





```
In [ ]: # At index I will use isin to substitute the loop and get just the values with more than 1$
crostab_eda = pd.crostab(index=df_train['device.deviceCategory'], # at this line, I am using the isin to
                        columns=df_train[df_train['device.operatingSystem']]
                        .isin(df_train['device.operatingSystem'])
                        .value_counts()[1:6].index.values))['device.operatingSystem'])

# Plotting the crostab that we did above
crostab_eda.plot(kind="bar", # select the bar to plot the count of categoricals
                figsize=(14,7), # adjusting the size of graphs
                stacked=True) # code to unstack
plt.title("Most frequent OS's by Device Category of users", fontsize=22) # adjusting title and font size
plt.xlabel("Device Name", fontsize=19) # adjusting x label and font size
plt.ylabel("Count Device x OS", fontsize=19) # adjusting y label and font size
plt.xticks(rotation=45) # Adjust the xticks, rotating the labels
plt.legend(loc=1, prop={'size': 12}) # code to unstack

plt.show() # rendering
```

Very interesting values.

## SubContinent

```
In [ ]: # the top 8 of browsers represent % of total
print("Description of SubContinent count: ")
print(df_train['geoNetwork.subContinent'].value_counts()[1:8]) # printing the top 7 percentage of browser

# setting the graph size
plt.figure(figsize=(16,7))

# let explore the browser used by users
sns.countplot(df_train[df_train['geoNetwork.subContinent']]
              .isin(df_train['geoNetwork.subContinent'])
              .value_counts()[1:15].index.values))['geoNetwork.subContinent'], palette="hls")
# It's a module to color the category's

plt.title("TOP 15 most frequent SubContinents", fontsize=20) # setting the title size
plt.xlabel("SubContinent Names", fontsize=18) # setting the x label size
plt.ylabel("SubContinent Count", fontsize=18) # setting the y label size
plt.xticks(rotation=45) # Adjust the xticks, rotating the labels
plt.legend(loc=1, prop={'size': 12}) # code to unstack

plt.show() #use plt.show to render the graph that we did above
```

WoW, We have a very high number of users from North America.

TOP 5 regions are equivalent of almost 70% + of total

TOP 1 => Northern America - 44.18%  
TOP 2 => Southeast Asia - 8.29%  
TOP 3 => Northern Europe - 6.73%  
TOP 4 => South-east Asia - 6.33%  
TOP 5 => Western Europe - 6.23%

## Let's cross the SubContinent by Browser

```
In [ ]: ## I will use the crostab to explore two categorical values

# At index I will use isin to substitute the loop and get just the values with more than 1$
crostab_eda = pd.crostab(index=df_train[df_train['geoNetwork.subContinent']]
                        .isin(df_train['geoNetwork.subContinent'])
                        .value_counts()[1:10].index.values))['geoNetwork.subContinent'],
                        .value_counts()[1:5].index.values))

# Plotting the crostab that we did above
crostab_eda.plot(kind="bar", # select the bar to plot the count of categoricals
                figsize=(16,7), # adjusting the size of graphs
                stacked=True) # code to unstack
plt.title("TOP 10 Most frequent Subcontinents by Browsers used", fontsize=22) # adjusting title and font size
plt.xlabel("SubContinent Name", fontsize=19) # adjusting x label and font size
plt.ylabel("Count Subcontinent", fontsize=19) # adjusting y label and font size
plt.xticks(rotation=45) # Adjust the xticks, rotating the labels
plt.legend(loc=1, prop={'size': 12}) # code to unstack

plt.show() # rendering
```

Nice, this graph is very insightful. The North America have a low ratio of Safari x Chrome... I thought that it was the contrary

Firefox have a relative high presence in North America too.

```
In [ ]: print('train date', min(df_train['date']), 'to', max(df_train['date']))

In [ ]: year = df_train['year'].value_counts() # counting the Year with value counts
month = df_train['month'].value_counts() # counting months
weekday = df_train['weekday'].value_counts() # Counting weekday
day = df_train['day'].value_counts() # counting day
date = df_train['date'].value_counts() # Counting date

# creating the first trace with the necessary parameters
trace = go.Scatter(x=dates, y=dates, astype=str), y=dates, temp, visits,
                  opacity = 0.8, line = dict(color = color_op[3]), name= 'Visits by day')

# Below we will get the total values by Transaction Revenue Log by date
dates_temp_sum = df_train.groupby('date')['totals.transactionRevenue'].sum().to_frame().reset_index()

# using the new dates_temp_sum we will create the second trace
trace2 = go.Scatter(x=dates_temp_sum.date, astype=str), line = dict(color = color_op[1]), name= "Revenue Log by day",
                  y=dates_temp_sum['totals.transactionRevenue'], opacity = 0.8)

# getting the total values by Transactions by each date
dates_temp_count = df_train[df_train['totals.transactionRevenue'] > 0].groupby('date')['totals.transactionRevenue'].count().to_frame().reset_index()

# creating the third trace
trace3 = go.Scatter(x=dates_temp_count.date, astype=str), line = dict(color = color_op[5]), name= "Sellings by day",
                  y=dates_temp_count['totals.transactionRevenue'], opacity = 0.8)

# creating the layout will allow us to give an title and
# give us some interesting options to handle with the outputs of graphs
layout = dict(
    title="Informations by Date",
    xaxis=dict(
        rangeslider=dict(
            buttons=list([
                dict(count=1, label='1m', step='month', stepmode='backward'),
                dict(count=3, label='3m', step='month', stepmode='backward'),
                dict(count=6, label='6m', step='month', stepmode='backward'),
                dict(step='all')
            ])
        ),
        rangeslider=dict(visible = True),
        type='date'
    )
)

# creating figure with the both traces and layout
fig = dict(data=[trace, trace2, trace3], layout=layout)

#rendering the graphs
iplot(fig) #it's an equivalent to plt.show()
```

Creating a Sophisticated interactive graphics to better understanding of date features

To see the code click in 'code'.

## SELECT THE OPTION:

```
In [ ]: # Setting the first trace
trace1 = go.Histogram(x=df_train["year"],
                    name='Year Count')

# Setting the second trace
trace2 = go.Histogram(x=df_train["month"],
                    name='Month Count')

# Setting the third trace
trace3 = go.Bar(y=day, values,
               x=day.index.values,
               name='Day Count')

# Setting the fourth trace
trace4 = go.Bar(y=weekday, values,
               x=weekday.index.values,
               name='Weekday Count')

# putting all traces in the same "array of graphics" to we render it below
data = [trace1, trace2, trace4, trace3]

#Creating the options to be possible we use in our
updatemenus = list([
    dict(active=1,
          x=0.15,
          buttons=list([
              dict(
                  label = 'Years Count',
                  method = 'update',
                  args = [{'visible': [True, False, False, False, False]},
                        {'title': 'Count of Year'}]),
              dict(
                  label = 'Months Count',
                  method = 'update',
                  args = [{'visible': [False, True, False, False, False]},
                        {'title': 'Count of Months'}]),
              dict(
                  label = 'WeekDays Count',
                  method = 'update',
                  args = [{'visible': [False, False, True, False, False]},
                        {'title': 'Count of WeekDays'}]),
              dict(
                  label = 'Days Count ',
                  method = 'update',
                  args = [{'visible': [False, False, False, True, False]},
                        {'title': 'Count of Day'}]) ]
          )
    ])

layout = dict(title="The percentageal Distributions of Date Features (Select from Dropdown)",
              showlegend=False,
              updatemenus=updatemenus,
              xaxis = dict(
                  type="category"
              ),
              bargmode="group"
              )

fig = dict(data=data, layout=layout)

print("SELECT BELOW: ")
iplot(fig)

*** How can I set order to my year, months and days? ***
```

Very Cool graphs.

WE can see that the number of access are clearly downing by the through the time.

- The months with highest accesses are October and November.
- On the Weekend the traffic is lower than other days
- The 5 days with highest number of accesses is 1 and 5
- Considering the full count of dates, we can see that the days with highest accesses are almost all in november/2016

Let's investigate the VisitHour and weekday to see if we can find some interesting patterns

```
In [ ]: date_sales = ['visitHour', 'weekday'] #setting the desired
cm = sns.light_palette("green", as_cmap=True)
pd.crostab(df_train[date_sales[0]], df_train[date_sales[1]],
           values=df_train['totals.transactionRevenue'], aggfunc=np.sum).style.background_gradient(cmap = cm)

# tab.columns.levels[1] = ["Sun", "Mon", "Thu", "Wed", "Thi", "Fri", "Sat"]
```

Very interesting, we can see that from 17 to 20 hour we have the highest numbers of

## I will use a interesting graphic called Squarify

- I will apply it in feature Country to discover where the user access the store

```
In [ ]: number_of_colors = 20 # total number of different colors that we will use
# Here I will generate a bunch of hexadecimal colors
color = ['#'+''.join(random.choice('0123456789ABCDEF') for i in range(6))
         for i in range(number_of_colors)]

In [ ]: country_tree = df_train["geoNetwork.country"].value_counts() #counting the values of Country
print("Description most frequent countries: ")
print(country_tree[1:15]) #printing the 15 top most

country_tree = round((df_train["geoNetwork.country"].value_counts()[1:30] \
                    / len(df_train["geoNetwork.country"])* 100),2)

plt.figure(figsize=(14,5))
g = squarify.plot(sizes=country_tree.values, label=country_tree.index,
                  alpha=.4, color=color)
g.set_title("TOP 30 Countries - % size of total", fontsize=20)
g.set_axis_off()
plt.show()
```

USA have a very highest value than another country.

Below I will take a look on cities and find for the highest revenues from them

## Now, I will look on City feature and see the principal cities in the dataset

```
In [ ]: df_train.loc[df_train["geoNetwork.city"] == "not available in demo dataset", 'geoNetwork.city'] = np.nan

In [ ]: city_tree = df_train["geoNetwork.city"].value_counts() #counting
print("Description most frequent Cities: ")
print(city_tree[1:15])

city_tree = round((city_tree[1:30] / len(df_train["geoNetwork.city"])* 100),2)

plt.figure(figsize=(14,5))
g = squarify.plot(sizes=city_tree.values, label=city_tree.index,
                  alpha=.4, color=color)
g.set_title("TOP 30 Cities - % size of total", fontsize=20)
g.set_axis_off()
plt.show()
```

Nicely distributed clients that accessed the store. (non set) have 3.81% of total, so I don't will consider in top five, but it was the top 2 most frequent.

The top 5 are:

- Mountain View
- New York
- San Francisco
- Sunnyvale
- London

And in terms of money, how the Countries and Cities are ?

## Creating a function with the plotted to better investigate the dataset

- Click in 'code' to see the commented code

```
In [ ]: def PieChart(df_colunm, title, limit=15):
    """
    This function helps to investigate the proportion of visits and total of transaction revenue
    by each category.
    """
    count_trace = df_train[df_colunm].value_counts()[1:limit].to_frame().reset_index()
    rev_trace = df_train.groupby(df_colunm)['totals.transactionRevenue'].sum().nlargest(10).to_frame().reset_index()

    trace1 = go.Pie(labels=count_trace['index'], values=count_trace[df_colunm], name= "% Accesses", hole=.5,
                    hoverinfo="label+percent+name", showlegend=True, domain= {'x': [0, .48]}),
    marker=dict(colors=color))

    trace2 = go.Pie(labels=rev_trace[df_colunm],
                    values=rev_trace['totals.transactionRevenue'], name= "% Revenue", hole=.5,
                    hoverinfo="label+percent+name", showlegend=False, domain= {'x': [.52, 1]})

    layout = dict(title=title, height=450, font=dict(size=15),
                  annotations = [
                      dict(
                          x=.25, y=.5,
                          text="Visits",
                          showarrow=False,
                          font=dict(size=20)
                      ),
                      dict(
                          x=.80, y=.5,
                          text="Revenue",
                          showarrow=False,
                          font=dict(size=20)
                      )
                  ])

    fig = dict(data=[trace1, trace2], layout=layout)
    iplot(fig)
```

## Device Category feature

```
In [ ]: PieChart("device.deviceCategory", "Device Category")
```

## I will apply the Prie Chart in Country's again

```
In [ ]: # call the function
PieChart("geoNetwork.country", "Top Cities by Accesses and Revenue", limit=12)

• New York is responsible by 14% of visits and 31% of revenues.
• Mountain view have 19% in visits but just 16% of revenues
• Chicago have just 3.5% of visits but have a high significance in revenues
```

## Seeing again Channel Grouping more specified

```
In [ ]: PieChart("channelGrouping", "Channel Grouping Visits and Revenues")

It's interesting to note that Referral have a less number of Visits but is responsible for almost 40% of revenues"
```

## Months in pizza graph

## Let's see the NetWork Domain

- I will plot visits and revenues by each category, including the non-set and unknown accesses and revenues

```
In [ ]: PieChart("geoNetwork.networkDomain", "Network Domain")

Wow, another very cool information.

• (not set) domain have almost 50% of total visits and 62% of Revenues.
• Unknown is responsible by 28% of visits but just 2.70% of Revenues
• comcast.net have 5.5% of visits and 7.4% Revenues.
```

Let's take a look on Mobile and Browser proportions

```
In [ ]: PieChart("device.deviceCategory", "Device Category")

The absolutely high part of revenues are from Desktop Devices
```

## Traffic Source Medium

```
In [ ]: PieChart("trafficSource.medium", "Traffic Source - Medium")

• Organic have highest number of visits but is the third in revenues
• Referral have almost 40% in both Visits and Revenues
• The none category have almost 16% of visits but almost 40% of revenues
```

Now I will take a look on trafficSource source, the Source to access the store

```
In [ ]: PieChart("trafficSource.source", "Visits and Revenue by TOP Sources", limit=8)

We have a high number of visits from youtube but 0 sales.
the mail.googleplex is have a low number of access but have the highest value in revenues
```

## I will continue this notebook! Votes up the kernel and stay tuned to next updates

```
In [ ]: df_train.corr()['totals.transactionRevenue']

Seeing the crostab with heatmap
```

```
In [ ]: country_repayment = ['channelGrouping', 'weekday'] #setting the desired
cm = sns.light_palette("green", as_cmap=True)
pd.crostab(df_train[country_repayment[0]], df_train[country_repayment[1]],
           values=df_train['totals.transactionRevenue'], aggfunc=np.sum).style.background_gradient(cmap = cm)
# tab.columns.levels[1] = ["Sun", "Mon", "Thu", "Wed", "Thi", "Fri", "Sat"]
```

## Geolocation plot to visually understand the data

```
In [ ]: # Counting total visits by country
countMaps = pd.DataFrame(df_train['geoNetwork.country'].value_counts()).reset_index()
countMaps.columns = ['country', 'counts'] #renaming columns
sumRevMaps = countMaps.reset_index().drop('index', axis=1) #reseting index and dropping the column

data = [ dict(
    type = 'choropleth',
    locations = countMaps['country'],
    locationmode = 'country names',
    z = countMaps['counts'],
    text = countMaps['country'],
    autocolorscale = False,
    marker = dict(
        line = dict(
            color = 'rgb(180,180,180)',
            width = 0.5
        ),
        colorbar = dict(
            autotick = False,
            tickprefix = '',
            title = 'Number of Visits'),
        )
    )

layout = dict(
    title = 'Counting Visits Per Country',
    geo = dict(
        showframe = False,
        showcoastlines = True,
        projection = dict(
            type = 'Mercator'
        )
    )
)

figure = dict(data=data, layout=layout)
iplot(figure, validate=False, filename='map-countrys-count')
```

## Total Revenues by Country

```
In [ ]: # I will create a variable of Revenues by country sum
sumRevMaps = df_train[df_train['totals.transactionRevenue'] > 0].groupby("geoNetwork.country")["totals.transactionRevenue"].count().to_frame().reset_index()
sumRevMaps.columns = ['country', 'count sales'] # renaming columns
sumRevMaps = sumRevMaps.reset_index().drop('index', axis=1) #reseting index and drop index column

data = [ dict(
    type = 'choropleth',
    locations = sumRevMaps['country'],
    locationmode = 'country names',
    z = sumRevMaps['count sales'],
    text = sumRevMaps['country'],
    autocolorscale = False,
    marker = dict(
        line = dict(
            color = 'rgb(180,180,180)',
            width = 0.5
        ),
        colorbar = dict(
            autotick = False,
            tickprefix = '',
            title = 'Count of Sales'),
        )
    )

layout = dict(
    title = 'Total Sales by Country',
    geo = dict(
        showframe = False,
        showcoastlines = True,
        projection = dict(
            type = 'Mercator'
        )
    )
)

figure = dict(data=data, layout=layout)
iplot(figure, validate=False, filename='map-countrys-total')
```

## Some tests that I am doing to try find interesting feature engineering approaches

```
In [ ]: df_train['month_unique_user_count'] = df_train.groupby(['month'])['fullVisitorId'].transform('unique')
df_train['day_unique_user_count'] = df_train.groupby(['day'])['fullVisitorId'].transform('unique')
df_train['weekday_unique_user_count'] = df_train.groupby(['weekday'])['fullVisitorId'].transform('unique')

df_train['traf_sourc_browser_count'] = df_train.groupby(['trafficSource.medium', 'device.browser'])['totals.pageviews'].transform('unique')
df_train['id_browser_pageviews_sumprod'] = df_train.groupby(['fullVisitorId', 'device.browser'])['total.pageviews'].transform('sumprod')
df_train['id_browser_hits_sumprod'] = df_train.groupby(['fullVisitorId', 'device.browser'])['totals.hits'].transform('sumprod')
df_train['id_browser_hits_sumprod_mob'] = df_train.groupby(['fullVisitorId', 'device.browser', 'device.isMobile'])['totals.hits'].transform('sum')

df_train['id_networkDomain_hits'] = df_train.groupby(['fullVisitorId', 'geoNetwork.networkDomain'])['totals.hits'].transform('var')
df_train['id_networkDomain_country_hits'] = df_train.groupby(['fullVisitorId', 'geoNetwork.networkDomain', 'geoNetwork.country'])['totals.hits'].transform('unique')
```

```
In [ ]: df_train['totals.transactionRevenue', 'id_browser_hits_sumprod', 'id_networkDomain_hits', 'id_networkDomain_country_hits', 'id_browser_hits_sumprod_mob'].corr()
```

## Preprocessing the fulldataset and creating new features

```
In [ ]: aggs = {
    'date': ['min', 'max'],
    'totals.hits': ['sum', 'min', 'max', 'mean', 'median'],
    'totals.pageviews': ['sum', 'min', 'max', 'mean', 'median'],
    'totals.bounces': ['sum', 'mean', 'median'],
    'totals.newVisits': ['sum', 'mean', 'median']
}

# Previous applications categorical features
cat_aggregations = {}

for cat in dummy_features:
    cat_aggregations[cat] = ['min', 'max', 'mean']

prev_agg = df_train.groupby(['fullVisitorId']).agg(**aggs)
prev_agg.columns = pd.Index(['Agg_1' + e[0] + "_" + e[1].upper() for e in prev_agg.columns.tolist()])
```

prev\_agg

```
In [ ]: new_columns = {
    k + "_" + agg for k in aggs.keys() for agg in aggs[k]
}
new_columns

In [ ]: dummy_features

In [ ]: ## Testing some grouping approaches

In [ ]: df_train['cumcount'] = df_train.groupby(['fullVisitorId']).cumcount() + 1

Some tests to feature engineering
```

```
In [ ]: aggs = {
    'date': ['min', 'max'],
    'totals.transactionRevenue': ['sum', 'size'],
    'totals.hits': ['sum', 'min', 'max', 'count', 'median'],
    'totals.pageviews': ['sum', 'min', 'max', 'mean', 'median'],
    'totals.bounces': ['sum', 'mean', 'median'],
    'totals.newVisits': ['sum', 'mean', 'median']
}

# Previous applications categorical features
cat_aggregations = {}

for cat in dummy_features:
    cat_aggregations[cat] = ['min', 'max', 'mean']

prev_agg = df_train.groupby(['fullVisitorId']).agg(**aggs)
prev_agg.head()
```

## I will continue working on this kernel, stay tuned

\*\* Please, if you liked this kernel don't forget to votes up and give your feedback \*\*

```
In [ ]: prev_agg.columns = ["_".join(x) for x in prev_agg.columns.ravel()]

In [ ]: prev_agg.head()

In [ ]:
```