

```
#import packages
import pandas as pd
import seaborn as sns
import scipy.stats as ss
import statsmodels.api as sm
import matplotlib.pyplot as plt
pd.set_option('display.max_columns', None)
from matplotlib.pyplot import figure

In [2]:
# read data
T3 = pd.read_csv('t3_user_active_mins.csv')
T4 = pd.read_csv('t4_user_active_mins_pre.csv')
pd.set_option('display.max_columns', None)
from matplotlib.pyplot import figure

In [9]:
# look at data
T1.head(10)

Out [9]:
   uid      dt  active_mins
0  0  2019-02-22      5.0
1  0  2019-03-11      5.0
2  0  2019-03-18      3.0
3  0  2019-03-22      4.0
4  0  2019-04-03      9.0
5  0  2019-04-06      1.0
6  0  2019-04-17      1.0
7  0  2019-05-07      3.0
8  0  2019-05-14      1.0
9  0  2019-05-19      1.0

In [4]:
T2.head(10)

Out [4]:
   uid  variant_number      dt  signup_date
0  0  0  2019-02-06  2019-09-24
1  1  0  2019-02-06  2016-11-07
2  2  0  2019-02-06  2018-09-17
3  3  0  2019-02-06  2018-03-04
4  4  0  2019-02-06  2017-03-09
5  5  0  2019-02-06  2018-06-25
6  6  0  2019-02-06  2017-01-22
7  7  0  2019-02-06  2016-08-12
8  8  0  2019-02-06  2019-01-18
9  9  0  2019-02-06  2019-05-02

In [5]:
T3.head(10)

Out [9]:
   uid      dt  active_mins
0  0  2018-09-24      3.0
1  0  2018-11-08      4.0
2  0  2018-11-24      3.0
3  0  2018-11-28      6.0
4  0  2018-12-02      6.0
5  0  2018-12-04      1.0
6  0  2018-12-07      8.0
7  0  2018-12-09      8.0
8  0  2018-12-14      8.0
9  0  2018-12-15      2.0

In [6]:
T4.head(10)

Out [9]:
   uid  gender  user_type
0  0  male  non_reader
1  1  male  reader
2  2  male  non_reader
3  3  male  non_reader
4  4  male  non_reader
5  5  female  non_reader
6  6  female  non_reader
7  7  male  non_reader
8  8  male  new_user
9  9  female  non_reader

Hypothesis: Changes implemented for group B has increased user engagement
I will also break down whether
1) Active usership time increased or decreased
2) Frequency of visits increased or decreased
3) Do the above two points for old/new users, male/female users, and if they are readers/non-readers.

In [7]:
# see number of control group users and test group users.
T2['variant_number'].value_counts()
# 80% in control(40k), 20% in treatment/test group (10k)

Out [7]:
0    40000
Name: variant_number, dtype: int64

In [8]:
# load users into control and treatment/test group arrays.
user_control = []
user_treatment = []
for index, row in T2.iterrows():
    if row['variant_number'] == 0:
        user_control.append(row['uid'])
    else:
        user_treatment.append(row['uid'])

In [9]:
# merge attributes table and variant table.
var_gen_ageinqua = pd.merge(T2, T4, how = "inner", on = "uid")

In [10]:
var_gen_ageinqua.head()

Out [10]:
   uid  variant_number      dt  signup_date  gender  user_type
0  0  0  2019-02-06  2018-09-24  male  non_reader
1  1  0  2019-02-06  2016-11-07  male  reader
2  2  0  2019-02-06  2018-09-17  male  non_reader
3  3  0  2019-02-06  2018-03-04  male  non_reader
4  4  0  2019-02-06  2017-03-09  male  non_reader

In [11]:
# split var_gen_ageinqua into control and treatment groups.
var_gen_ageinqua_control = var_gen_ageinqua[var_gen_ageinqua['variant_number'] == 0]
var_gen_ageinqua_treatment = var_gen_ageinqua[var_gen_ageinqua['variant_number'] == 1]
cat_cols = ['gender', 'user_type']
cat_cols = ['days_since_joined_till_AB_tested']

In [12]:
# sort control and treatment groups into arrays based on if they are in treatment or control and date they joined before ab test.
uid_control = list(var_gen_ageinqua_control['uid'])
uid_treatment = list(var_gen_ageinqua_treatment['uid'])
before_AB_control = T3[T3['uid'].isin(uid_control)]
before_AB_treatment = T3[T3['uid'].isin(uid_treatment)]
after_AB_control = T1[T1['uid'].isin(uid_control)]
after_AB_treatment = T1[T1['uid'].isin(uid_treatment)]

In [13]:
# looks at created arrays
after_AB_control.head()

Out [13]:
   uid      dt  active_mins
88697  40000  2019-02-13      3.0
88698  40000  2019-03-02     18.0
88699  40000  2019-03-12      4.0
88690  40001  2019-02-14     16.0
88661  40001  2019-02-17      5.0

In [14]:
after_AB_control.head()

Out [14]:
   uid      dt  active_mins
0  0  2019-02-22      5.0
1  0  2019-03-11      5.0
2  0  2019-03-18      3.0
3  0  2019-03-22      4.0
4  0  2019-04-03      9.0

In [15]:
before_AB_treatment.head()

Out [15]:
   uid      dt  active_mins
88649  40001  2018-08-12      1.0
88650  40001  2018-08-21      1.0
88651  40001  2018-08-21      5.0
88652  40001  2018-09-23      3.0
88653  40001  2019-10-03      1.0

In [16]:
before_AB_control.head()

Out [16]:
   uid      dt  active_mins
0  0  2018-09-24      3.0
1  0  2018-11-08      4.0
2  0  2018-11-24      3.0
3  0  2018-11-28      6.0
4  0  2018-12-02      6.0

In [17]:
# merge the control/test and date tables with the attributes table
before_AB_control_base_table = pd.merge(before_AB_control, var_gen_ageinqua, how = "inner", on = "uid")
before_AB_treatment_base_table = pd.merge(before_AB_treatment, var_gen_ageinqua, how = "inner", on = "uid")
after_AB_control_base_table = pd.merge(after_AB_control, var_gen_ageinqua, how = "inner", on = "uid")
after_AB_treatment_base_table = pd.merge(after_AB_treatment, var_gen_ageinqua, how = "inner", on = "uid")

In [18]:
# renaming some columns
before_AB_control_base_table.rename(columns = {"dt_x":"date", "dt_y":"AB_Date", inplace = True)
before_AB_treatment_base_table.rename(columns = {"dt_x":"date", "dt_y":"AB_Date", inplace = True)
after_AB_control_base_table.rename(columns = {"dt_x":"date", "dt_y":"AB_Date", inplace = True)
after_AB_treatment_base_table.rename(columns = {"dt_x":"date", "dt_y":"AB_Date", inplace = True)

Out [18]:
   uid      date  active_mins  variant_number  AB_Date  signup_date  gender  user_type
0  0  2018-09-24      3.0      0  2019-02-06  2018-09-24  male  non_reader
1  0  2018-11-08      4.0      0  2019-02-06  2018-09-24  male  non_reader
2  0  2018-11-24      3.0      0  2019-02-06  2018-09-24  male  non_reader
3  0  2018-11-28      6.0      0  2019-02-06  2018-09-24  male  non_reader
4  0  2018-12-02      6.0      0  2019-02-06  2018-09-24  male  non_reader

In [19]:
before_AB_treatment_base_table_nout = before_AB_treatment_base_table[before_AB_treatment_base_table.active_mins < np.percentile(before_AB_treatment_base_table.active_mins, 90)]
after_AB_treatment_base_table_nout = after_AB_treatment_base_table[after_AB_treatment_base_table.active_mins < np.percentile(after_AB_treatment_base_table.active_mins, 90)]

In [20]:
before_AB_treatment_base_table_nout.head()

Out [20]:
   uid      date  active_mins  variant_number  AB_Date  signup_date  gender  user_type
0  40001  2018-06-12      1.0      1  2019-02-06  2017-04-29  male  non_reader
1  40001  2018-06-21      1.0      1  2019-02-06  2017-04-29  male  non_reader
2  40001  2018-09-21      3.0      1  2019-02-06  2017-04-29  male  non_reader
3  40001  2019-09-23      5.0      1  2019-02-06  2017-04-29  male  non_reader
4  40001  2019-10-03      1.0      1  2019-02-06  2017-04-29  male  non_reader

In [21]:
# plot of male readers for before (orange) and after (blue) data
f, axes = plt.subplots(1, 2)
plt.style.use('ggplot')
pd.set_option('display.max_rows', None)
sns.displot(after_AB_treatment_base_table_nout[(before_AB_treatment_base_table_nout["gender"] == "female")&(before_AB_treatment_base_table_nout["user_type"] == "reader", color="red", axlabel = "Active Mins, Female Readers", axes=0)])
sns.displot(after_AB_treatment_base_table_nout[(before_AB_treatment_base_table_nout["gender"] == "female")&(after_AB_treatment_base_table_nout["user_type"] == "reader", color="blue", axlabel = "Active Mins, Female Readers", axes=0)])
sns.displot(before_AB_treatment_base_table_nout[(before_AB_treatment_base_table_nout["gender"] == "female")&(before_AB_treatment_base_table_nout["user_type"] == "reader", color="red", axlabel = "Active Mins, Male Non-Readers", axes=1)])
sns.displot(before_AB_treatment_base_table_nout[(before_AB_treatment_base_table_nout["gender"] == "female")&(after_AB_treatment_base_table_nout["user_type"] == "reader", color="blue", axlabel = "Active Mins, Male Non-Readers", axes=1)])
plt.show()

C:\Users\vaandanaconda\venv\Data_Scili\lib\site-packages\seaborn\distributions.py:2557: FutureWarning: 'displot' is a deprecated function and will be removed in a future version. Please adapt your code to use either 'displot' (a figure-level function with similar flexibility) or 'kdeplot' (an axes-level function for kernel density plots).
warnings.warn(msg, FutureWarning)
C:\Users\vaandanaconda\venv\Data_Scili\lib\site-packages\seaborn\distributions.py:2557: FutureWarning: 'displot' is a deprecated function and will be removed in a future version. Please adapt your code to use either 'displot' (a figure-level function with similar flexibility) or 'kdeplot' (an axes-level function for kernel density plots).
warnings.warn(msg, FutureWarning)
C:\Users\vaandanaconda\venv\Data_Scili\lib\site-packages\seaborn\distributions.py:2557: FutureWarning: 'displot' is a deprecated function and will be removed in a future version. Please adapt your code to use either 'displot' (a figure-level function with similar flexibility) or 'kdeplot' (an axes-level function for kernel density plots).
warnings.warn(msg, FutureWarning)
C:\Users\vaandanaconda\venv\Data_Scili\lib\site-packages\seaborn\distributions.py:2557: FutureWarning: 'displot' is a deprecated function and will be removed in a future version. Please adapt your code to use either 'displot' (a figure-level function with similar flexibility) or 'kdeplot' (an axes-level function for kernel density plots).
warnings.warn(msg, FutureWarning)
C:\Users\vaandanaconda\venv\Data_Scili\lib\site-packages\seaborn\distributions.py:2557: FutureWarning: 'displot' is a deprecated function and will be removed in a future version. Please adapt your code to use either 'displot' (a figure-level function with similar flexibility) or 'kdeplot' (an axes-level function for kernel density plots).
warnings.warn(msg, FutureWarning)

In [22]:
# plot of male readers for before (orange) and after (blue) data
f, axes = plt.subplots(1, 2)
plt.style.use('ggplot')
pd.set_option('display.max_rows', None)
sns.displot(after_AB_treatment_base_table_nout[(before_AB_treatment_base_table_nout["gender"] == "male")&(before_AB_treatment_base_table_nout["user_type"] == "reader", color="red", axlabel = "Active Mins, Male Readers", axes=0)])
sns.displot(after_AB_treatment_base_table_nout[(before_AB_treatment_base_table_nout["gender"] == "male")&(after_AB_treatment_base_table_nout["user_type"] == "reader", color="blue", axlabel = "Active Mins, Male Readers", axes=0)])
sns.displot(before_AB_treatment_base_table_nout[(before_AB_treatment_base_table_nout["gender"] == "male")&(before_AB_treatment_base_table_nout["user_type"] == "reader", color="red", axlabel = "Active Mins, Female Non-Readers", axes=1)])
sns.displot(before_AB_treatment_base_table_nout[(before_AB_treatment_base_table_nout["gender"] == "male")&(after_AB_treatment_base_table_nout["user_type"] == "reader", color="blue", axlabel = "Active Mins, Female Non-Readers", axes=1)])
plt.show()

C:\Users\vaandanaconda\venv\Data_Scili\lib\site-packages\seaborn\distributions.py:2557: FutureWarning: 'displot' is a deprecated function and will be removed in a future version. Please adapt your code to use either 'displot' (a figure-level function with similar flexibility) or 'kdeplot' (an axes-level function for kernel density plots).
warnings.warn(msg, FutureWarning)
C:\Users\vaandanaconda\venv\Data_Scili\lib\site-packages\seaborn\distributions.py:2557: FutureWarning: 'displot' is a deprecated function and will be removed in a future version. Please adapt your code to use either 'displot' (a figure-level function with similar flexibility) or 'kdeplot' (an axes-level function for kernel density plots).
warnings.warn(msg, FutureWarning)
C:\Users\vaandanaconda\venv\Data_Scili\lib\site-packages\seaborn\distributions.py:2557: FutureWarning: 'displot' is a deprecated function and will be removed in a future version. Please adapt your code to use either 'displot' (a figure-level function with similar flexibility) or 'kdeplot' (an axes-level function for kernel density plots).
warnings.warn(msg, FutureWarning)
C:\Users\vaandanaconda\venv>Data_Scili\lib\site-packages\seaborn\distributions.py:2557: FutureWarning: 'displot' is a deprecated function and will be removed in a future version. Please adapt your code to use either 'displot' (a figure-level function with similar flexibility) or 'kdeplot' (an axes-level function for kernel density plots).
warnings.warn(msg, FutureWarning)

In [23]:
user_count_before_AB = pd.DataFrame(before_AB_treatment_base_table_nout['uid'].value_counts())
user_count_before_AB.reset_index(inplace = True)
user_count_after_AB = pd.DataFrame(after_AB_treatment_base_table_nout['uid'].value_counts())
user_count_after_AB.reset_index(inplace = True)
user_count_after_AB.rename(columns = {"index":"before_count", "index":"uid", inplace = True)
user_count_after_AB.rename(columns = {"index":"after_count", "index":"uid", inplace = True)

In [24]:
fig = plt.figure(figsize=(10,6))
plt.style.use('ggplot')
pd.set_option('display.max_rows', None)
sns.displot(user_count_before_AB['before_count'], hist=False, color='red', axlabel = "Frequency of Visit")
sns.displot(user_count_after_AB['after_count'], hist=False, color='blue', axlabel = "Frequency of Visit")
plt.legend()

C:\Users\vaandanaconda\venv\Data_Scili\lib\site-packages\seaborn\distributions.py:2557: FutureWarning: 'displot' is a deprecated function and will be removed in a future version. Please adapt your code to use either 'displot' (a figure-level function with similar flexibility) or 'kdeplot' (an axes-level function for kernel density plots).
warnings.warn(msg, FutureWarning)
C:\Users\vaandanaconda\venv>Data_Scili\lib\site-packages\seaborn\distributions.py:2557: FutureWarning: 'displot' is a deprecated function and will be removed in a future version. Please adapt your code to use either 'displot' (a figure-level function with similar flexibility) or 'kdeplot' (an axes-level function for kernel density plots).
warnings.warn(msg, FutureWarning)
C:\Users\vaandanaconda\venv>Data_Scili\lib\site-packages\seaborn\distributions.py:2557: FutureWarning: 'displot' is a deprecated function and will be removed in a future version. Please adapt your code to use either 'displot' (a figure-level function with similar flexibility) or 'kdeplot' (an axes-level function for kernel density plots).
warnings.warn(msg, FutureWarning)
C:\Users\vaandanaconda\venv>Data_Scili\lib\site-packages\seaborn\distributions.py:2557: FutureWarning: 'displot' is a deprecated function and will be removed in a future version. Please adapt your code to use either 'displot' (a figure-level function with similar flexibility) or 'kdeplot' (an axes-level function for kernel density plots).
warnings.warn(msg, FutureWarning)

In [25]:
after_AB_control_base_table_nout = after_AB_control_base_table[after_AB_control_base_table.active_mins < np.percentile(after_AB_control_base_table.active_mins, 90)]
after_AB_treatment_base_table_nout = after_AB_treatment_base_table[after_AB_treatment_base_table.active_mins < np.percentile(after_AB_treatment_base_table.active_mins, 90)]

In [26]:
after_AB_treatment_base_table_nout.head()

Out [26]:
   uid      date  active_mins  variant_number  AB_Date  signup_date  gender  user_type
0  40001  2019-02-13      3.0      1  2019-02-06  2019-02-04  male  new_user
1  40000  2019-03-02     18.0      1  2019-02-06  2019-02-04  male  new_user
2  40000  2019-03-12      4.0      1  2019-02-06  2019-02-04  male  new_user
3  40001  2019-02-14     16.0      1  2019-02-06  2017-04-29  male  non_reader
4  40001  2019-02-17      5.0      1  2019-02-06  2017-04-29  male  non_reader

In [27]:
after_AB_control_base_table_nout.head()

Out [27]:
   uid      date  active_mins  variant_number  AB_Date  signup_date  gender  user_type
0  0  2019-02-22      5.0      0  2019-02-06  2018-09-24  male  non_reader
1  0  2019-03-11      5.0      0  2019-02-06  2018-09-24  male  non_reader
2  0  2019-03-18      3.0      0  2019-02-06  2018-09-24  male  non_reader
3  0  2019-03-22      4.0      0  2019-02-06  2018-09-24  male  non_reader
4  0  2019-04-03      9.0      0  2019-02-06  2018-09-24  male  non_reader

In [28]:
# plot of male/female readers for control and treatment data after the ab date
f, axes = plt.subplots(1, 2)
plt.style.use('ggplot')
pd.set_option('display.max_rows', None)
sns.displot(after_AB_control_base_table_nout[(after_AB_control_base_table_nout["gender"] == "female")&(after_AB_control_base_table_nout["user_type"] == "reader", color="red", axlabel = "Active Mins, Male Readers", axes=0)])
sns.displot(after_AB_control_base_table_nout[(after_AB_control_base_table_nout["gender"] == "female")&(after_AB_control_base_table_nout["user_type"] == "reader", color="blue", axlabel = "Active Mins, Male Readers", axes=0)])
sns.displot(before_AB_control_base_table_nout[(before_AB_control_base_table_nout["gender"] == "female")&(before_AB_control_base_table_nout["user_type"] == "reader", color="red", axlabel = "Active Mins, Female Non-Readers", axes=1)])
sns.displot(before_AB_control_base_table_nout[(before_AB_control_base_table_nout["gender"] == "female")&(after_AB_control_base_table_nout["user_type"] == "reader", color="blue", axlabel = "Active Mins, Female Non-Readers", axes=1)])
plt.show()

C:\Users\vaandanaconda\venv\Data_Scili\lib\site-packages\seaborn\distributions.py:2557: FutureWarning: 'displot' is a deprecated function and will be removed in a future version. Please adapt your code to use either 'displot' (a figure-level function with similar flexibility) or 'kdeplot' (an axes-level function for kernel density plots).
warnings.warn(msg, FutureWarning)
C:\Users\vaandanaconda\venv>Data_Scili\lib\site-packages\seaborn\distributions.py:2557: FutureWarning: 'displot' is a deprecated function and will be removed in a future version. Please adapt your code to use either 'displot' (a figure-level function with similar flexibility) or 'kdeplot' (an axes-level function for kernel density plots).
warnings.warn(msg, FutureWarning)
C:\Users\vaandanaconda\venv>Data_Scili\lib\site-packages\seaborn\distributions.py:2557: FutureWarning: 'displot' is a deprecated function and will be removed in a future version. Please adapt your code to use either 'displot' (a figure-level function with similar flexibility) or 'kdeplot' (an axes-level function for kernel density plots).
warnings.warn(msg, FutureWarning)
C:\Users\vaandanaconda\venv>Data_Scili\lib\site-packages\seaborn\distributions.py:2557: FutureWarning: 'displot' is a deprecated function and will be removed in a future version. Please adapt your code to use either 'displot' (a figure-level function with similar flexibility) or 'kdeplot' (an axes-level function for kernel density plots).
warnings.warn(msg, FutureWarning)
C:\Users\vaandanaconda\venv>Data_Scili\lib\site-packages\seaborn\distributions.py:2557: FutureWarning: 'displot' is a deprecated function and will be removed in a future version. Please adapt your code to use either 'displot' (a figure-level function with similar flexibility) or 'kdeplot' (an axes-level function for kernel density plots).
warnings.warn(msg, FutureWarning)

In [29]:
# plot of male/female readers for control and treatment data after the ab date
f, axes = plt.subplots(1, 2)
plt.style.use('ggplot')
pd.set_option('display.max_rows', None)
sns.displot(after_AB_control_base_table_nout[(after_AB_control_base_table_nout["gender"] == "female")&(after_AB_control_base_table_nout["user_type"] == "reader", color="red", axlabel = "Active Mins, Male Readers", axes=0)])
sns.displot(after_AB_control_base_table_nout[(after_AB_control_base_table_nout["gender"] == "female")&(after_AB_control_base_table_nout["user_type"] == "reader", color="blue", axlabel = "Active Mins, Male Readers", axes=0)])
sns.displot(before_AB_control_base_table_nout[(before_AB_control_base_table_nout["gender"] == "female")&(before_AB_control_base_table_nout["user_type"] == "reader", color="red", axlabel = "Active Mins, Female Non-Readers", axes=1)])
sns.displot(before_AB_control_base_table_nout[(before_AB_control_base_table_nout["gender"] == "female")&(after_AB_control_base_table_nout["user_type"] == "reader", color="blue", axlabel = "Active Mins, Female Non-Readers", axes=1)])
plt.show()

C:\Users\vaandanaconda\venv>Data_Scili\lib\site-packages\seaborn\distributions.py:2557: FutureWarning: 'displot' is a deprecated function and will be removed in a future version. Please adapt your code to use either 'displot' (a figure-level function with similar flexibility) or 'kdeplot' (an axes-level function for kernel density plots).
warnings.warn(msg, FutureWarning)
C:\Users\vaandanaconda\venv>Data_Scili\lib\site-packages\seaborn\distributions.py:2557: FutureWarning: 'displot' is a deprecated function and will be removed in a future version. Please adapt your code to use either 'displot' (a figure-level function with similar flexibility) or 'kdeplot' (an axes-level function for kernel density plots).
warnings.warn(msg, FutureWarning)
C:\Users\vaandanaconda\venv>Data_Scili\lib\site-packages\seaborn\distributions.py:2557: FutureWarning: 'displot' is a deprecated function and will be removed in a future version. Please adapt your code to use either 'displot' (a figure-level function with similar flexibility) or 'kdeplot' (an axes-level function for kernel density plots).
warnings.warn(msg, FutureWarning)
C:\Users\vaandanaconda\venv>Data_Scili\lib\site-packages\seaborn\distributions.py:2557: FutureWarning: 'displot' is a deprecated function and will be removed in a future version. Please adapt your code to use either 'displot' (a figure-level function with similar flexibility) or 'kdeplot' (an axes-level function for kernel density plots).
warnings.warn(msg, FutureWarning)

In [30]:
# creating dummy variables for gender
col_to_drop = ['date', 'AB_Date', 'signup_date']

def mapMaleFemale(df):
    return df['gender'].map({'female':0, 'male':1})

In [31]:
# add gender to arrays of attributes and active minutes for control and treatment
before_AB_control_base_table['gender'] = mapMaleFemale(before_AB_control_base_table)
after_AB_control_base_table['gender'] = mapMaleFemale(after_AB_control_base_table)
after_AB_treatment_base_table['gender'] = mapMaleFemale(after_AB_treatment_base_table)
after_AB_control_base_table_nout['gender'] = mapMaleFemale(after_AB_control_base_table_nout)
after_AB_treatment_base_table_nout['gender'] = mapMaleFemale(after_AB_treatment_base_table_nout)

Out [32]:
   uid      date  active_mins  variant_number  AB_Date  signup_date  gender  user_type
0  0  2019-02-22      5.0      0  2019-02-06  2018-09-24  1.0  non_reader
1  0  2019-03-11      5.0      0  2019-02-06  2018-09-24  1.0  non_reader
2  0  2019-03-18      3.0      0  2019-02-06  2018-09-24  1.0  non_reader
3  0  2019-03-22      4.0      0  2019-02-06  2018-09-24  1.0  non_reader
4  0  2019-04-03      9.0      0  2019-02-06  2018-09-24  1.0  non_reader

In [33]:
total_control = pd.merge(before_AB_control_base_table, before_AB_control_base_table, how = "inner", on = "uid")
total_treatment = pd.merge(before_AB_treatment_base_table, after_AB_treatment_base_table, how = "inner", on = "uid")

In [34]:
# Function to order tables how I would like them for my analysis
def finalizeDataTables(df, col_to_drop, generateFig):
    df_dum_1 = pd.get_dummies(df, columns=col_to_drop, prefix='user')
    df_dum_2 = pd.get_dummies(df, columns=col_to_drop, prefix='user')
    df = pd.concat([df_dum_1, df_dum_2], axis=1)
    df.reset_index(inplace = True)
    a = pd.DataFrame(df['user_type'].value_counts())
    a.reset_index(inplace = True)
    df = pd.merge(df, a, how = "inner", on = "user_type")
    df.to_groupby = list(df.columns)
    col_to_group_by = remove('active_mins')
    df = pd.DataFrame(df.groupby(col_to_group_by)['active_mins'].sum())
    df.reset_index(inplace = True)
    df.drop(columns = ['user_type'], inplace = True)
    df['avg_time_per_user'] = df['active_mins']/df['appearances']
    if generateFig:
        df.to_csv(df['user_type'].index + '.csv', index = False)
    return df

In [35]:
# passing tables through above function
before_AB_control_base_table = finalizeDataTables(before_AB_control_base_table, col_to_drop, False)
after_AB_control_base_table = finalizeDataTables(after_AB_control_base_table, col_to_drop, False)
before_AB_treatment_base_table = finalizeDataTables(before_AB_treatment_base_table, col_to_drop, False)
after_AB_treatment_base_table = finalizeDataTables(after_AB_treatment_base_table, col_to_drop, False)

In [36]:
before_AB_control_base_table.describe()

Out [36]:
   uid  variant_number  gender  contributor  new_user  non_reader  reader  appearances  active_mins  avg_time_per_user
mean  33655.000000  33655.0  33655.000000  33655.000000  33655.000000  33655.000000  33655.000000  33655.000000  33655.000000  33655.000000
min  11543.409661  0.0  0.000000  0.000000  0.000000  0.000000  0.000000  0.000000  1.000000  1.000000e+00
25%  15997.500000  0.0  0.000000  0.000000  0.000000  0.000000  0.000000  6.000000  15.000000  2.616667
50%  19967.000000  0.0  0.000000  0.000000  0.000000  0.000000  0.000000  10.000000  25.000000  4.000000
75%  30011.500000  0.0  0.000000  0.000000  0.000000  0.000000  0.000000  15.000000  35.000000  6.388889
max  39999.000000  0.0  0.000000  0.000000  0.000000  0.000000  0.000000  33.000000  33996.000000  70.971667

In [37]:
# filter outliers
before_AB_control_base_table = before_AB_control_base_table[(np.abs(ss.zscore(before_AB_control_base_table['avg_time_per_user'])) < 4)]
after_AB_control_base_table = before_AB_control_base_table[(np.abs(ss.zscore(after_AB_control_base_table['avg_time_per_user'])) < 4)]
before_AB_treatment_base_table = before_AB_treatment_base_table[(np.abs(ss.zscore(after_AB_treatment_base_table['avg_time_per_user'])) < 4)]
after_AB_treatment_base_table = before_AB_treatment_base_table[(np.abs(ss.zscore(after_AB_treatment_base_table['avg_time_per_user'])) < 4)]

Out [37]:
   uid  variant_number  gender  contributor  new_user  non_reader  reader  appearances  active_mins  avg_time_per_user
count  33605.000000  33605.0  33605.000000  33605.000000  33605.000000  33605.000000  33605.000000  33605.000000  33605.000000  33605.000000
mean  19986.213002  0.0  0.657621  0.023300  0.003330  0.003330  0.718947  0.176981  25.493736  466.116202
min  11543.757525  0.0  0.000000  0.000000  0.000000  0.000000  0.450489  0.383139  17.233079  1588.811397
25%  15997.000000  0.0  0.000000  0.000000  0.000000  0.000000  0.000000  6.000000  15.000000  2.616667
50%  19963.000000  0.0  0.000000  0.000000  0.000000  0.000000  0.000000  10.000000  25.000000  4.000000
75%  30011.000000  0.0  0.000000  0.000000  0.000000  0.000000  0.000000  15.000000  35.000000  6.388889
max  39999.000000  0.0  0.000000  0.000000  0.000000  0.000000  0.000000  33.000000  33996.000000  70.971667

In [38]:
after_AB_treatment_base_table.head()

Out [38]:
   uid  variant_number  gender  contributor  new_user  non_reader  reader  appearances  active_mins  avg_time_per_user
0  40001  1  1.0  0  0  1  0  0  32  29.0  0.933333
1  40000  1  1.0  0  0  1  0  3  25.0  8.666667
2  40000  1  1.0  0  0  1  0  26  183.0  7.039602
3  40004  1  0.0  0  0  1  0  9  56.0  6.222222
4  40005  1  0.0  0  0  1  0  36  289.0  8.027778

In [39]:
after_AB_control_base_table.describe()

Out [39]:
   uid  variant_number  gender  contributor  new_user  non_reader  reader  appearances  active_mins  avg_time_per_user
count  33605.000000  33605.0  33605.000000  33605.000000  33605.000000  33605.000000  33605.000000  33605.000000  33605.000000  33605.000000
mean  19986.213002  0.0  0.657621  0.023300  0.003330  0.003330  0.718947  0.176981  25.493736  466.116202
min  11543.757525  0.0  0.000000  0.000000  0.000000  0.000000  0.450489  0.383139  17.233079  1588.811397
25%  15997.000000  0.0  0.000000  0.000000  0.000000  0.000000  0.000000  6.000000  15.000000  2.616667
50%  19963.000000  0.0  0.000000  0.000000  0.000000  0.000000  0.000000  10.000000  25.000000  4.000000
75%  30011.000000  0.0  0.000000  0.000000  0.000000  0.000000  0.000000  15.000000  35.000000  6.388889
max  39999.000000  0.0  0.000000  0.000000  0.000000  0.000000  0.000000  33.000000  33996.000000  70.971667

In [40]:
after_AB_control_base_table.describe()

Out [40]:
   uid  variant_number  gender  contributor  new_user  non_reader  reader  appearances  active_mins  avg_time_per_user
count  31700.000000  31700.0  31700.000000  31700.000000  31700.000000  31700.000000  31700.000000  31700.000000  31700.000000  31700.000000
mean  19989.442732  0.0  0.661640  0.024983  0.009779  0.009779  0.186804  25.493736  466.116202  8.666667
min  11543.409661  0.0  0.000000  0.000000  0.000000  0.000000  0.000000  6.000000  15.000000  1.000000
25%  15997.000000  0.0  0.000000  0.000000  0.000000  0.000000  0.000000  10.000000  25.000000  2.500000
50%  19963.000000  0.0  0.000000  0.000000  0.000000  0.000000  0.000000  15.000000  35.000000  3.750000
75%  29845.250000  0.0  0.000000  0.000000  0.000000  0.000000  0.000000  34.000000  50.000000  7.693702
max  39999.000000  0.0  0.000000  0.000000  0.000000  0.000000  0.000000  110.000000  26887.000000  265.967
```