot = pd.pivot_table(device_purchases, values=['count'], columns=['devi

```

In [40]:

```
1 device_pivot.reset_index()
```

Out[40]:

|        | date                      | count   |        |
|--------|---------------------------|---------|--------|
| device |                           | Desktop | Mobile |
| 0      | 2015-01-01 00:00:00+00:00 | 493     | 219    |
| 1      | 2015-01-02 00:00:00+00:00 | 484     | 237    |
| 2      | 2015-01-03 00:00:00+00:00 | 507     | 253    |
| 3      | 2015-01-04 00:00:00+00:00 | 474     | 239    |
| 4      | 2015-01-05 00:00:00+00:00 | 483     | 271    |
| ...    | ...                       | ...     | ...    |
| 115    | 2015-04-26 00:00:00+00:00 | 529     | 263    |
| 116    | 2015-04-27 00:00:00+00:00 | 509     | 270    |
| 117    | 2015-04-28 00:00:00+00:00 | 476     | 260    |

In [41]:

```
1 device_pivot.columns=device_pivot.columns.droplevel(level=0)
```

In [42]:

```
1 device_pivot=device_pivot.reset_index()
```

In [69]:

```
1 sex_pivot
```

Out[69]:

| sex | date                      | Female | Male |
|-----|---------------------------|--------|------|
| 0   | 2015-01-01 00:00:00+00:00 | 3      | 3    |
| 1   | 2015-01-02 00:00:00+00:00 | 2      | 3    |
| 2   | 2015-01-03 00:00:00+00:00 | 3      | 4    |
| 3   | 2015-01-04 00:00:00+00:00 | 2      | 5    |
| 4   | 2015-01-05 00:00:00+00:00 | 6      | 5    |
| ... | ...                       | ...    | ...  |
| 115 | 2015-04-26 00:00:00+00:00 | 1      | 0    |
| 116 | 2015-04-27 00:00:00+00:00 | 1      | 0    |
| 117 | 2015-04-28 00:00:00+00:00 | 1      | 0    |
| 118 | 2015-04-29 00:00:00+00:00 | 0      | 0    |
| 119 | 2015-04-30 00:00:00+00:00 | 0      | 1    |

120 rows × 3 columns

In [72]:

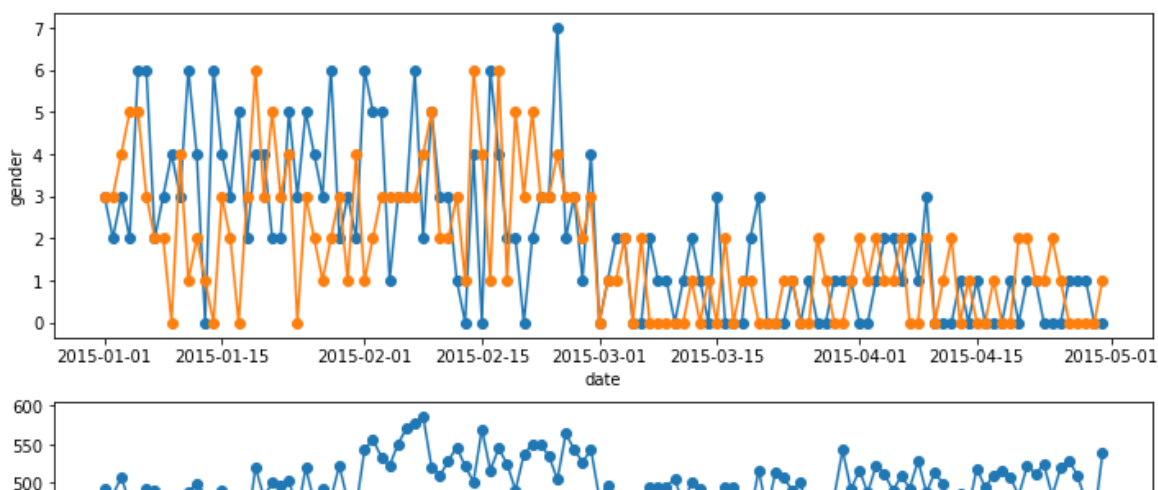
```

1 plt.set_cmap('jet')
2 fig,ax=plt.subplots(2,1,figsize=(12,8))
3 ax[0].plot(sex_pivot['date'],sex_pivot[['Female','Male']],marker='o')
4 ax[0].set_ylabel('gender')
5 ax[0].set_xlabel('date')
6 ax[1].plot(device_pivot['date'],device_pivot[['Desktop','Mobile']],linestyle='-',
7 ax[1].set_ylabel('device')
8 ax[1].set_xlabel('date')
9 fig.suptitle(f'Gender& Device Daily Payment Confirmations', fontsize=24)
10 plt.show()

```

&lt;Figure size 432x288 with 0 Axes&gt;

## Gender& Device Daily Payment Confirmations



## Initial A/B testing

### Resonse variable

- Conversion rate
  - time frame: one day
  - Definition: #of payment confirmation/ # of home page
  - importance: strong measure of growth

In [73]:

```

1 df_purchases=round(df_purchase['sum'].mean()*1000,2)
2 df_views=round(df_purchase['count'].mean()*1000,2)
3 print('daily purchases=', df_purchases)
4 print('daily views=', df_views)

```

daily purchases= 3766.67

daily views= 753333.33

In [74]:

```
1 df_purchase
```

Out[74]:

|     | date                      | sum | count |
|-----|---------------------------|-----|-------|
| 0   | 2015-01-01 00:00:00+00:00 | 6   | 712   |
| 1   | 2015-01-02 00:00:00+00:00 | 5   | 721   |
| 2   | 2015-01-03 00:00:00+00:00 | 7   | 760   |
| 3   | 2015-01-04 00:00:00+00:00 | 7   | 713   |
| 4   | 2015-01-05 00:00:00+00:00 | 11  | 754   |
| ... | ...                       | ... | ...   |
| 115 | 2015-04-26 00:00:00+00:00 | 1   | 792   |
| 116 | 2015-04-27 00:00:00+00:00 | 1   | 779   |
| 117 | 2015-04-28 00:00:00+00:00 | 1   | 736   |
| 118 | 2015-04-29 00:00:00+00:00 | 0   | 713   |
| 119 | 2015-04-30 00:00:00+00:00 | 1   | 801   |

120 rows × 3 columns

### Test sensitivity

In [206]:

```
1 df
```

Out[206]:

|       | user_id | date                      | device  | sex    | group   | home_page | search_page | payment_page | payment_conf |
|-------|---------|---------------------------|---------|--------|---------|-----------|-------------|--------------|--------------|
| 0     | 450007  | 2015-02-28 00:00:00+00:00 | Desktop | Female | test    | 1         | 0           | 0            |              |
| 1     | 756838  | 2015-01-13 00:00:00+00:00 | Desktop | Male   | control | 1         | 0           | 0            |              |
| 2     | 568983  | 2015-04-09 00:00:00+00:00 | Desktop | Male   | test    | 1         | 1           | 0            |              |
| 3     | 190794  | 2015-02-18 00:00:00+00:00 | Desktop | Female | test    | 1         | 1           | 0            |              |
| 4     | 537909  | 2015-01-15 00:00:00+00:00 | Desktop | Male   | test    | 1         | 0           | 0            |              |
| ...   | ...     | ...                       | ...     | ...    | ...     | ...       | ...         | ...          | ...          |
| 90395 | 307667  | 2015-03-30 00:00:00+00:00 | Desktop | Female | test    | 1         | 0           | 0            |              |

In [76]:

```
1 df['payment_confirmation'].value_counts()
```

Out[76]:

```
0    89948
```

```
1      452
```

Name: payment\_confirmation, dtype: int64

In [77]:

```
1 total_subs_count=np.sum(df['payment_confirmation'])
2 total_subs_count
```

Out[77]:

```
452
```

In [78]:

```
1 users_count=len(df['user_id'].unique())
2 users_count
```

Out[78]:

```
90400
```

In [89]:

```
1 conversion_rate=round(total_subs_count/users_count*100,2)
2 std=round(df['payment_confirmation'].std(),2)
3 print('conversion_rate:%s%%'%(conversion_rate))
4 print('std:',std)
```

conversion\_rate:0.5%

std: 0.07

In [100]:

```
1 ##small_sensitivity
2 small_sensitivity=0.1
3 small_conversion_rate=(conversion_rate/100)*(1+small_sensitivity)
4 small_purchases=df_views*small_conversion_rate
5 purchase_lift=small_purchases-df_purchases
6
7 print('small_conversion_rate: %s%%' %(small_conversion_rate*100))
8 print('small_purchases: %s' %(round(small_purchases,2)))
9 print('purchase_lift: %s' %(round(purchase_lift,2)))
```

small\_conversion\_rate: 0.55%

small\_purchases: 4143.33

purchase\_lift: 376.66

In [99]:

```

1  ##medium_sensitivity
2  medium_sensitivity=0.2
3  medium_conversion_rate=(conversion_rate/100)*(1+medium_sensitivity)
4  medium_purchases=df_views*medium_conversion_rate
5  purchase_lift=medium_purchases-df_purchases
6
7  print('medium_conversion_rate: %s%%' %(medium_conversion_rate*100))
8  print('medium_purchases: %s' %(round(medium_purchases,2)))
9  print('purchase_lift: %s' %(round(purchase_lift,2)))

```

medium\_conversion\_rate: 0.6%  
medium\_purchases: 4520.0  
purchase\_lift: 753.33

In [98]:

```

1  ##large_sensitivity
2  large_sensitivity=0.5
3  large_conversion_rate=(conversion_rate/100)*(1+large_sensitivity)
4  large_purchases=df_views*large_conversion_rate
5  purchase_lift=large_purchases-df_purchases
6
7  print('medium_conversion_rate: %s%%' %(large_conversion_rate*100))
8  print('medium_purchases: %s' %(round(large_purchases,2)))
9  print('purchase_lift: %s' %(round(purchase_lift,2)))

```

medium\_conversion\_rate: 0.75%  
medium\_purchases: 5650.0  
purchase\_lift: 1883.33

```

1  We choose the sensitivity to 0.6%

```

## Sample Size

In [101]:

```

1  purchase_mean=df['payment_confirmation'].mean()
2  purchase_std=df['payment_confirmation'].std()

```

Out[101]:

0.005

In [105]:

```

1  test_n=len(df[df['group']=='test'])
2  control_n=len(df[df['group']=='control'])
3  print(test_n,control_n)
4  sizes=[test_n,control_n]
5  ratio=max(sizes)/min(sizes)
6  print(ratio)

```

44810 45590  
1.0174068288328497



In [106]:

```
1  ## Sensitivity 0.6%, power=0.8, significance level 0.05
2  effect_size=(0.006-purchase_mean)/purchase_std
3  power=0.8
4  alpha=0.05
```

In [118]:

```
1  from statsmodels.stats import power as pwr
2  analysis=pwr.TTestIndPower()
3  ssresult=analysis.solve_power(effect_size=effect_size,
4                                power=power,
5                                nobsl=None,
6                                alpha=alpha,
7                                ratio=ratio)
8  print('Sample Size:',round(ssresult))
```

Sample Size: 77430

## Effect Size

In [119]:

```
1  sample_size=round(ssresult)
2  analysis=pwr.TTestIndPower()
3  esresult=analysis.solve_power(effect_size=None,
4                                power=power,
5                                nobsl=sample_size,
6                                alpha=alpha,
7                                ratio=ratio)
8  print('effect Size:',round(esresult,2))
```

effect Size: 0.01

## Statistical Power

In [122]:

```
1  effect_size=esresult
2  analysis=pwr.TTestIndPower()
3  pwresult=analysis.solve_power(effect_size=effect_size,
4                                power=None,
5                                nobsl=sample_size,
6                                alpha=alpha,
7                                ratio=ratio)
8  print('Power:',round(pwresult,3))
```

Power: 0.8

## Analyzing the A/B testing

In [137]:

```
1 results=df.groupby(by=['group']).agg({'user_id':pd.Series.nunique})
2 print(results)
```

```
      user_id
group
control    45590
test       44810
```

In [138]:

```
1 results=pd.DataFrame(results)
```

In [139]:

```
1 results['unique_users']=len(df.user_id.unique())
```

In [146]:

```
1 results['percentage']=round(results['user_id']/results['unique_users']*100,2)
```

In [147]:

```
1 results
```

Out[147]:

|         | user_id | unique_users | percentage |
|---------|---------|--------------|------------|
| group   |         |              |            |
| control | 45590   | 90400        | 50.43      |
| test    | 44810   | 90400        | 49.57      |

In [149]:

```
1 results2=df.groupby(by=['group','device','sex']).agg({'user_id':pd.Series.nunique})
2 results2=pd.DataFrame(results2)
3 results2['unique_users']=len(df.user_id.unique())
4 results2['percentage']=round(results2['user_id']/results2['unique_users']*100,2)
5 print(results2)
```

|         |         |        | user_id | unique_users | percentage |
|---------|---------|--------|---------|--------------|------------|
| group   | device  | sex    |         |              |            |
| control | Desktop | Female | 15227   | 90400        | 16.84      |
|         |         | Male   | 15146   | 90400        | 16.75      |
|         | Mobile  | Female | 7515    | 90400        | 8.31       |
|         |         | Male   | 7702    | 90400        | 8.52       |
| test    | Desktop | Female | 14770   | 90400        | 16.34      |
|         |         | Male   | 15057   | 90400        | 16.66      |
|         | Mobile  | Female | 7563    | 90400        | 8.37       |
|         |         | Male   | 7420    | 90400        | 8.21       |

In [163]:

```
1 test=df[df.group=='test']
2 control=df[df.group=='control']
```

In [164]:

```
1 print('test_size:',len(test))
2 print('con_size:',len(control))
```

```
test_size: 44810
con_size: 45590
```

In [165]:

```
1 test_mean=test['payment_confirmation'].mean()
2 cntrol_mean=control['payment_confirmation'].mean()
3 test_std=test['payment_confirmation'].std()
4 cntrol_std=control['payment_confirmation'].std()
```

In [168]:

```
1 print('test_conversion_rate:%s%%' %round(test_mean*100,2))
2 print('control_conversion_rate:%s%%' %round(cntrol_mean*100,2))
3 print('test_std:',round(test_std,2))
4 print('control_std:',round(cntrol_std,2))
```

```
test_conversion_rate:0.49%
control_conversion_rate:0.51%
test_std: 0.07
control_std: 0.07
```

In [169]:

```
1 test_result=df[df.group=='test']['payment_confirmation']
2 control_result=df[df.group=='control']['payment_confirmation']
3 ttest=stats.ttest_ind(test_result,control_result)
4 statistic=ttest[0]
5 p_value=ttest[1]
6 print('statistic',statistic)
7 print('p_value',p_value)
8 if p_value>=0.05:
9     print('Not significant')
10 else:
11     print('Significant!')
```

```
statistic -0.38195831967542987
p_value 0.7024931787844451
Not significant
```

In [ ]:

```
1
```