

Exploratory Data Analysis Lab

Estimated time needed: 30 minutes

In this module you get to work with the cleaned dataset from the previous module.

In this assignment you will perform the task of exploratory data analysis. You will find out the distribution of data, presence of outliers and also determine the correlation between different columns in the dataset.

Objectives

In this lab you will perform the following:

- Identify the distribution of data in the dataset.
- Identify outliers in the dataset.
- Remove outliers from the dataset.
- Identify correlation between features in the dataset.

Hands on Lab

Import the pandas module.

```
In [96]: import pandas as pd
import matplotlib inline
import plotly.express as px
import plotly.figure_factory as ff
import numpy as np
import seaborn as sns
```

Load the dataset into a dataframe.

```
In [2]: df = pd.read_csv("https://cf-courses-data.s3.us.cloud-object-storage.appdomain.cloud/IBM-DA0321EN-SkillsNetwork/LargeData/m2_survey_data.csv")
```

Distribution

Determine how the data is distributed

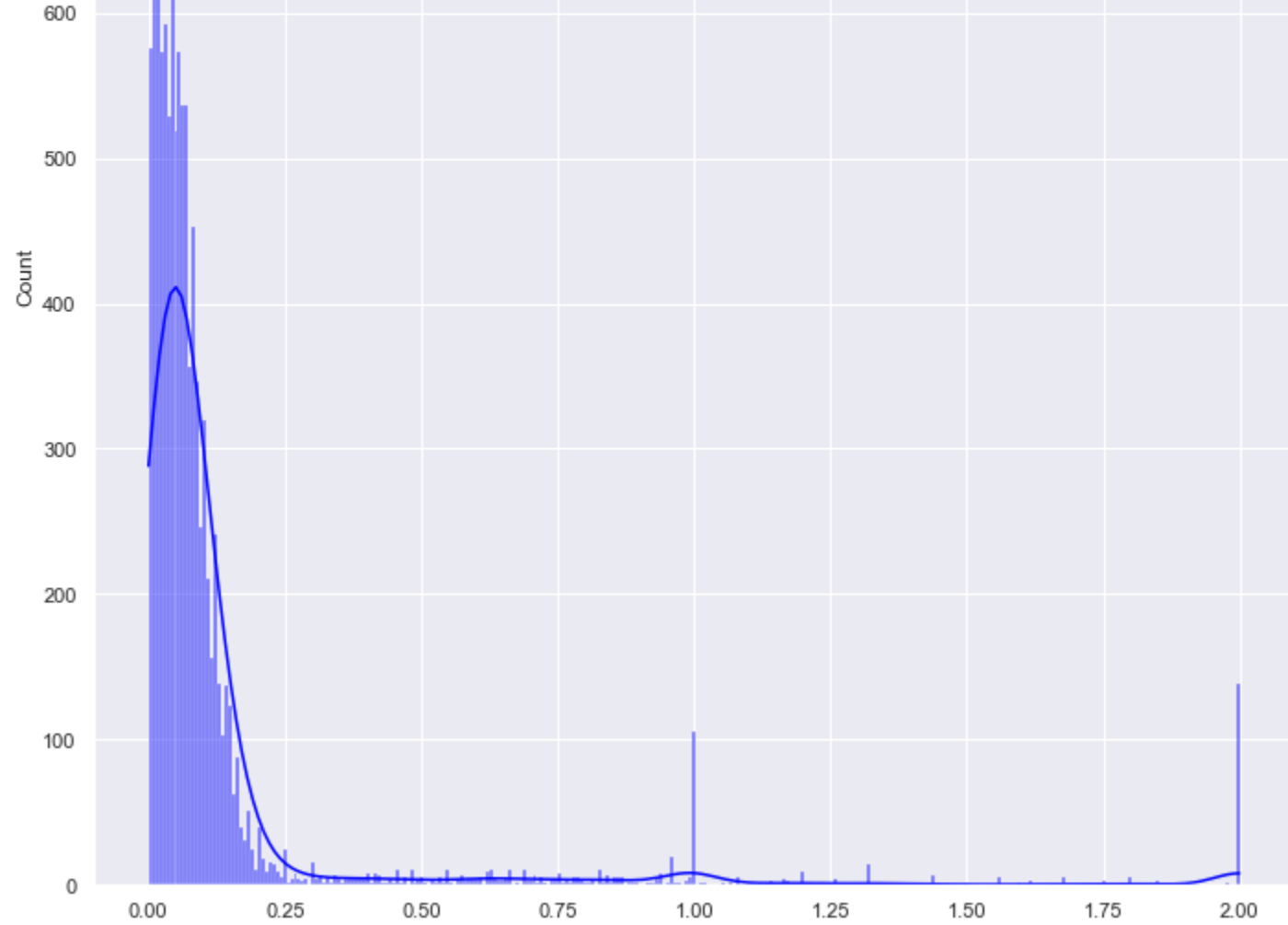
The column `ConvertedComp` contains Salary converted to annual USD salaries using the exchange rate on 2019-02-01.

This assumes 12 working months and 50 working weeks.

Plot the distribution curve for the column `ConvertedComp`.

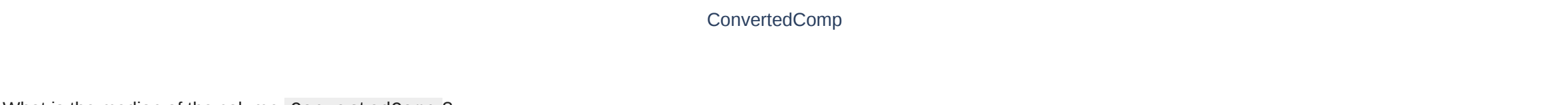
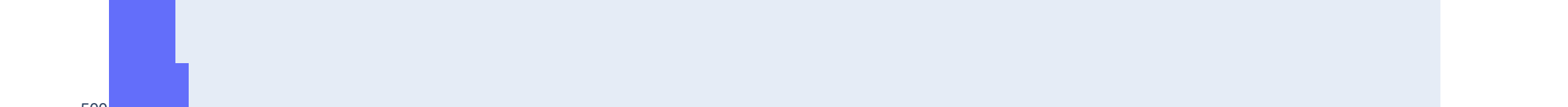
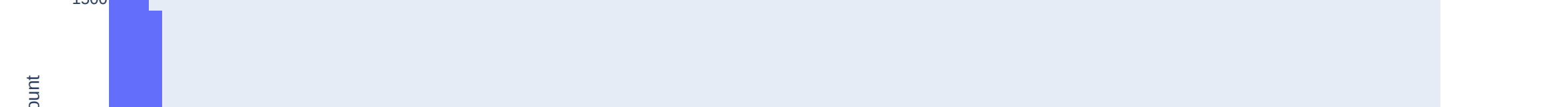
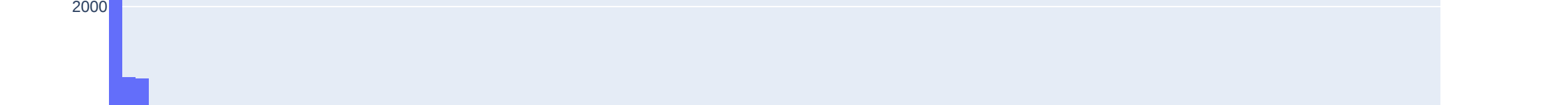
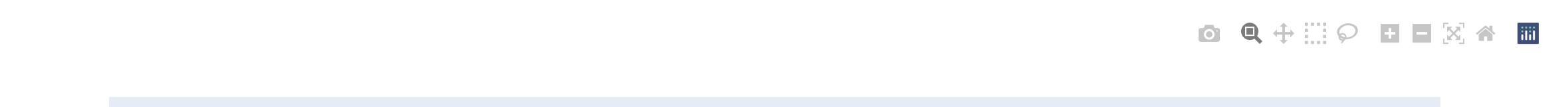
```
In [21]: sns.set_theme()
sns.displot(df['ConvertedComp'], kde = True, color='blue', height = 10)
```

```
Out[21]: <seaborn.axisgrid.FacetGrid at 0x21b40fbabe0>
```



Plot the histogram for the column `ConvertedComp`.

```
In [4]: # your code goes here
fig = px.histogram(df, x='ConvertedComp')
fig.show()
```



What is the median of the column `ConvertedComp`?

```
In [5]: # your code goes here
df.ConvertedComp.median()
```

```
Out[5]: 57745.0
```

```
In [6]: df.Age.median()
```

```
Out[6]: 29.0
```

How many responders identified themselves only as a `Man`?

```
In [7]: # your code goes here
df.Gender.value_counts()
```

```
Out[7]: Man                10480
Woman                731
Non-binary, genderqueer, or gender non-conforming         63
Man;Non-binary, genderqueer, or gender non-conforming      26
Woman;Non-binary, genderqueer, or gender non-conforming    14
Woman;Man                9
Woman;Man;Non-binary, genderqueer, or gender non-conforming 2
Name: Gender, dtype: int64
```

Find out the median `ConvertedComp` of responders identified themselves only as a `Woman`?

```
In [8]: # your code goes here
df.loc[df['Gender'] == 'Woman', ['ConvertedComp']].median()
```

```
Out[8]: ConvertedComp    57708.0
dtype: float64
```

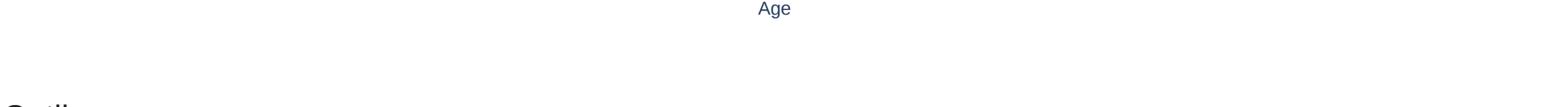
Give the five number summary for the column `Age`?

```
In [9]: # your code goes here
df['Age'].describe()
```

```
Out[9]: count    11111.000000
mean         30.778895
std          7.533886
min          16.000000
25%          25.000000
50%          29.000000
75%          35.000000
max          99.000000
Name: Age, dtype: float64
```

Plot a histogram of the column `Age`.

```
In [10]: # your code goes here
histAge = px.histogram(df, x="Age")
histAge.show()
```



Outliers

Finding outliers

Find out if outliers exist in the column `ConvertedComp` using a box plot?

```
In [11]: # your code goes here
sns.boxplot(x=df['ConvertedComp'])
sns.boxplot(x=df, x='ConvertedComp', points='all')
boxplot.show()
```

