

▼ DATA VISUALIZATION

by Dr Juan H Klopper

- Research Fellow
- School for Data Science and Computational Thinking
- Stellenbosch University



▼ INTRODUCTION

Visualizing data is not only a pleasing activity, but it can provide an even richer understanding of the data than summary statistics. Data visualisation is also part and parcel of data communication, a key activity in Data Science.

▼ PACKAGES USED IN THIS NOTEBOOK

In this notebook, we take a look at one of the myriad plotting packages in Python called plotly. We choose it for this course due to its ability to produce modern, multi-use plots, ready for online reports and print documents such as research papers. Matplotlib, seaborn, and altair are some of the other established plotting packages.

```
1 import pandas as pd
```

Two of the main plotting modules in plotly are the `graph_objects` and the `express` modules. We import these with the commonly used namespace abbreviations `go` and `px`. We also import the `io` module to set a theme for our plots. We choose `plotly_white` since this notebook uses a white background. You can type `pio.templates` which will return a list of all the themes.

```
1 import plotly.graph_objects as go
```

```

2 import plotly.express as px
3 import plotly.io as pio
4 pio.templates.default = 'plotly_white'

1 pio.templates # A list of the available plot themes

```

```

Templates configuration
-----

```

```

    Default template: 'plotly_white'
    Available templates:
      ['ggplot2', 'seaborn', 'simple_white', 'plotly',
       'plotly_white', 'plotly_dark', 'presentation', 'xgridoff',
       'ygridoff', 'gridon', 'none']

```

This notebook was created on a computer with a retina display. In order to plot high resolution plots, we use the %config magic command below.

```

1 %config InlineBackend.figure_format = "retina" # For Retina type displays

```

Finally, the %load_ext magic command is used to set the way tables are printed to the screen in Google Colab.

```

1 %load_ext google.colab.data_table

```

The drive function allows us to connect to the files in our Google Drive.

```

1 from google.colab import drive # Connect to Google Drive

```

▼ DATA IMPORT

Since our data is in the data folder in Google Drive, we need to mount the drive.

```

1 # Log on and list files in the DATA directory of our Google Drive
2 drive.mount('/gdrive')
3 %cd '/gdrive/My Drive/Stellenbosch University/School for Data Science and Comput
4 %ls

```

```

Mounted at /gdrive
/gdrive/My Drive/Stellenbosch University/School for Data Science and Computat
australia_rain.csv      crops.csv              DefaultMissingData.csv
bitcoin_ethereum.csv  customers.csv          kaggle_survey_2020_responses.csv
breast_cancer.csv     data.csv               MissingData.csv
client_data.csv       DatesTimes.csv        montague_gardens_construction.csv

```

The data set contains data on banking customers and is used in machine learning to predict

```
1 df = pd.read_csv('customers.csv')
```

Later, we will also use data about rainfall in Australia.

```
1 rain = pd.read_csv('australia_rain.csv')
```

```
1 df # Viewing a table of the customer data
```

1 to 25 of 10127 entries Filter ?

index	Attrition_Flag	Customer_Age	Dependent_count	Education_Level	Marital_Status	Income_
0	Existing Customer	49	4	Uneducated	Single	R80K - F
1	Attrited Customer	48	2	Graduate	Married	R60K - F
2	Existing Customer	44	3	Graduate	Unknown	R80K - F
3	Existing Customer	48	3	Graduate	Single	R80K - F
4	Existing Customer	49	1	Graduate	Married	R120K +

We consider some meta data about the dataframe object.

```
1 df.shape # Number of observations (rows) and variables (columns)
```

```
(10127, 14)
```

2	Attrited	51	3	College	Single	R120K +
---	----------	----	---	---------	--------	---------

```
1 df.columns # Statistical variables (column names)
```

```
Index(['Attrition_Flag', 'Customer_Age', 'Dependent_count', 'Education_Level',
       'Marital_Status', 'Income_Category', 'Card_Category', 'Months_on_book',
       'Total_Relationship_Count', 'Months_Inactive_12_mon',
       'Contacts_Count_12_mon', 'Credit_Limit', 'Total_Revolving_Bal',
       'Avg_Open_To_Buy'],
      dtype='object')
```

13	Attrited	40	3	College	Single	R80K - F
----	----------	----	---	---------	--------	----------

```
1 # Data types of each column
```

```
2 df.dtypes
```

```
Attrition_Flag      object
Customer_Age        int64
Dependent_count     int64
Education_Level     object
Marital_Status      object
Income_Category     object
Card_Category       object
Months_on_book      int64
Total_Relationship_Count  int64
Months_Inactive_12_mon  int64
Contacts_Count_12_mon  int64
Credit_Limit        object
Total_Revolving_Bal  int64
Avg_Open_To_Buy     object
dtype: object
```

44	Customer	52	2	Graduate	Single	R80K - F
----	----------	----	---	----------	--------	----------

A method that we have not used before is the `info` method. It returns a pandas dataframe object with a column of all the variables, a column of the number of non missing data, and a column of the data types.

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

```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 10127 entries, 0 to 10126
Data columns (total 14 columns):
#   Column                                Non-Null Count  Dtype
---  -
0   Attrition_Flag                        10127 non-null  object
1   Customer_Age                          10127 non-null  int64
2   Dependent_count                       10127 non-null  int64
3   Education_Level                       10127 non-null  object
4   Marital_Status                        10127 non-null  object
5   Income_Category                       10127 non-null  object
6   Card_Category                         10127 non-null  object
7   Months_on_book                        10127 non-null  int64
8   Total_Relationship_Count              10127 non-null  int64
9   Months_Inactive_12_mon                10127 non-null  int64
10  Contacts_Count_12_mon                 10127 non-null  int64
11  Credit_Limit                          10127 non-null  object
12  Total_Revolving_Bal                   10127 non-null  int64
13  Avg_Open_To_Buy                       10127 non-null  object
dtypes: int64(7), object(7)
memory usage: 1.1+ MB

```

We start our journey introducing only the most commonly used plots. In fact, the plotly package is enormous and we cannot possibly look at all the plots that it can create. There is a lot more information at the [plotly](https://plotly.com) website.

▼ BAR PLOT

Bar plots are great for indicating frequency and relative frequency. In other words, counting the sample space elements of categorical or discrete data.

One axis, usually the horizontal axis is reserved to indicate the sample space elements of the categorical variable. The other axis is used to show the frequency or relative frequency, i.e. the height of a bar. There are spaces in between the bars to indicate that the sample space elements are indeed not a continuity.

Below, we use the `value_counts` method on the series object of the `Attrition_Flag` variable.

```

1 # Calculate the frequency count
2 df.Attrition_Flag.value_counts()

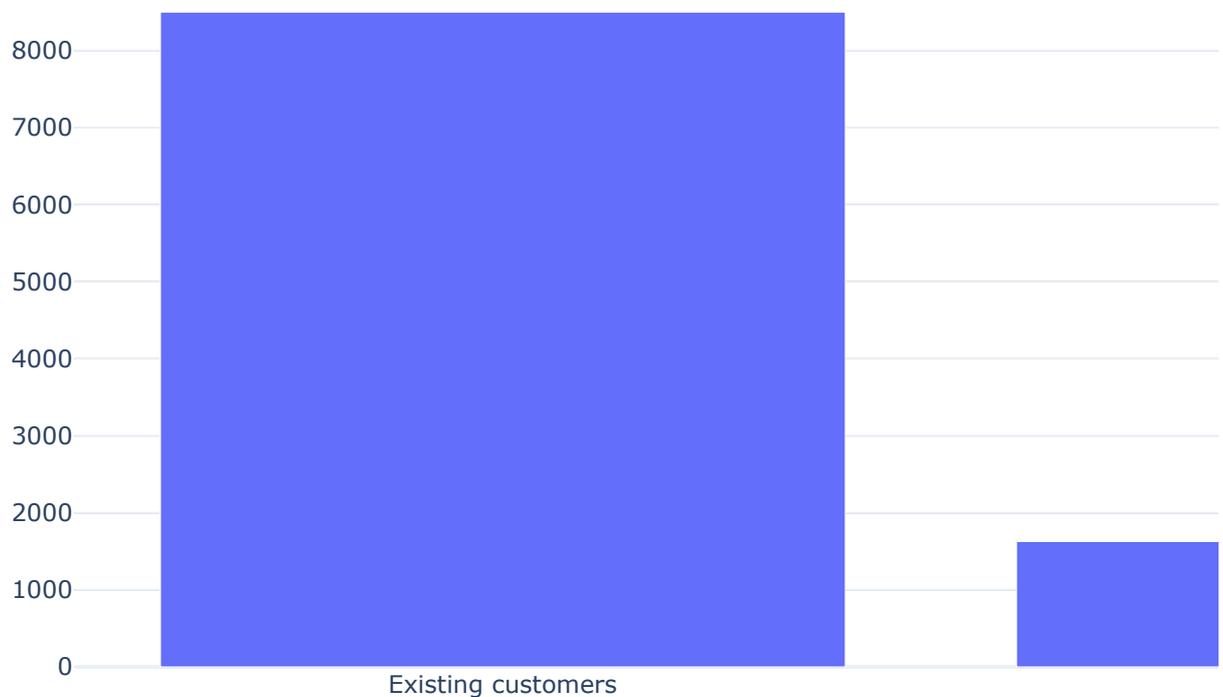
Existing Customer      8500
Attrited Customer     1627
Name: Attrition_Flag, dtype: int64

```

A total of 8500 customers are still with the bank (`Existing Customer`) and 1627 have left the bank (`Attrited Customer`). We visualize this data as a bar plot. The `Figure` function in the

`graph_objects` module is assigned to a computer variable. We then use the `add_trace` method on this figure to add a bar chart using the `Bar` function. We set the `x` axis (sample space elements of the categorical variable as Python a list object) and the `y` axis values (frequencies also as a Python list object).

```
1 churn_fig = go.Figure() # Simple bar chart
2
3 churn_fig.add_trace(
4     go.Bar(
5         x=['Existing customers', 'Lost customers'],
6         y=[8500, 1627]
7     )
8 )
```



We entered the `x` and `y` values by hand. While this is easy enough for smaller sample element numbers, it is not always the case. We can get the sample space elements using the `unique()` method. It will be returned in the order in which the method discovers them in the relevant pandas series object.

```
1 df.Education_Level.unique()

array(['Uneducated', 'Graduate', 'College', 'Unknown', 'High School',
```

```
'Post-Graduate', 'Doctorate'], dtype=object)
```

This might not be the order in which we want the bars to appear in.

With the default argument values of the `.value_counts()` method, we get the frequencies in descending order.

```
1 df.Education_Level.value_counts()

Graduate      3128
High School   2013
Unknown       1519
Uneducated    1487
College       1013
Post-Graduate  516
Doctorate     451
Name: Education_Level, dtype: int64
```

This might also not be useful as in the case of a statistical variable such as *Month*. The bottom line is that we have to be careful with our code when designing plots.

The `Bar` function requires the frequency variable to be a list of values. We extract that below using the `value` property and the `tolist` method.

```
1 totals = df.Education_Level.value_counts().values.tolist()
2 totals

[3128, 2013, 1519, 1487, 1013, 516, 451]
```

Since we are using a series object, with the sample space elements as the index, we can use the `index` attribute. We can also convert it to a list using the `tolist` method.

```
1 levels = df.Education_Level.value_counts().index.tolist()
2 levels

['Graduate',
 'High School',
 'Unknown',
 'Uneducated',
 'College',
 'Post-Graduate',
 'Doctorate']
```

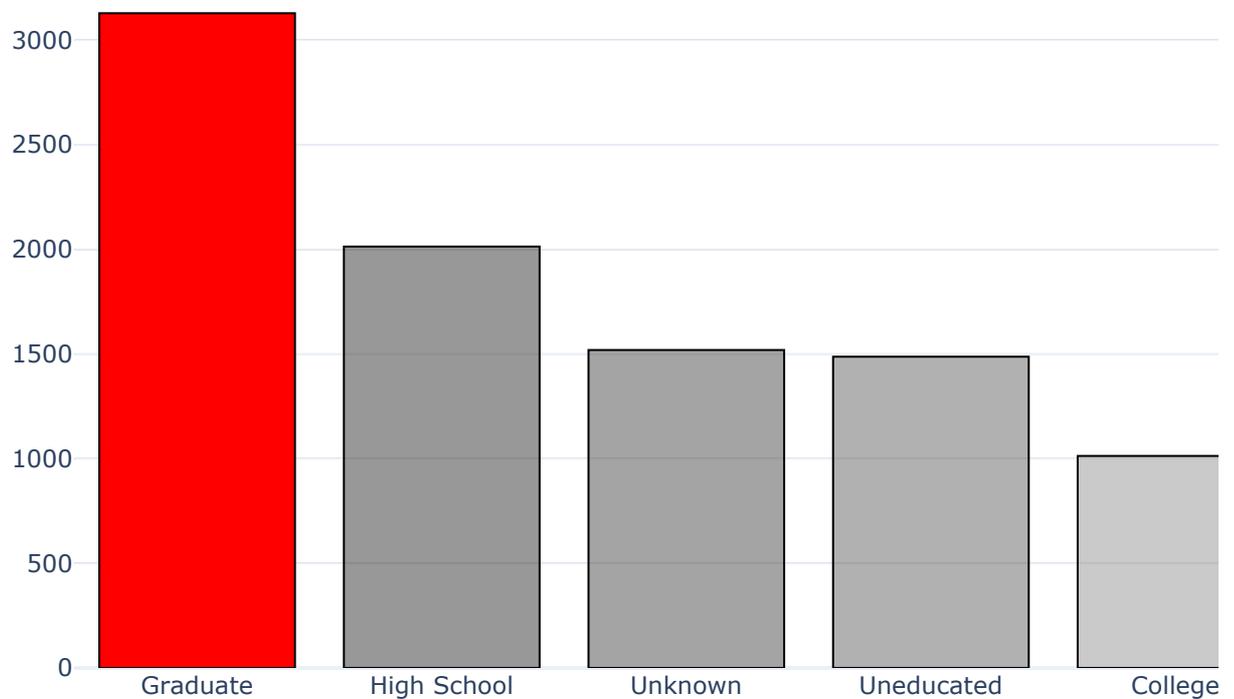
With these lists saved, we can use them for our `x` and `y` arguments in the `Bar` function. We also use the `marker` argument and set its value using a dictionary. We specify seven colors and also a bar outline color and line thickness. Colors can be specified by their name such as `red` or

blue or by RGBA values. The latter uses four values. The first three are for red, green, and blue

```

1 # Simple bar plot
2 churn_fig = go.Figure()
3
4 churn_fig.add_trace(go.Bar(
5     x=levels,
6     y=totals,
7     marker={'color':['red', 'rgba(50, 50, 50, 0.5)', 'rgba(75, 75, 75, 0.5)', '1
8         'line':{'color':'black', 'width':1}}
9 ))
10
11 churn_fig.show()

```



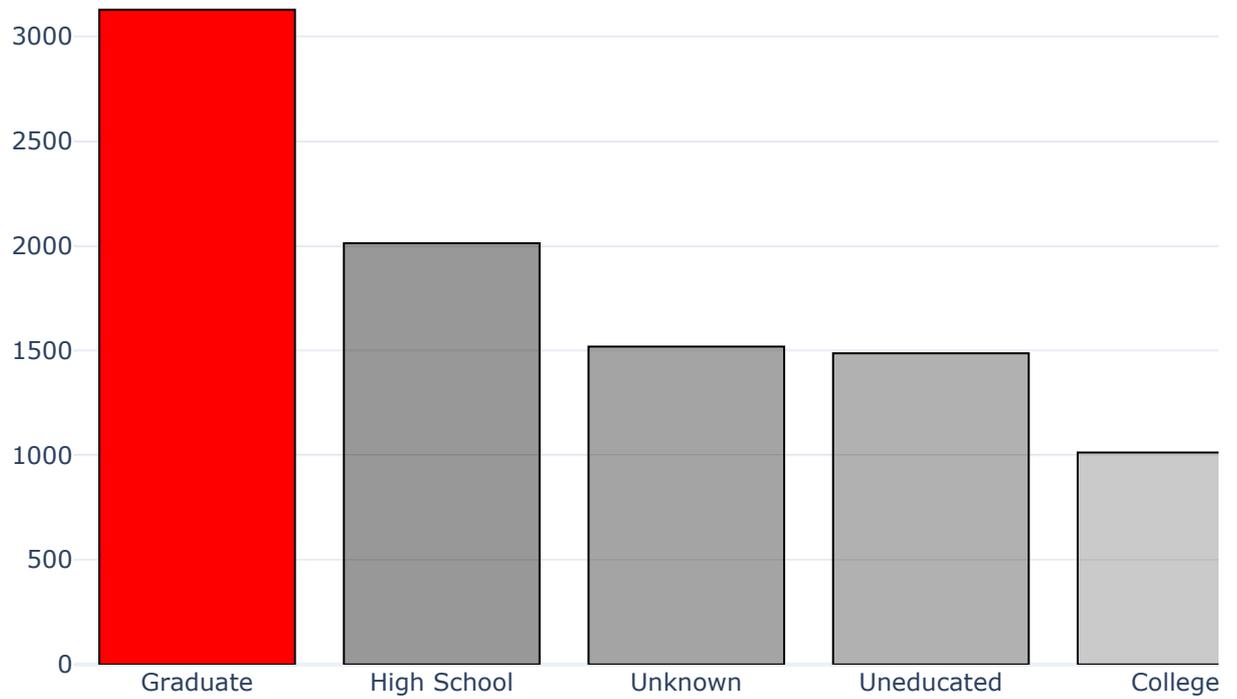
The figure is assigned to a computer variable. We can use the `update_layout` method to add a title and axes labels.

```

1 # Add a title
2 churn_fig.update_layout(title='Number of customers in each education level')
3 churn_fig.show()

```

Number of customers in each education level



```
1 # Add axes labels
2 churn_fig.update_layout(xaxis=dict(title='Education level'),
3                           yaxis=dict(title='Counts'))
4 churn_fig.show()
```

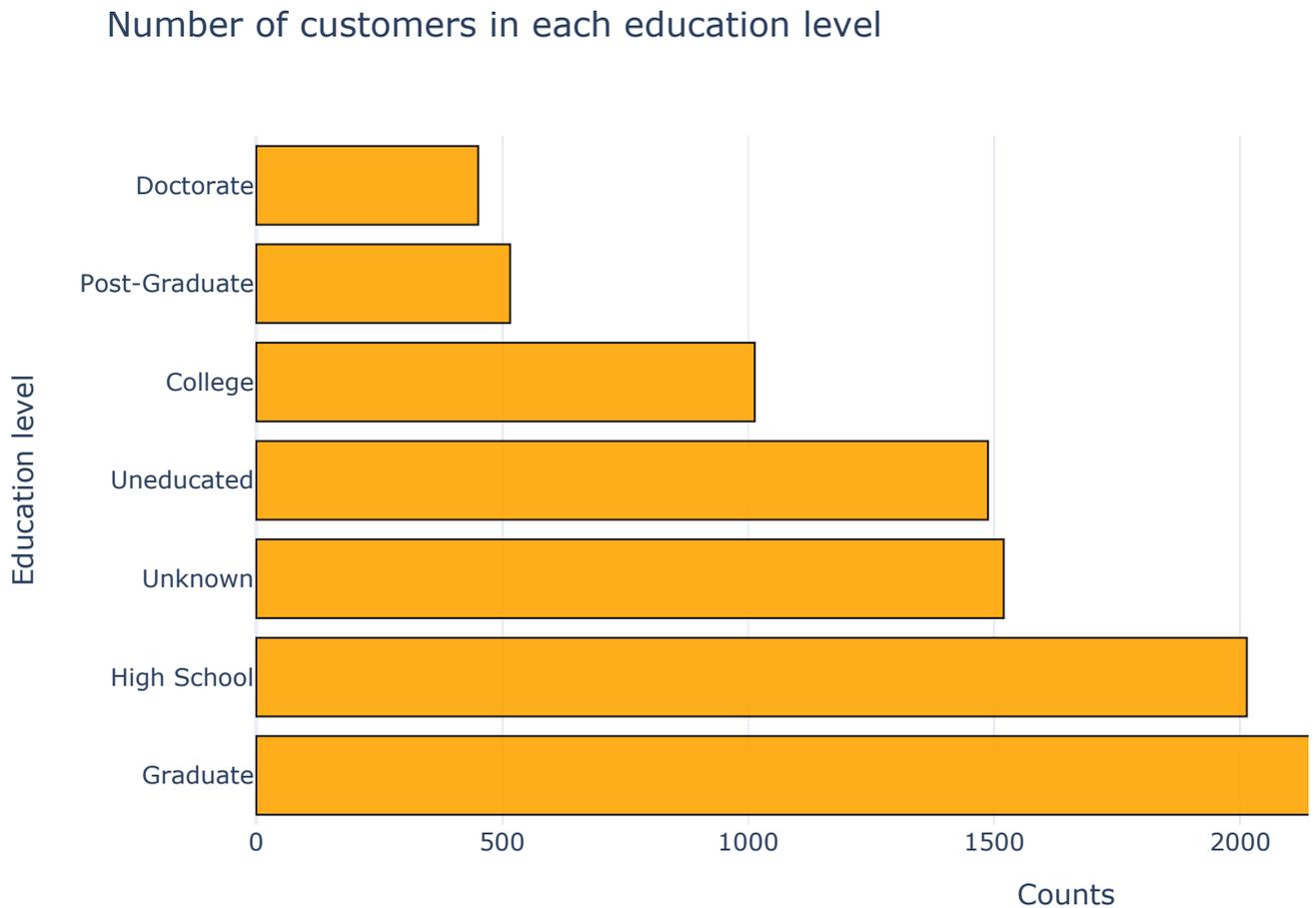
Number of customers in each education level

In the bar plot below, we add many more properties as illustration of what is possible.

```
1 churn_fig = go.Figure()  
2  
3 churn_fig.add_trace(go.Bar(  
4     x=levels,  
5     y=totals,  
6     text=levels,  
7     textposition='outside',  
8     marker={'color':'deepskyblue',  
9             'line':{'color':'black', 'width':1},  
10            'opacity':0.9}  
11 ))  
12  
13 churn_fig.update_layout(title='Number of customers in each education level')  
14  
15 churn_fig.update_layout(xaxis = dict(title='Education level'),  
16                           xaxis_tickangle=-25,  
17                           yaxis=dict(title='Counts'))  
18  
19 churn_fig.show()
```

We can also create horizontal bar plots using the `orientation` argument set to `h`, remembering to switch the axes information.

```
1 churn_fig = go.Figure()  
2  
3 churn_fig.add_trace(go.Bar(  
4     y=levels,  
5     x=totals,  
6     marker={'color':'orange',  
7           'line':{'color':'black', 'width':1},  
8           'opacity':0.9},  
9     orientation = 'h'  
10 ))  
11  
12 churn_fig.update_layout(title='Number of customers in each education level')  
13  
14 churn_fig.update_layout(yaxis = dict(title='Education level'),  
15                          xaxis=dict(title='Counts'))  
16  
17 churn_fig.show()
```



We can group data by another categorical variable. Below, we plot the frequencies of the different educational levels but divided for each of the customers groups (lost versus existing). We start with the `crosstab` function in pandas. It returns a dataframe object. We pass two categorical variables. The variable listed first results in the rows and the sample space elements of the second variable becomes the column names

```
1 attr_edu = pd.crosstab(df.Attrition_Flag, df.Education_Level)
2 attr_edu
```

1 to 2 of 2 entries

Attrition_Flag	College	Doctorate	Graduate	High School	Post-Graduate	Uneducated	Unkn
Attrited Customer	154	95	487	306	92	237	
Existing Customer	859	356	2641	1707	424	1250	

Show per page

Note that the order of the educational levels is now different from when we used the `value_counts` method on a series object. We save the educational level values as a list.

```
1 levels = attr_edu.columns.tolist()
2 levels
```

```
['College',
 'Doctorate',
 'Graduate',
 'High School',
 'Post-Graduate',
 'Uneducated',
 'Unknown']
```

Since our object is a dataframe, we can use indexing. Below, we do this by using `iloc` indexing. With it we specify the row and the column values we require. Remember that the colon symbol is short-hand for including all the values (columns in this case).

```
1 attr_values = attr_edu.iloc[0, :].tolist()
2 attr_values
```

```
[154, 95, 487, 306, 92, 237, 256]
```

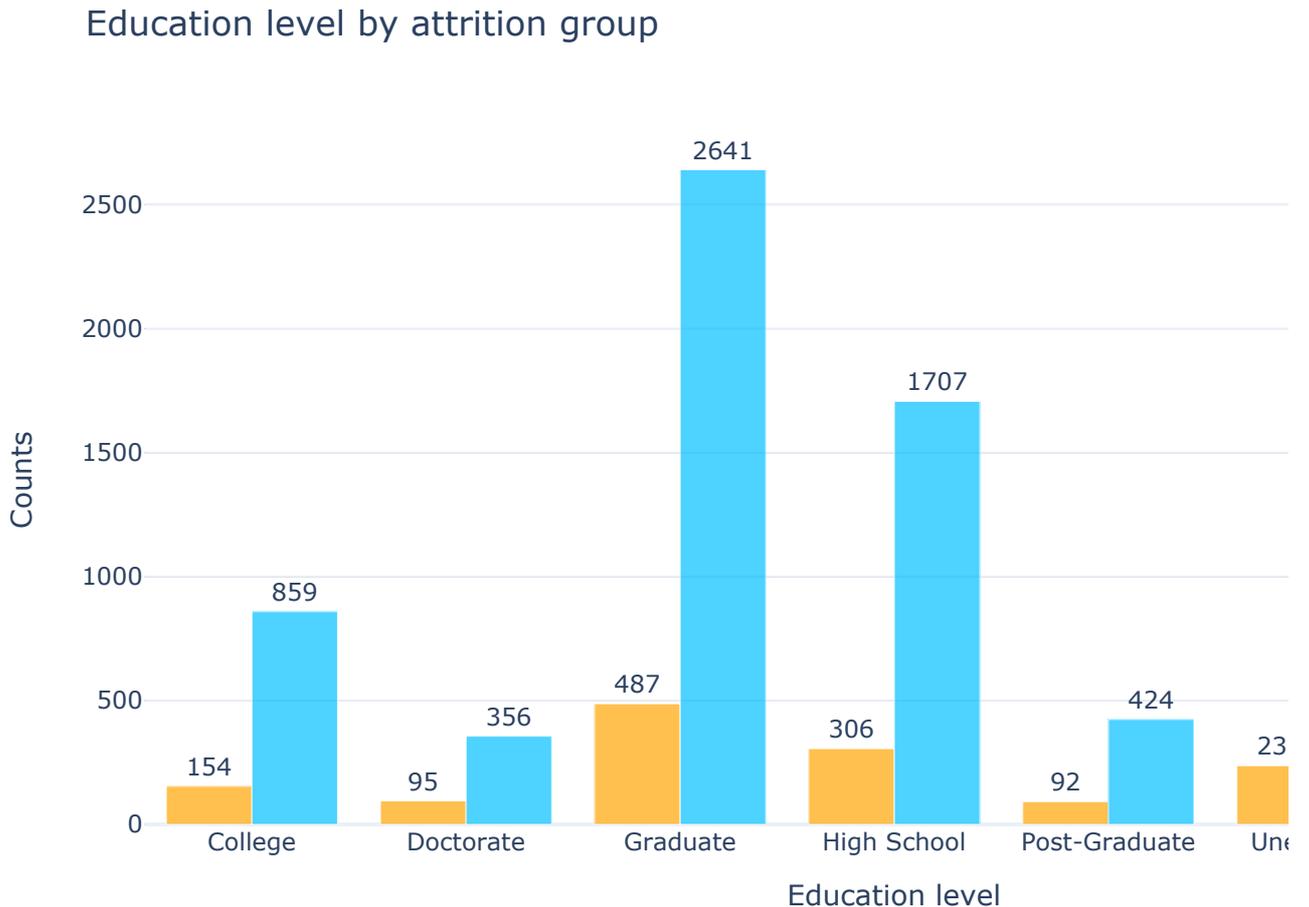
```
1 exis_values = attr_edu.iloc[1, :].tolist()
2 exis_values
```

```
[859, 356, 2641, 1707, 424, 1250, 1263]
```

Now we can generate two traces using the `add_trace` function twice. We also set `barmode` argument to `group` in the `update_layout` method to create a grouped bar chart.

```
1 churn_edu_fig = go.Figure()
2
```

```
3 churn_edu_fig.add_trace(go.Bar(  
4     x=levels,  
5     y=attr_values,  
6     text=attr_values,  
7     textposition='outside',  
8     name='Lost customers',  
9     marker={'color':'orange', 'opacity':0.7}  
10 ))  
11  
12 churn_edu_fig.add_trace(go.Bar(  
13     x=levels,  
14     y=exis_values,  
15     text=exis_values,  
16     textposition='outside',  
17     name='Existing customers',  
18     marker={'color':'deepskyblue', 'opacity':0.7}  
19 ))  
20  
21 churn_edu_fig.update_layout(title='Education level by attrition group')  
22  
23 churn_edu_fig.update_layout(xaxis = dict(title='Education level'),  
24                             yaxis=dict(title='Counts'),  
25                             bargroup='group')  
26  
27 churn_edu_fig.show()
```



▼ Exercise (advanced)

Create the bar plot above, but group the horizontal axis by the customer groups.

▼ Solution

One way to generate a plot for this exercise is to use the idea of creating arrays and lists so that we can loop over their elements.

We start by generating a dataframe object, using the reverse order of the categorical variables in the `crosstab` function.

```
1 attr_edu = pd.crosstab(df.Education_Level, df.Attrition_Flag)
2 attr_edu
```

1 to 7 of 7 entries 

Education_Level	Attrited Customer	Existing Customer
College	154	859
Doctorate	95	356
Graduate	487	2641
High School	306	1707
Post-Graduate	92	424
Uneducated	237	1250
Unknown	256	1263

Show per page

We save the column values again.

```
1 groups = attr_edu.columns.tolist()
2 groups

['Attrited Customer', 'Existing Customer']
```

We can extract all the dataframe values as an array.

```
1 values = attr_edu.values
2 values

array([[ 154,  859],
       [  95,  356],
       [ 487, 2641],
       [ 306, 1707],
       [  92,  424],
```

```
[ 237, 1250],
[ 256, 1262111]
```

We also still need a list of all the education levels. This is stored as the index of the dataframe object.

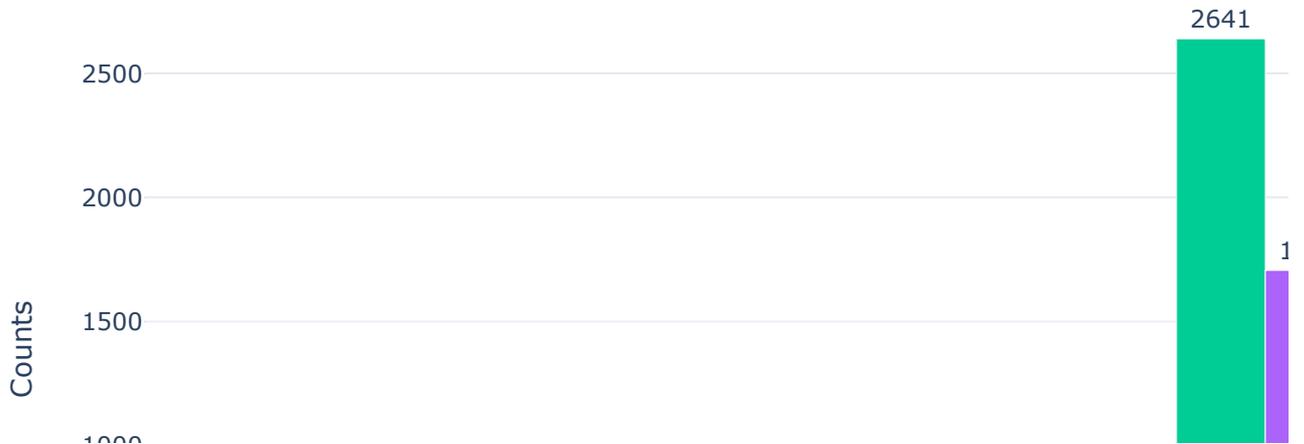
```
1 levels = attr_edu.index.tolist()
2 levels
```

```
['College',
 'Doctorate',
 'Graduate',
 'High School',
 'Post-Graduate',
 'Uneducated',
 'Unknown']
```

Now we generate seven traces by looping over the seven elements in the (values) array and the (education level) list above. We could also have added each trace manually. The idea is to use the power of the Python language to do all the hard work.

```
1 churn_edu_fig = go.Figure()
2
3 # Loop over elements in grps array object and nms list object
4 for i in range(len(values)):
5     churn_edu_fig.add_trace(go.Bar(
6         x=groups,
7         y=values[i],
8         text=values[i],
9         textposition='outside',
10        name=levels[i],
11    ))
12
13 churn_edu_fig.update_layout(title='Educational level frequencies by customer group')
14
15 # Using alternate dictionary syntax
16 churn_edu_fig.update_layout({'xaxis':{'title':'Customer group'},
17                             'yaxis':{'title':'Counts'},
18                             'barmode':'group'})
19
20 churn_edu_fig.show()
```

Educational level frequencies by customer group



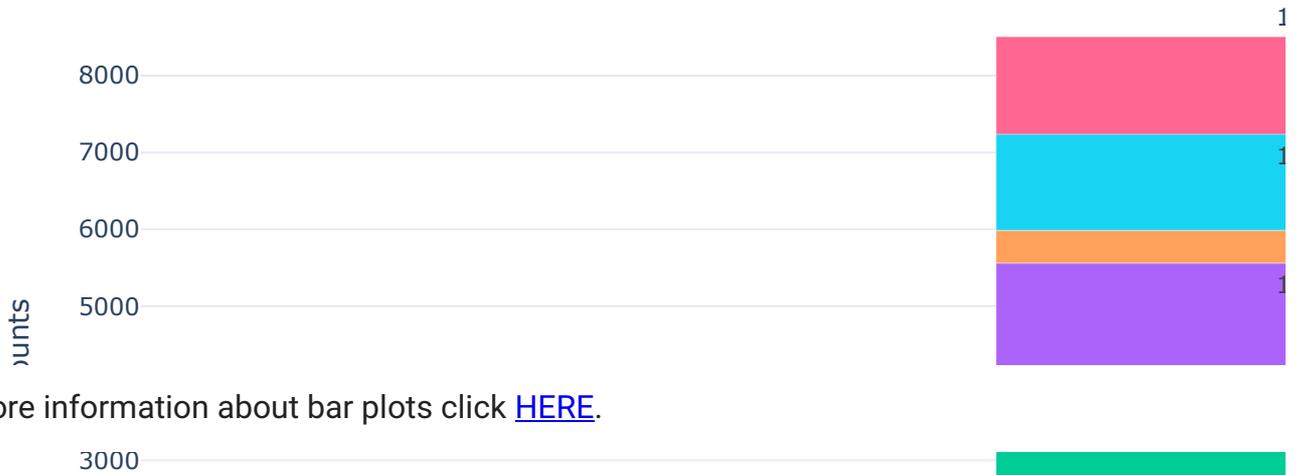
Below, we turn this into a stacked bar chart.

```

1 churn_edu_fig = go.Figure()
2
3 # Loop over elements in grps array object and nms list object
4 for i in range(len(values)):
5     churn_edu_fig.add_trace(go.Bar(
6         x=groups,
7         y=values[i],
8         text=values[i],
9         textposition='outside',
10        name=levels[i],
11    ))
12
13 churn_edu_fig.update_layout(title='Educational level frequencies by customer group')
14
15 # Using alternate dictionary syntax
16 churn_edu_fig.update_layout({'xaxis':{'title':'Customer group'},
17                             'yaxis':{'title':'Counts'},
18                             'barmode':'stack'})
19
20 churn_edu_fig.show()

```

Educational level frequencies by customer group



For more information about bar plots click [HERE](#).

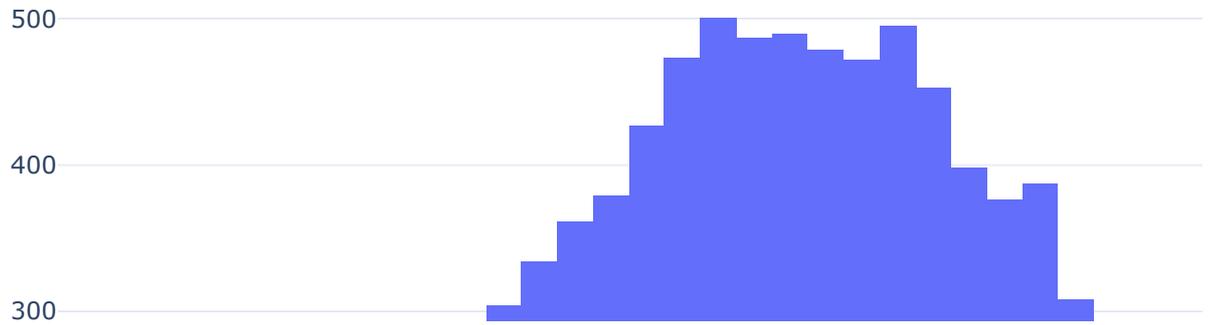
▼ HISTOGRAM



A histogram is used for continuous numerical variables. By creating bins, we can count how many times a value in that interval appears. Binning is a method of dividing the range of values for the variable up into intervals and counting how many times values *fall within* each interval.

Histograms are great at showing us the distribution of the data. Below, we plot a histogram of the age of our customers. This time, we make use of the plotly express library. The first argument is the dataframe object and the `x` argument is assigned to the `Customer_Age` column in the dataframe.

```
1 age_hist = px.histogram(df,
2                       x='Customer_Age')
3 age_hist.show()
```



The bin size was set by default at one year increments. Note that there are no gaps between the bars as with a bar chart. This indicates the fact that the variable is continuous numerical.

200

We can create a **stacked** histogram by using the `color` argument to point to another categorical variable, `Attrition_Flag` in this case. The `labels` argument is set to a dictionary. It can be used to change text aspects of the axes labels of the plot. First, we see a plot without this argument and then with it to see the differences.

```

1 age_hist = px.histogram(
2     df,
3     x='Customer_Age',
4     color='Attrition_Flag',
5     title='Histogram of customer ages',
6     opacity=0.7,
7     marginal='rug'
8 )
9
10 age_hist.show()

```

Histogram of customer ages



Note the `Customer_Age` title of the `x` axis, taken from the columns name. The dictionary is used below to replace this value. This is only available in the express module.

```

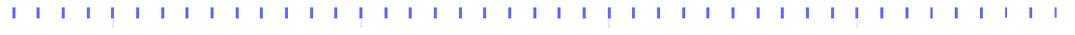
1 age_hist = px.histogram(
2     df,
3     x='Customer_Age',
4     color='Attrition_Flag',
5     title='Histogram of customer ages',
6     opacity=0.7,
7     marginal='rug',
8     labels={
9         'Attrition_Flag': 'Customer group',
10        'Customer_Age': 'Customer age'
11    }
12 )
13
14 age_hist.show()

```

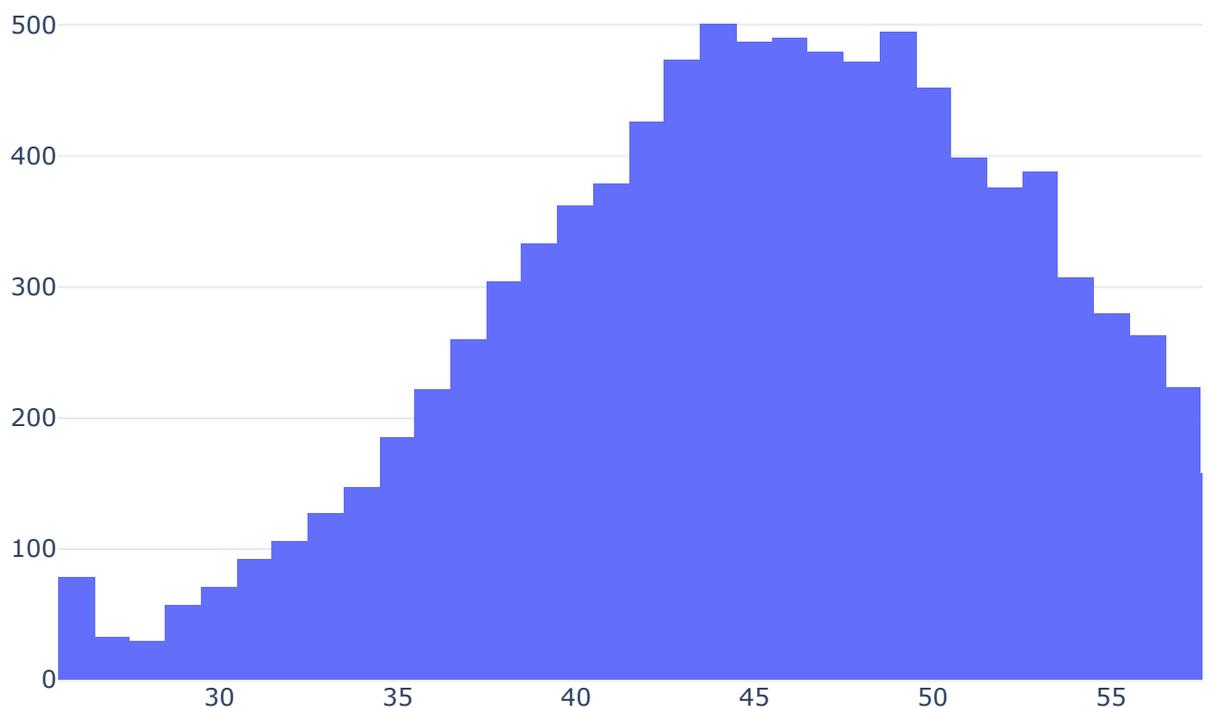
The rug plot shows the actual values as small marks at the top of the plot. Since the values overlap (many customers with the same age), we see that the rug plot does not give us a good idea of the distribution of the data.



The graph objects module of plotly provides more flexibility.



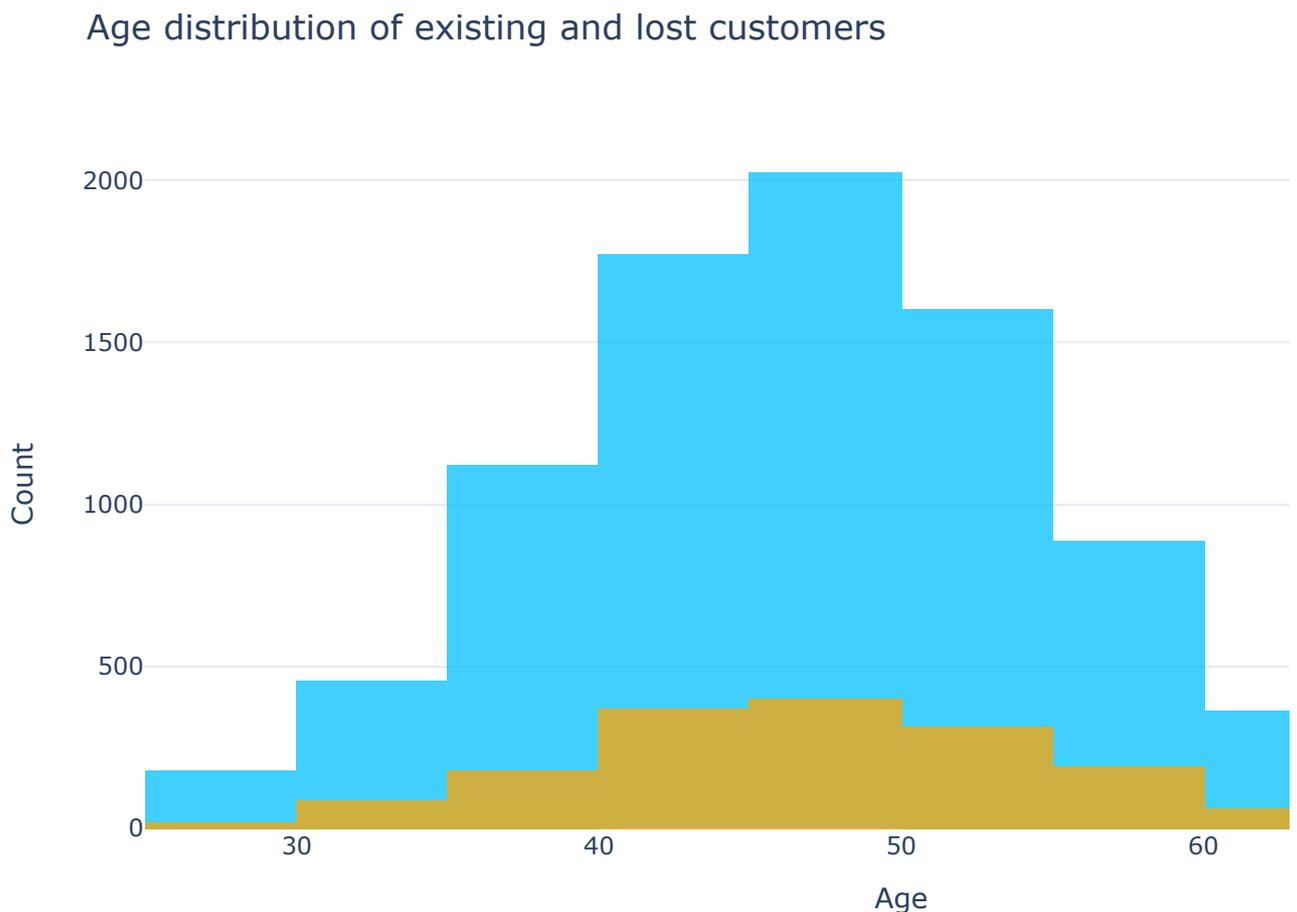
```
1 age_hist = go.Figure()
2
3 age_hist.add_trace(go.Histogram(
4     x=df.Customer_Age
5 ))
6
7 age_hist.show()
```



Let's look at the age distribution for each of the customer groups. Using the `barmode='overlay'` option, we create a **non-stacked** histogram. Below we set a start and end interval for the x axis and also a bins size of five year increments.

```
1 ages_cust = go.Figure()
```

```
2
3 ages_cust.add_trace(go.Histogram(
4     x=df[df.Attrition_Flag == 'Existing Customer']['Customer_Age'],
5     name='Existing customer',
6     marker_color='deepskyblue',
7     xbins=dict(start=25,
8                 end=80,
9                 size=5)
10 ))
11 ages_cust.add_trace(go.Histogram(
12     x=df[df.Attrition_Flag == 'Attrited Customer']['Customer_Age'],
13     name='Lost customer',
14     marker_color='orange',
15     xbins=dict(start=25,
16                 end=80,
17                 size=5)
18 ))
19
20 ages_cust.update_layout(barmode='overlay',
21                          title='Age distribution of existing and lost customers
22                          xaxis=dict(title='Age'),
23                          yaxis=dict(title='Count'))
24
25 ages_cust.update_traces(opacity=0.75)
26
27 ages_cust.show()
```



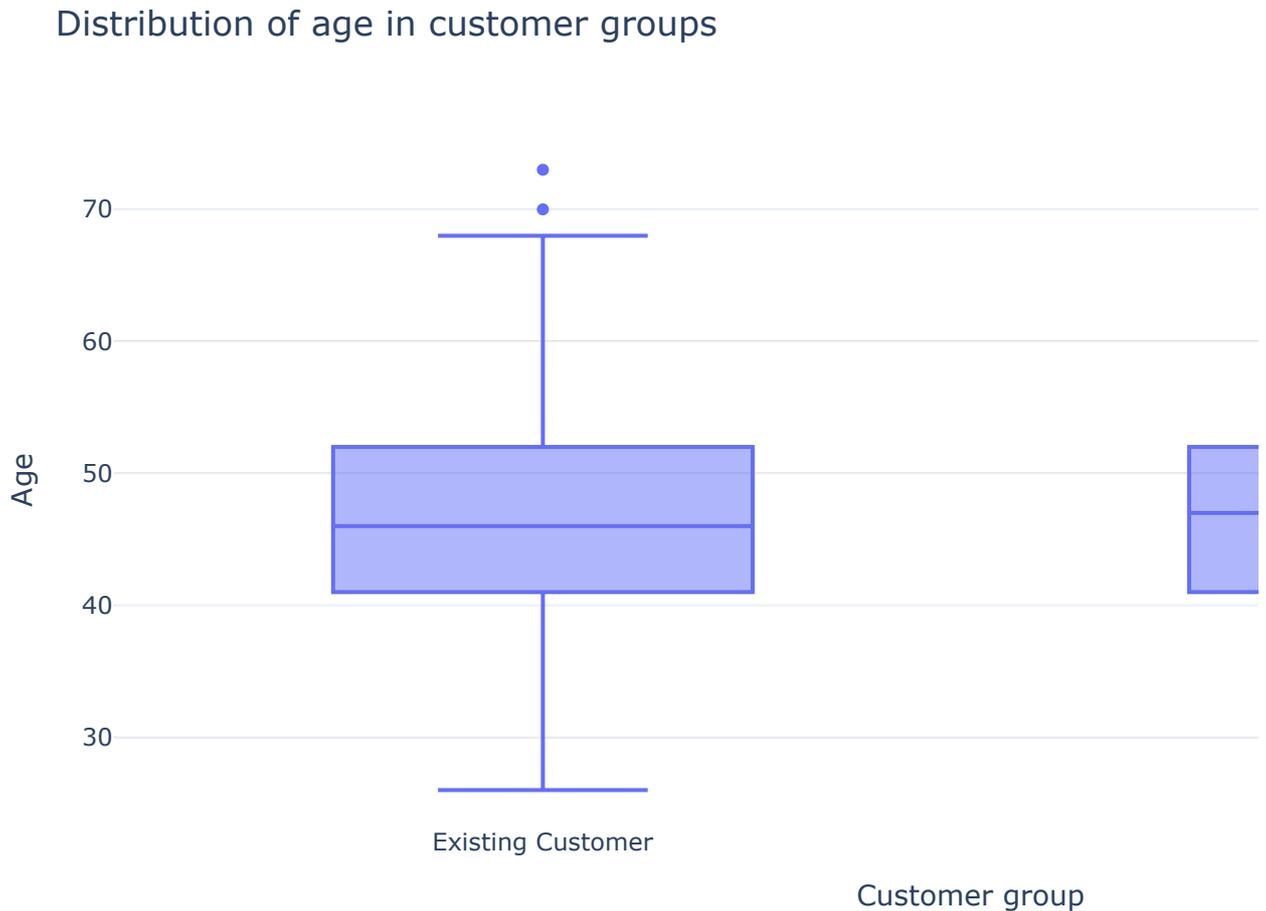
For more information on histograms from plotly click [HERE](#).

▼ BOX PLOT

Box-and-whisker plots are another visual representation of the distribution of a continuous numerical variable.

Below, we create a box plot of the ages of each customer group.

```
1 # Simple box plots using express
2 ages_churn_box_px = px.box(
3     df,
4     x='Attrition_Flag',
5     y='Customer_Age',
6     title='Distribution of age in customer groups',
7     labels={'Customer_Age':'Age', 'Attrition_Flag':'Customer group'})
8 ages_churn_box_px.show()
```



With the `Box` function in the `graph_objects` module, we can take more control over each trace

```

1 # Extracting list objects
2 exis_age = df[df.Attrition_Flag == 'Existing Customer']['Customer_Age'].to_list()
3 churn_age = df[df.Attrition_Flag == 'Attrited Customer']['Customer_Age'].to_list()

```

Now we create separate traces for each group. Code comments highlight the arguments.

```

1 # Adding separate traces and configuration
2 ages_churn_box = go.Figure()
3
4 ages_churn_box.add_trace(go.Box(
5     y=exis_age,
6     name='Existing customers',
7     marker_color='green',
8     boxmean=True, # Add a mean as a dotted line
9     boxpoints='all' # Add all the age values as dots next to the box
10 ))
11
12 ages_churn_box.add_trace(go.Box(
13     y=churn_age,
14     name='Lost customers',
15     marker_color='red',
16     boxmean='sd', # Add a mean and standard deviation as dotted lines
17     boxpoints='all'
18 ))
19
20 ages_churn_box.update_layout(title='Distribution of ages',
21                             xaxis={'title':'Group'},
22                             yaxis={'title':'Count'})
23
24 ages_churn_box.show()

```

Distribution of ages



You can learn more about box-and-whisker plots [HERE](#).

▼ SCATTER PLOT



Scatter plots allow us to view the difference between observations with respect to continuous numerical variables. If we restrict ourselves to two variables, each dot on a plane with an x and a y axis can represent the values for a pair of continuous numerical variables for each observation.

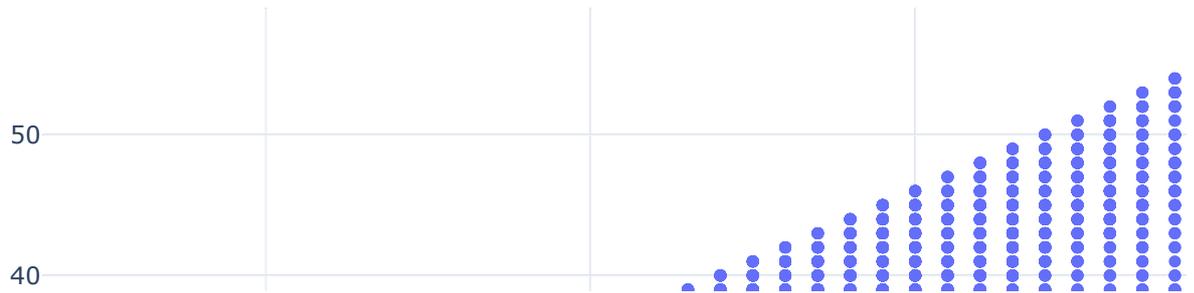


- Create a scatter plot of age (x axis) vs systolic blood pressure (y axis).

```

1 age_mob = go.Figure()
2
3 age_mob.add_trace(go.Scatter(
4     x=df.Customer_Age,
5     y=df.Months_on_book,
6     mode='markers'
7 ))
8
9 age_mob.update_layout(title="Age vs month on books",
10                        xaxis=dict(title="Age"),
11                        yaxis=dict(title="Months"))
12
13 age_mob.show()
```

Age vs month on books



We can add a third variable in the form of the size of the marker. Below, we do just that and compare the age and months on books of the customers. To this, we add the number of dependents as the size of the marker. We also split the customers by the sample space elements of the `Attrition_Flag` variable.



We can even further and add box-and-whisker plots, together with a linear model (using ordinary-least-squares).

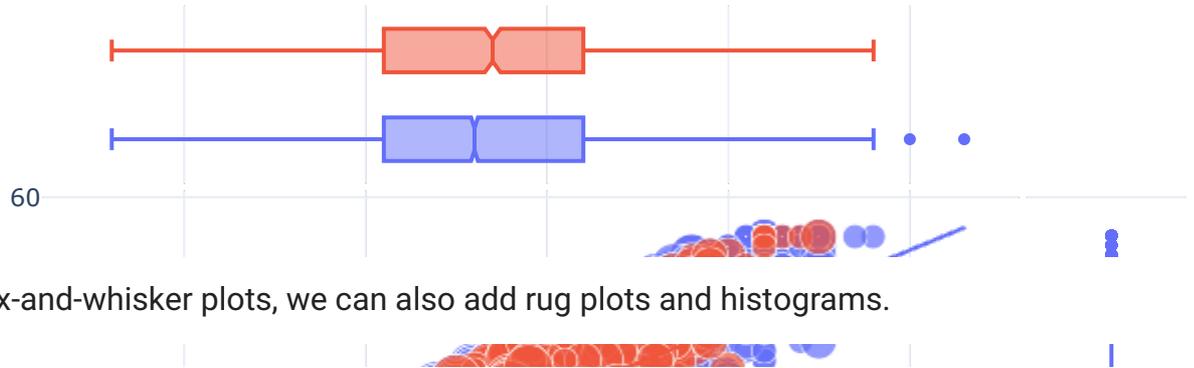
```

1 # Box and whisker plots of the variables
2 age_mob_group_px = px.scatter(
3     df,
4     x='Customer_Age',
5     y='Months_on_book',
6     size='Dependent_count', # Determines size of markers
7     color='Attrition_Flag', # Group by this variable
8     marginal_y='box',
9     marginal_x='box',
10    trendline='ols',
11    title='Age vs months on books',
12    labels={'Customer_Age':'Age', 'Months_on_book':'Months'}) # Over-write colu
13 age_mob_group_px.show()

```

/usr/local/lib/python3.7/dist-packages/statsmodels/tools/_testing.py:19: FutureWarning: pandas.util.testing is deprecated. Use the functions in the public API at pandas.testing instead.

Age vs months on books



Instead of box-and-whisker plots, we can also add rug plots and histograms.

```

1 # Rug plot and histogram as marginal plots
2 px.scatter(
3     df,
4     x='Customer_Age',
5     y='Months_on_book',
6     size='Dependent_count', # Determines size of markers
7     color='Attrition_Flag',
8     marginal_y='histogram',
9     marginal_x='rug',
10    trendline='ols',
11    title='Comparing age vs months on books for each of the two groups',
12    labels={'Customer_Age': 'Age', 'Months_on_book': 'Months'}).show()

```

Comparing age vs months on books for each of the two groups

Instead of the size of the marker as visual indicator of the third variable, we can also use color. Below, we also add the `facet_col` argument. This creates individual plots based on the sample space elements of a categorical variable and positions them as columns.

- Separate scatter plots
- *Third dimension* by color scale

```

1 px.scatter(
2     df,
3     x='Customer_Age',
4     y='Months_on_book',
5     color='Dependent_count',
6     facet_col='Attrition_Flag',
7     trendline='ols',
8     title='Sperate scatter plots per group',
9     labels={'Attrition_Flag':'Group', 'Months_on_book':'Months', 'Customer_Age':
10    color_continuous_scale=px.colors.sequential.Viridis).show();

```

(Hover on the trendline to see the regression model and the coefficient of determination, R^2 . We will learn more about this when we look at linear regression.)

You can learn more about scatter plots [HERE](#).

▼ TIME SERIES DATA

Time series contain dates or sequential data of some form. Here, we use the Australian rainfall data set. We start by investigating the difference in maximum daily temperatures between Darwin and Hobart.

We start by using the `info` method to learn more about the dataframe object.

```
1 rain.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 145460 entries, 0 to 145459
Data columns (total 23 columns):
 #   Column                Non-Null Count  Dtype
---  -
 0   Date                  145460 non-null object
 1   Location              145460 non-null object
 2   MinTemp               143975 non-null float64
 3   MaxTemp               144199 non-null float64
 4   Rainfall              142199 non-null float64
 5   Evaporation           82670 non-null  float64
 6   Sunshine              75625 non-null  float64
 7   WindGustDir           135134 non-null object
 8   WindGustSpeed         135197 non-null float64
 9   WindDir9am           134894 non-null object
10  WindDir3pm            141232 non-null object
11  WindSpeed9am          143693 non-null float64
12  WindSpeed3pm          142398 non-null float64
13  Humidity9am           142806 non-null float64
14  Humidity3pm           140953 non-null float64
15  Pressure9am           130395 non-null float64
16  Pressure3pm           130432 non-null float64
17  Cloud9am              89572 non-null  float64
18  Cloud3pm              86102 non-null  float64
19  Temp9am               143693 non-null float64
20  Temp3pm               141851 non-null float64
21  RainToday             142199 non-null object
22  RainTomorrow          142193 non-null object
dtypes: float64(16), object(7)
memory usage: 25.5+ MB
```

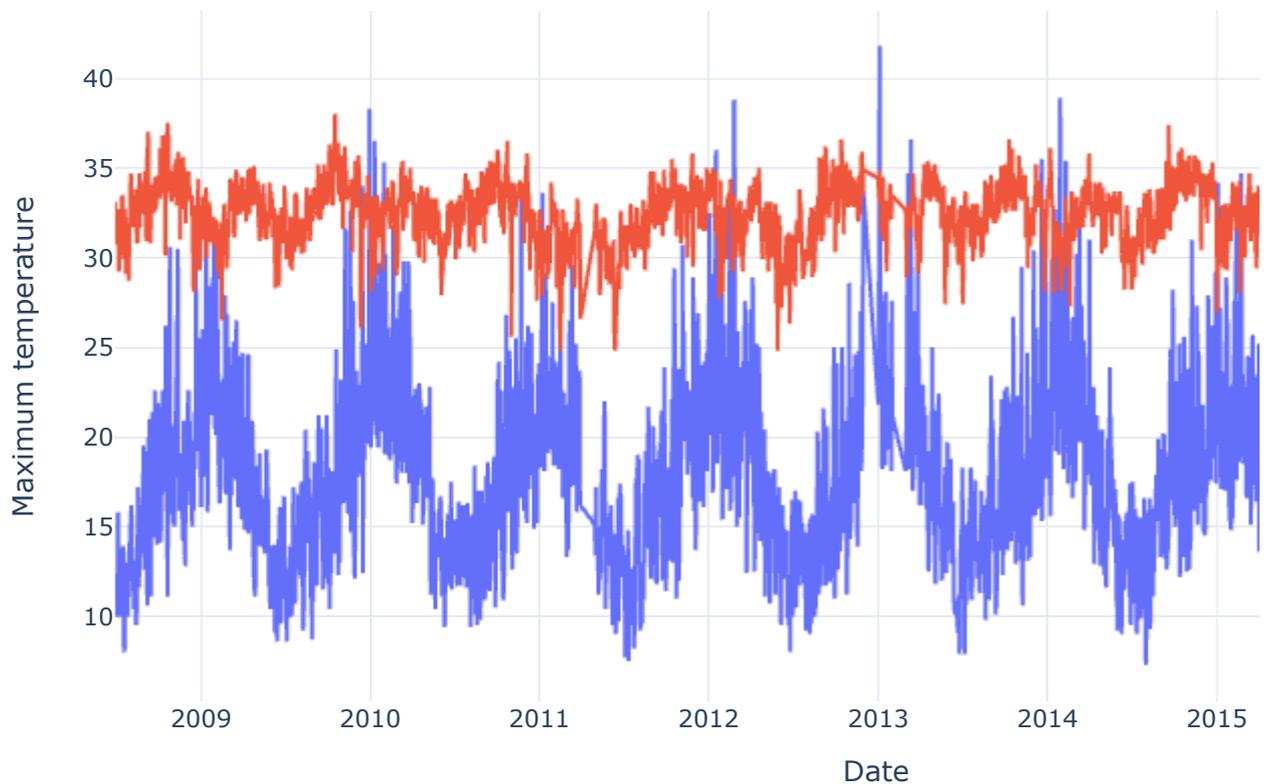
The data points for `MaxTemp`, the maximum temperature, are connect by line segments using the `line` function in the `express` module.

```

1 px.line(
2     rain[(rain.Location == 'Darwin') | (rain.Location == 'Hobart')],
3     x='Date',
4     y='MaxTemp',
5     color='Location',
6     title='Maximum temperature in Darwin and Hobart between 2009 and 2017',
7     labels={'MaxTemp':'Maximum temperature'}
8 )

```

Maximum temperature in Darwin and Hobart between 2009 and 2017



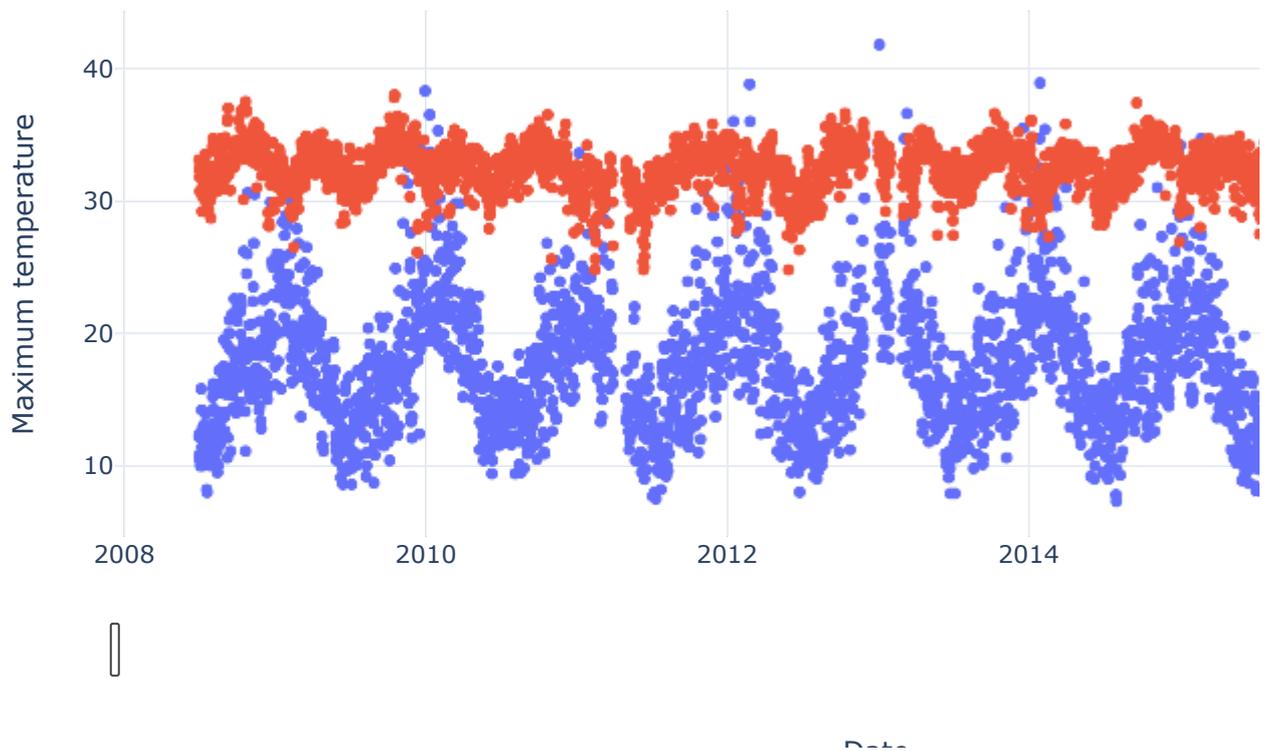
We can also view the individual data points as a scatter plot. We also add a range slider at the bottom of the plot.

```

1 px.scatter(
2     rain[(rain.Location == 'Darwin') | (rain.Location == 'Hobart')],
3     x='Date',
4     y='MaxTemp',
5     color='Location',
6     title='Maximum temperature in Hobart and Darwin between 2009 and 2017',
7     labels={'MaxTemp':'Maximum temperature'}
8 ).update_xaxes(rangeslider_visible=True)

```

Maximum temperature in Hobart and Darwin between 2009 and 2017



Lastly, let's look at the difference in daily rainfall between Perth in the West and Sydney in the East.

```

1 px.scatter(
2     rain[(rain.Location == 'Perth') | (rain.Location == 'Sydney')],
3     x='Date',
4     y='Rainfall',
5     color='Location',
6     title='Daily rainfall in Sydney and Perth between 2009 and 2017'
7 ).update_axes(rangeslider_visible=True)

```

Daily rainfall in Sydney and Perth between 2009 and 2011



▼ CONCLUSION



There is so much more to learn about plotly and we will see new types of plots in later notebooks. Visit the homepage [HERE](#). For references on all the arguments click [HERE](#).

||

||

1

Date

✓ 0s completed at 13:56

● ×