

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

from datetime import datetime
# Display figures inline in Jupyter notebook
import matplotlib.pyplot as plt
import seaborn as sns
# Use seaborn style defaults and set the default figure size
sns.set(rc={'figure.figsize':(11, 4)})
```

```
In [2]: # !pip install mxlend
from mxlend.frequent_patterns import apriori, association_rules
```

Reading the data files

```
In [3]: event=pd.read_csv('events.csv')
```

```
In [4]: events.head(5)
```

	timestamp	visitorid	event	itemid	transactionid
0	1433221322117	257597	view	355908	NaN
1	1433224214164	992329	view	248676	NaN
2	1433221998927	111016	view	318965	NaN
3	1433221955914	483717	view	253185	NaN
4	1433221337106	951259	view	367447	NaN

```
In [5]: category_tree=pd.read_csv('category_tree.csv')
```

```
In [6]: category_tree.head(5)
```

	categoryid	parentid
0	1016	2130
1	809	1690
2	570	90
3	1691	8850
4	536	16910

```
In [7]: category_tree.shape
```

```
Out[7]: (1669, 2)
```

```
In [8]: category_tree.parentid.nunique()
```

```
Out[8]: 362
```

```
In [9]: item_prop=pd.read_csv('item_properties_part1.csv')
item_prop2=pd.read_csv('item_properties_part2.csv')
```

```
In [10]: item_props=pd.concat([item_prop,item_prop2],ignore_index=True)
```

```
In [11]: item_props.head(5)
```

	timestamp	itemid	property	value
0	1435460400000	460429	categoryid	1338
1	1441508400000	206783	888	1116713 960601 n277200
2	1439089203000	395014	400	n552.000 639502 n700.000 424566
3	1431226803000	59481	790	n15360.000
4	1431831600000	156781	917	828513

```
In [12]: item_props.dtypes
```

```
Out[12]: timestamp    int64
itemid            int64
property         object
value            object
dtype: object
```

```
In [13]: item_props=item_props.sort_values(by='timestamp')
```

Change timestamp to date-time format i tried to convert it directly but it gave me DateTime in the future so i divided the value by 1000 before converting it to be reasonable. maybe it's wrong but i did this just to be readable

```
In [14]: item_props['timestamp']=item_props['timestamp'].apply( lambda x: datetime.fromtimestamp(x/1000))
```

```
In [15]: item_props.shape
```

```
Out[15]: (20275902, 4)
```

```
In [16]: item_props.itemid.nunique()
```

```
Out[16]: 417053
```

```
In [17]: item_props
```

	timestamp	itemid	property	value
5903679	2015-05-10 05:00:00	317951	790	n32880.000
5668465	2015-05-10 05:00:00	422842	480	1133979
11314219	2015-05-10 05:00:00	310185	776	103591
15170323	2015-05-10 05:00:00	110973	112	679677
15170323	2015-05-10 05:00:00	179597	available	0
...
9472889	2015-09-13 05:00:00	364708	928	769062
9473006	2015-09-13 05:00:00	231604	888	561561 1055803 447378 n12.000 1135780 1284577 ..
2196967	2015-09-13 05:00:00	161357	888	12762 16970 145048 237874 1229126 784581 12977 ..
9472887	2015-09-13 05:00:00	267142	available	0
20275901	2015-09-13 05:00:00	275768	888	888666 n10800.000 746840 1318567

20275902 rows x 4 columns

Exploring properties of items over time

```
In [18]: item_props=item_props.set_index('timestamp')
```

```
In [19]: item_props=item_props.sort_values(by=['timestamp', 'itemid','property', 'value'])
```

```
In [20]: item_props
```

	timestamp	itemid	property	value
2015-05-10 05:00:00	0	159		519769
2015-05-10 05:00:00	0	283		66094 372274 478989
2015-05-10 05:00:00	0	6		1152934 1238769
2015-05-10 05:00:00	0	678		372274
2015-05-10 05:00:00	0	790		n91200.000
...
2015-09-13 05:00:00	466864	790		n111840.000
2015-09-13 05:00:00	466864	813		1148082 353870 1262739
2015-09-13 05:00:00	466864	888		1262739 205682 1050016 1154859
2015-09-13 05:00:00	466864	available		0
2015-09-13 05:00:00	466865	888	150169 780351 820477 437265 951705 103274 1154..	

20275902 rows x 3 columns

Most frequently changed properties

from the below code we find properties **888,790,available,categoryid,451** the most frequently changed properties

```
In [21]: # item_props[item_props.duplicated(['itemid','property'])].property.value_counts()
```

```
In [22]: sn.barplot(x=item_props.property.value_counts().index[0:10], y=item_props.property.value_counts()[0:10])
```

```
Out[22]: <AxesSubplot:ylabel='property'>
```



get the changed items over time

```
In [23]: changed_items=item_props[item_props.duplicated(['itemid','property'])]
changed_items
```

	timestamp	itemid	property	value
2015-05-17 05:00:00	0	6		1152934 1238769
2015-05-17 05:00:00	1	790		n5760.000
2015-05-17 05:00:00	1	888		172646
2015-05-17 05:00:00	1	available		0
2015-05-17 05:00:00	3	283		138228 150169 1182824 327918 261419
...
2015-09-13 05:00:00	466864	790		n111840.000
2015-09-13 05:00:00	466864	813		1148082 353870 1262739
2015-09-13 05:00:00	466864	888		1262739 205682 1050016 1154859
2015-09-13 05:00:00	466864	available		0
2015-09-13 05:00:00	466865	888	150169 780351 820477 437265 951705 103274 1154..	

8272088 rows x 3 columns

listing the items that have different category over time

```
In [24]: changed_items[changed_items['property']!='categoryid']
```

```
Out[24]:
```

	timestamp	itemid	property	value
2015-05-17 05:00:00	25	categoryid	1509	
2015-05-17 05:00:00	94	categoryid	1147	
2015-05-17 05:00:00	130	categoryid	1277	
2015-05-17 05:00:00	149	categoryid	1120	
2015-05-17 05:00:00	168	categoryid	1006	
...
2015-09-13 05:00:00	466738	categoryid	1617	
2015-09-13 05:00:00	466767	categoryid	1114	
2015-09-13 05:00:00	466768	categoryid	1244	
2015-09-13 05:00:00	466783	categoryid	1114	
2015-09-13 05:00:00	466829	categoryid	438	

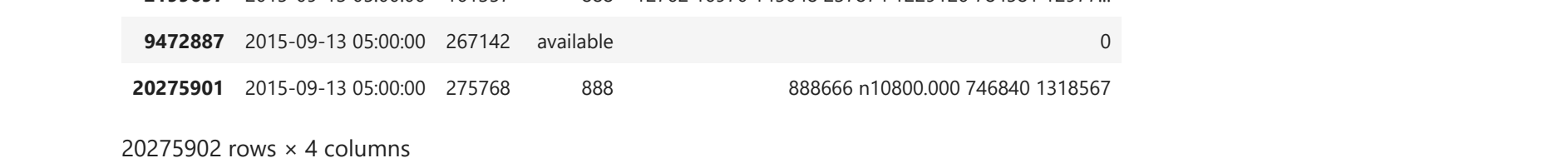
371161 rows x 3 columns

lets check frist item **393623** and property **categoryid**

```
In [25]: # item_props[item_props['itemid']==369722] & (item_props['property']!='categoryid')]['value']
```

```
In [26]: item_props[item_props['itemid']==393623] & (item_props['property']!='categoryid')]['value'].astype(int).plot(figsize=(10,5))
```

```
Out[26]: <AxesSubplot:xlabel='timestamp'>
```



the changes in items avaiability over time

```
In [27]: changed_items[changed_items['property']!='available']
```

```
Out[27]:
```

	timestamp	itemid	property	value
2015-05-17 05:00:00	1	available	0	
2015-05-17 05:00:00	6	available	1	
2015-05-17 05:00:00	15	available	1	
2015-05-17 05:00:00	16	available	0	
2015-05-17 05:00:00	19	available	1	
...
2015-09-13 05:00:00	466848	available	1	
2015-09-13 05:00:00	466853	available	0	
2015-09-13 05:00:00	466858	available	1	
2015-09-13 05:00:00	466861	available	1	
2015-09-13 05:00:00	466864	available	0	

1086586 rows x 3 columns

check availability of item 267142 and 379701 over time

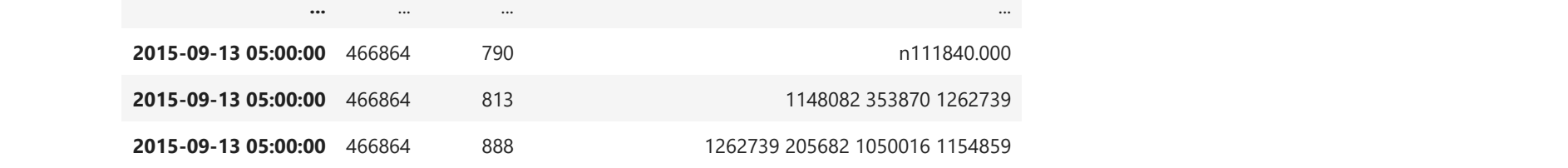
```
In [28]: item_props[item_props['itemid']==267142] & (item_props['property']!='available')]['value'].astype(int).plot(figsize=(10,5))
```

```
Out[28]: <AxesSubplot:xlabel='timestamp'>
```



```
In [29]: item_props[item_props['itemid']==379701] & (item_props['property']!='available')]['value'].astype(int).plot(figsize=(10,5))
```

```
Out[29]: <AxesSubplot:xlabel='timestamp'>
```



```
In [30]: # item_props[item_props['itemid']==267142] & (item_props['property']!='available')]['value']
```

item 267142 almost not available most of time it was available only once at 2015-08-02 05:00:00 in contrast to item 379701 which was available most of time

Num of unique properties per item

```
In [31]: propertiesperitem=item_props.groupby(['itemid'])['property'].nunique()
```

```
In [32]: propertiesperitem=propertiesperitem.to_frame( name='nuniqueproperty')
```

```
In [33]: propertiesperitem
```

	nuniqueproperty
0	28
1	35
2	24
3	29
4	25
...	...
466862	31
466863	23
466864	28
466865	27
466866	42

417053 rows x 1 columns

from the below output we find that the average number of properties per item amlost **28** and max is **59** and min **12**

```
In [34]: propertiesperitem.describe()
```

	nuniqueproperty
count	417053.000000
mean	28.782466
std	7.409326
min	12.000000
25%	24.000000
50%	27.000000
75%	31.000000
max	59.000000

```
In [35]: # item_props.groupby(['itemid','property'])['value'].nunique().to_frame(name='numofchanges').sort_values(by='numofchanges')
```

Properties that are constant over time

```
In [36]: valuesperproperty=item_props.groupby(['itemid','property'],as_index=False).agg(['value':['count','nunique']])
```

```
In [37]: # valuesperproperty=valuesperproperty.reset_index()
```

```
In [38]: valuesperproperty.columns = ['itemid','property','value_nunique','value_count']
```

```
In [39]: valuesperproperty[(valuesperproperty['value_nunique']==1)&(valuesperproperty['value_count']>1)]
```

```
Out[39]:
```

	itemid	property	value_nunique	value_count
--	--------	----------	---------------	-------------

Validating result by explore example of returned rows

- we can see from the below output /graph that the value is the same over time

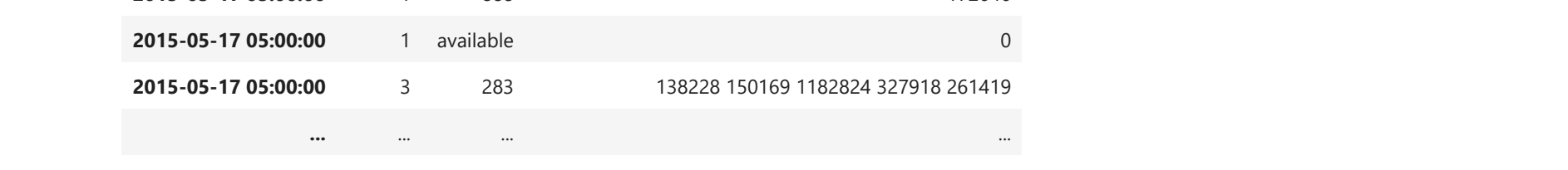
```
In [40]: item_props[item_props.itemid==40] & (item_props.property=='810')
```

```
Out[40]:
```

	timestamp	itemid	property	value
2015-05-24 05:00:00	40	810	n180.000 424566	
2015-05-31 05:00:00	40	810	n180.000 424566	
2015-06-07 05:00:00	40	810	n180.000 424566	
2015-06-14 05:00:00	40	810	n180.000 424566	
2015-06-28 05:00:00	40	810	n180.000 424566	
2015-07-05 05:00:00	40	810	n180.000 424566	
2015-07-12 05:00:00	40	810	n180.000 424566	
2015-07-19 05:00:00	40	810	n180.000 424566	
2015-08-02 05:00:00	40	810	n180.000 424566	
2015-08-09 05:00:00	40	810	n180.000 424566	
2015-08-16 05:00:00	40	810	n180.000 424566	
2015-08-23 05:00:00	40	810	n180.000 424566	
2015-08-30 05:00:00	40	810	n180.000 424566	
2015-09-06 05:00:00	40	810	n180.000 424566	
2015-09-13 05:00:00	40	810	n180.000 424566	

```
In [41]: item_props[item_props.itemid==40] & (item_props.property=='810')]['value'].value_counts().plot.bar(stacked=True)
```

```
Out[41]: <AxesSubplot>
```



- from the below code we can find that there are **1104** unique property that are constant over time for some items
- also we see that the count for some property is the same of the number of unique items so i think this properties may be the meta data about every item like **description** or **registerdate** because every item have it
- the top of this properties are **764,159,790,364** and **283**

```
In [42]: valuesperproperty.property.nunique()
```

```
Out[42]: 1104
```

```
In [43]: valuesperproperty.property.value_counts()
```

	value
764	417053
159	417053
790	417053
364	417053
283	417053
1591	1
1046	1
782	1
769	1
472	1
Name: property, Length: 1104, dtype: int64	

```
In [44]: item_props.itemid.nunique()
```

```
Out[44]: 417053
```

```
In [45]: sn.barplot(x=valuesperproperty.property.value_counts().index[0:20], y=valuesperproperty.property.value_counts()[0:20])
```

```
Out[45]: <AxesSubplot:ylabel='property'>
```



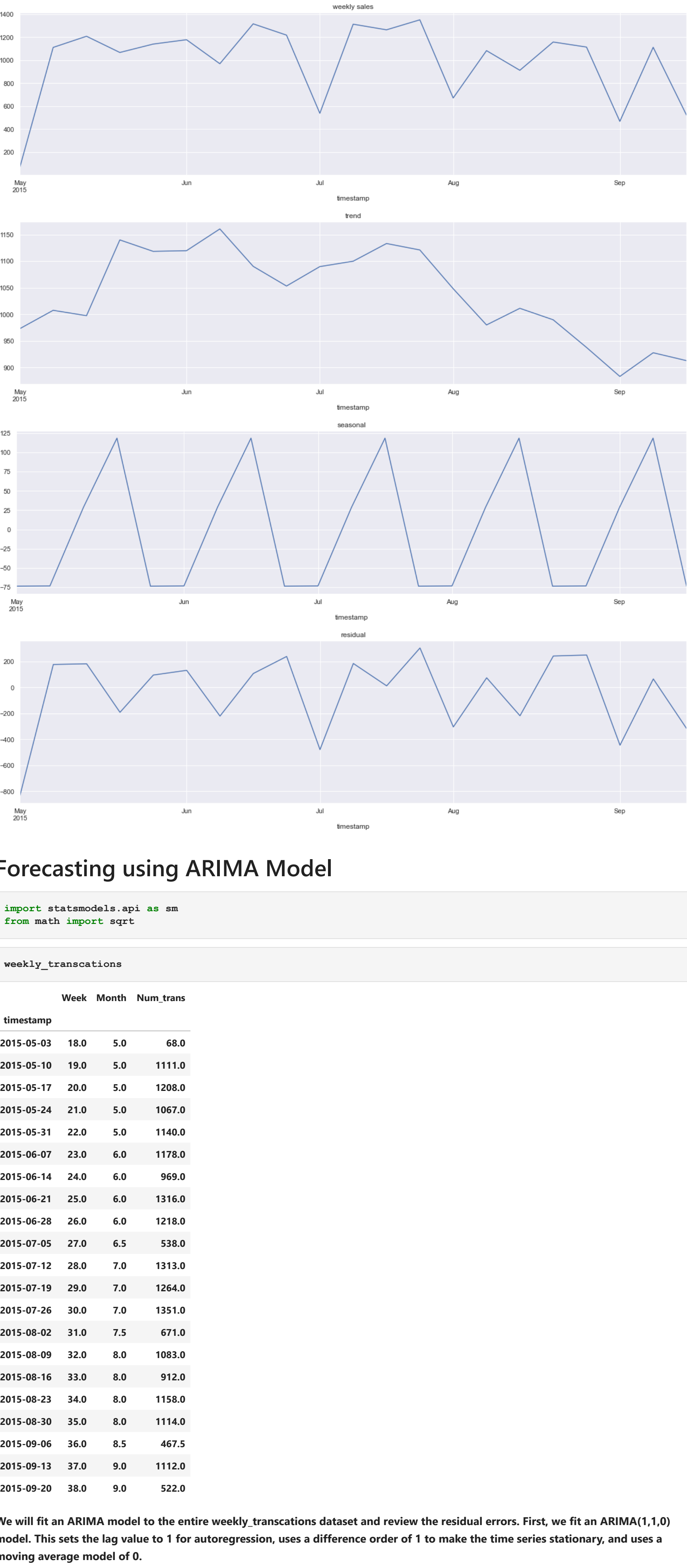
Snapshot merge

check the duplicated snapshot we have **7497165** duplicated rows of **20275902** which will reduce the data size by 40% after dropping it

```
In [46]: item_props[item_props.duplicated(['itemid','property','value'])]
```

```
Out[46]:
```

2015-09-13 05:00:00	466738	categoryid	1617
2015-09-13 05:00:00	466767	categoryid	1114
2015-09-13 05:00:00	466768	categoryid	1244



Forecasting using ARIMA Model

```
In [109]: import statsmodels.api as sm
from math import sqrt
```

```
In [110]: weekly_transactions
```

	Variance	std err	coef	std err	z	P> z	[0.025	0.975]
ar.L1	-0.5533	0.170	-3.259	0.001	-0.886	-0.221		
sigma2	1.602e+05	5.13e+04	3.121	0.002	5.96e+04	2.61e+05		
<hr/>								
tjung-box (L1) (Q):			0.00	Jarque-Bera (JB):		0.09		
Prob(Q):			0.95	Prob(JB):		0.96		
Heteroskedasticity (H):			0.45	Skew:		-0.16		
Prob(H) (two-sided):			0.32	Kurtosis:		3.07		

Warnings: [1] Covariance matrix calculated using the outer product of gradients (complex-step).

Split data to train and test and get RMSE

unning the example prints the prediction and expected value each iteration.

We can also calculate a final root mean squared error (RMSE) for the predictions

line plot is created showing the expected values (blue) compared to the rolling forecast predictions (red). We can see the values show some trend and are in the correct scale.

```
from sklearn.metrics import mean_squared_error
# split into train and test sets
X = weekly_transactions.Num_trans.values
size = int(len(X) * 0.66)
train, test = X[0:size], X[size:len(X)]
history = [x for x in train]
predictions = list()
# walk-forward validation
for t in range(len(test)):
    model = sm.tsa.ARIMA(history, order=(1,1,0))
    model_fit = model.fit()
    output = model_fit.forecast()
    yhat = output[0]
    predictions.append(yhat)
    obs = test[t]
    history.append(obs)
    print('predicted=%f, expected=%f' % (yhat, obs))
# evaluate forecasts
rmse = sqrt(mean_squared_error(test, predictions))
print('Test RMSE: %.3f' % rmse)
# plot forecasts against actual outcomes
plt.plot(test)
plt.plot(predictions, color='red')
plt.show()

predicted=1322.549499, expected=671.000000
predicted=912.161919, expected=1083.000000
predicted=912.998484, expected=912.000000
predicted=981.248073, expected=1158.000000
predicted=1051.634078, expected=1114.000000
predicted=1132.876437, expected=467.500000
predicted=731.553379, expected=1112.000000
predicted=790.378490, expected=522.000000
Test RMSE: 378.496
```

We will fit an ARIMA model to the entire weekly_transactions dataset and review the residual errors. First, we fit an ARIMA(1,1,0) model. This sets the lag value to 1 for autoregression, uses a difference order of 1 to make the time series stationary, and uses a moving average model of 0.

```
In [111]: # Fit model
model = sm.tsa.ARIMA(weekly_transactions['Num_trans'], order=(1,1,0))
model_fit = model.fit()
# summary of fit model
print(model_fit.summary())
# line plot of residuals
residuals = pd.DataFrame(model_fit.resid)
residuals.plot()
plt.show()
# density plot of residuals
# residuals.plot(kind='kde')
# plt.show()
# summary stats of residuals
# print(residuals.describe())
```

SARIMAX Results

Dep. Variable:	Num_trans	No. Observations:	21						
Model:	ARIMA(1, 1, 0)	Log Likelihood:	-148.503						
Date:	31 Dec 2021	AIC:	301.006						
Time:	18:35:35	BIC:	302.998						
Sample:	05-03-2015	HQIC:	301.395						
	09-20-2015								
Covariance Type:	opg								
	coef	std err	z	p> z	[0.025	0.975]			
ar.L1	-0.5533	0.170	-3.259	0.001	-0.886	-0.221			
sigma2	1.602e+05	5.13e+04	3.121	0.002	5.96e+04	2.61e+05			
Ljung-Box (L1) (Q):							0.00	Jarque-Bera (JB):	0.09
Prob(Q):							0.95	Prob(JB):	0.96
Heteroskedasticity (H):							0.45	B skew:	-0.16
Prob(H) (two-sided):							0.32	Rautova:	0.07

Warnings:

[1] Covariance matrix calculated using the outer product of gradients (complex-step).

Split data to train and test and get RMSE

Running the example prints the predicted and expected value each iteration.

We can also calculate a final root mean squared error (RMSE) for the predictions

A line plot is created showing the expected values (blue) compared to the rolling forecast predictions (red). We can see the values show some trend and are in the correct scale.

```
In [112]: from sklearn.metrics import mean_squared_error
# split into train and test sets
X = weekly_transactions.Num_trans.values
size = int(len(X) * 0.66)
train, test = X[0:size], X[size:len(X)]
history = [x for x in train]
predictions = list()
# walk-forward validation
for t in range(len(test)):
    model = sm.tsa.ARIMA(history, order=(1,1,0))
    model_fit = model.fit()
    output = model_fit.forecast()
    yhat = output[0]
    predictions.append(yhat)
    obs = test[t]
    history.append(obs)
    print('predicted=%f, expected=%f' % (yhat, obs))
# evaluate forecasts
rmse = sqrt(mean_squared_error(test, predictions))
print('Test RMSE: %.3f' % rmse)
# plot forecasts against actual outcomes
plt.plot(test)
plt.plot(predictions, color='red')
plt.show()
```

predicted=1322.549499, expected=671.000000
predicted=912.161919, expected=1083.000000
predicted=912.998484, expected=912.000000
predicted=983.248073, expected=1158.000000
predicted=1051.624078, expected=1114.000000
predicted=1132.876437, expected=467.500000
predicted=733.593378, expected=1112.000000
predicted=790.378490, expected=522.000000
Test RMSE: 378.496

Customer Segmentation with RFM & K-Means

```
In [114]: import pandas as pd
from datetime import datetime
```

```
# Display figures inline in Jupyter notebook
import matplotlib.pyplot as plt
import seaborn as sns
# Use seaborn style defaults and set the default figure size
sns.set(rc={'figure.figsize': (11, 4)})
import warnings
warnings.filterwarnings('ignore')
```

```
In [115]: events=pd.read_csv('events.csv')
```

Data Perpartition

```
In [116]: events['timestamp'] = events['timestamp'].apply(lambda x: datetime.fromtimestamp(x/1000))
```

```
In [117]: events
```

	timestamp	visitorid	event	itemid	transactionid
0	2015-06-02 07:02:12.117	257597	view	355908	NaN
1	2015-06-02 07:50:14.164	992329	view	248676	NaN
2	2015-06-02 07:13:19.827	111016	view	318965	NaN
3	2015-06-02 07:12:35.914	483717	view	253185	NaN
4	2015-06-02 07:02:17.106	951259	view	367447	NaN
...
2756096	2015-08-01 05:13:05.939	591435	view	261427	NaN
2756097	2015-08-01 05:30:13.142	762376	view	115946	NaN
2756098	2015-08-01 04:57:00.527	1251746	view	78144	NaN
2756099	2015-08-01 05:08:50.703	118451	view	283392	NaN
2756100	2015-08-01 05:36:03.914	199536	view	152913	NaN

2756101 rows x 5 columns

```
In [118]: events.dtypes
```

```
Out[118]: timestamp    datetime64[ns]
visitorid            int64
event               object
itemid              int64
transactionid        float64
dtype: object
```

```
In [119]: df=events.copy()
```

```
In [120]: df.isnull().sum()
```

```
Out[120]: timestamp    0
visitorid            0
event               0
itemid              0
transactionid        2733644
dtype: int64
```

```
In [121]: df.dropna(inplace = True)
df.isnull().sum()
```

```
Out[121]: timestamp    0
visitorid            0
event               0
itemid              0
transactionid        0
dtype: int64
```

```
In [122]: df.shape
```

```
Out[122]: (22457, 5)
```

```
In [123]: df.describe([0.05,0.01,0.25,0.5,0.75,0.80,0.90,0.95,0.99]).T
```

	count	mean	std	min	5%	25%	50%	75%	80%	90%
visitorid	22457	700481.127488	403657.194682	172.0	117676.00	76757.0	346997.0	712121.0	1062424.0	1133097.0
itemid	22457	236566.088614	134810.094990	15.0	4582.68	23221.8	120114.0	238804.0	352115.0	374962.0
transactionid	22457	8826.497796	5098.996290	0.0	181.56	887.8	4411.0	8813.0	13224.0	14104.8

```
In [124]: max(df.timestamp)
```

```
Out[124]: Timestamp('2015-09-18 04:43:12.017000')
```

Recency, Frequency & Monetary value calculation

we'll calculate the three key factors of RFM Analysis (recency, frequency, and monetary).

Recency: How recently customers made their purchase.

Frequency: For simplicity, we'll count the number of times each customer made a purchase.

Monetary: How much money they spent in total.

We are going to calculate these three key factors by grouping them by customers and taking "2015/09/18" as our reference end date since this is the last transaction date listed in our dataset.

```
In [125]: # Recency = Overall latest invoice date - individual customer's last invoice date
# Frequency = count of invoice no. of transaction(s)
# Monetary = Sum of Total amount for each customer
# Set 2015/09/18 as the overall last transaction date. This is to calculate recency in days.
```

```
In [126]: df
```

	timestamp	visitorid	event	itemid	transactionid
330	2015-06-02 07:17:56.276	599528	transaction	356475	4000.0
1304	2015-06-01 23:18:20.981	121688	transaction	15335	11117.0
418	2015-06-01 23:25:15.008	552148	transaction	81345	1554.0
814	2015-06-01 18:38:56.375	102019	transaction	150318	13666.0
843	2015-06-01 18:01:58.180	189384	transaction	310791	7244.0
...
2755294	2015-07-31 23:12:56.570	1050575	transaction	31640	8354.0
2755349	2015-07-31 23:17:58.779	861299	transaction	456602	3643.0
2755508	2015-07-31 17:48:50.123	855941	transaction	235771	4385.0
2755603	2015-07-31 17:12:40.300	548772	transaction	29167	13872.0
2755607	2015-07-31 18:09:49.163	1051054	transaction	312728	17579.0

22457 rows x 5 columns

```
In [127]: #Recency Metric
import datetime as dt
today_date = dt.datetime(2021,12,30)
temp_df = today_date.groupby('visitorid').agg({"timestamp":"max"})
temp_df.rename(columns={"timestamp": "Recency"}, inplace = True)
#Frequency Metric
recency_df = temp_df['Recency'].apply(lambda x: x.days)
#Monetary Metric
temp_df = df.groupby(['visitorid',"transactionid"]).agg({"transactionid":"count"})
freq_df = temp_df.groupby('visitorid').agg({"transactionid":"count"})
freq_df.rename(columns={"transactionid": "Frequency"}, inplace = True)
# Monetary Metric
## it should be item price but i dont have this field on Data
monetary_df = df.groupby('visitorid').agg({"itemid":"sum"})
monetary_df.rename(columns = {"itemid": "Monetary"}, inplace = True)
rfm = pd.concat([recency_df, freq_df, monetary_df], axis=1)
```

```
In [128]: df = rfm
df['RecencyScore'] = pd.qcut(df['Recency'], 5, labels = [5, 4, 3, 2, 1])
df['FrequencyScore'] = pd.qcut(df['Frequency'].rank(method = "first"), 5, labels = [1,2,3,4,5])
df['MonetaryScore'] = pd.qcut(df['Monetary'], 5, labels = [1,2,3,4,5])
df['RFM_SCORE'] = df['RecencyScore'].astype(str) + df['FrequencyScore'].astype(str) + df['MonetaryScore'].astype(str)
seq_map = {
    '[1-2][1-2]': 'Hibernating',
    '[1-2][3-4]': 'At Risk',
    '[1-2][5]': 'Can't Loose',
    '[3][1-2]': 'About to Sleep',
    '[3][3]': 'Need Attention!',
    '[3][4-5]': 'Loyal Customers',
    '[4][1]': 'Promising',
    '[5][1]': 'New Customers',
    '[4-5][2-3]': 'Potential Loyaltists',
    '[5][4-5]': 'Champions'
}
```

```
df['Segment'] = df['RecencyScore'].astype(str) + rfm['FrequencyScore'].astype(str)
df['Segment'] = df['Segment'].replace(seq_map, regex=True)
df.head()
```

	Recency	Frequency	Monetary	RecencyScore	FrequencyScore	MonetaryScore	RFM_SCORE	Segment
visitorid								
172	2328	1	475556	4	1	5	415	Promising
186	2331	1	49029	4	1	1	411	Promising
264	2305	1	621784	5	1	5	515	New Customers
419	2345	1	19278	4	1	1	411	Promising
539	2388	1	94371	2	1	1	211	Hibernating
...
1406787	2403	1	336832	2	5	3	253	Can't Loose
1406981	2416	1	436004	1	5	4	154	Can't Loose
1407070	2421	1	215596	1	5	2	152	Can't Loose
1407110	2338	1	360922	4	5	4	454	Loyal Customers
1407398	2367	1	218917	3	5	2	352	Loyal Customers

11719 rows x 8 columns

Clustering with the K-Means Algorithm

```
In [133]: from sklearn.cluster import KMeans
from sklearn.preprocessing import MinMaxScaler
```

```
In [134]: #scale
sc = MinMaxScaler((0,1))
df = sc.fit_transform(rfm)
```

```
In [135]: #kmeans
kmeans = KMeans(n_clusters = 10)
k_fit = kmeans.fit(df)
```

```
Out[135]: k_fit.n_clusters
```

```
Out[136]: k_fit.cluster_centers_
```

```
Out[137]: array([[9.50339184e-01, 3.31339983e-04, 2.72221896e-03],
[3.43601204e-01, 5.49583794e-04, 3.04244433e-03],
[6.51485534e-01, 6.29632685e-04, 3.03805874e-03],
[7.41322149e-02, 1.19005974e-03, 3.85068126e-03],
[8.59395075e-01, 6.73053231e-04, 3.0526279e-03],
[5.46087879e-01, 4.71346464e-04, 2.99596562e-03],
[2.09620319e-01, 6.10620846e-04, 3.13547813e-03],
[7.58668053e-01, 1.07667747e-03, 3.64626006e-03],
[4.46163173e-01, 7.71365635e-04, 3.48278091e-03],
[2.97068025e-01, 4.01397605e-04, 4.48925352e-03]])
```

```
In [137]: k_fit.labels_
```

```
Out[137]: array([6, 6, 3, ..., 0, 1, 5])
```

```
In [138]: df[0:5]
```

```
Out[139]: array([[2.46376812e-01, 0.00000000e+00, 3.74502469e-03],
[2.68115942e-01, 0.00000000e+00, 3.85928972e-04],
[7.9701449e-02, 0.00000000e+00, 4.89663724e-03],
[3.69565217e-01, 0.00000000e+00, 1.51626204e-04],
[6.81159420e-01, 0.00000000e+00, 7.43018014e-04]])
```

Determining the Optimum Number of Clusters

```
In [139]: kmeans = KMeans(n_clusters = 2)
k_fit = kmeans.fit(df)
ssd = []
K = range(1,30)
for k in K:
    kmeans = KMeans(n_clusters = k).fit(df)
    ssd.append(kmeans.inertia_)
```

```
Out[139]: Text(0.5, 1.0, 'Elbow method for Optimum number of clusters')
```



```
In [140]: # !pip install yellowbrick
```

```
In [141]: from yellowbrick.cluster import KElbowVisualizer
kmeans = KMeans()
vizu = KElbowVisualizer(kmeans, k = (2,20))
vizu.fit(df)
vizu.poof()
```



```
In [142]: kmeans = KMeans(n_clusters = 6).fit(df)
cluster = kmeans.labels_
pd.DataFrame({"visitorID": rfm.index, "cluster": cluster})
```

visitorID	cluster
0	172
1	186
2	264
3	419
4	539
...	...
11714	1406787
11715	1406981
11716	1407070
11717	1407110
11718	1407398

11719 rows x 2 columns