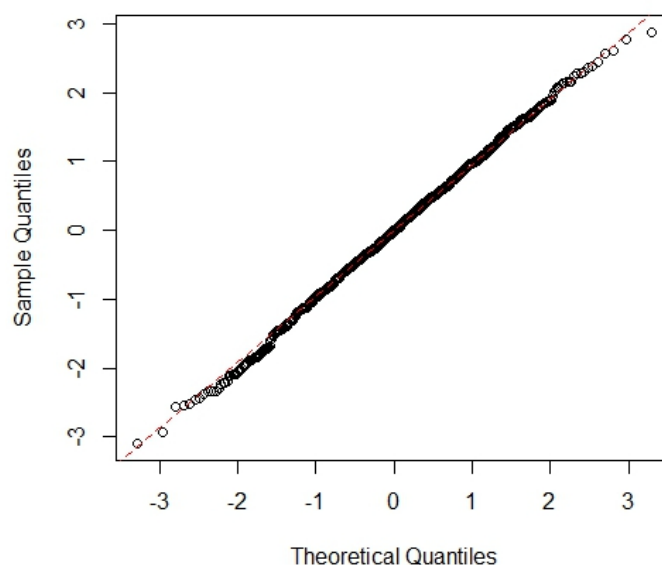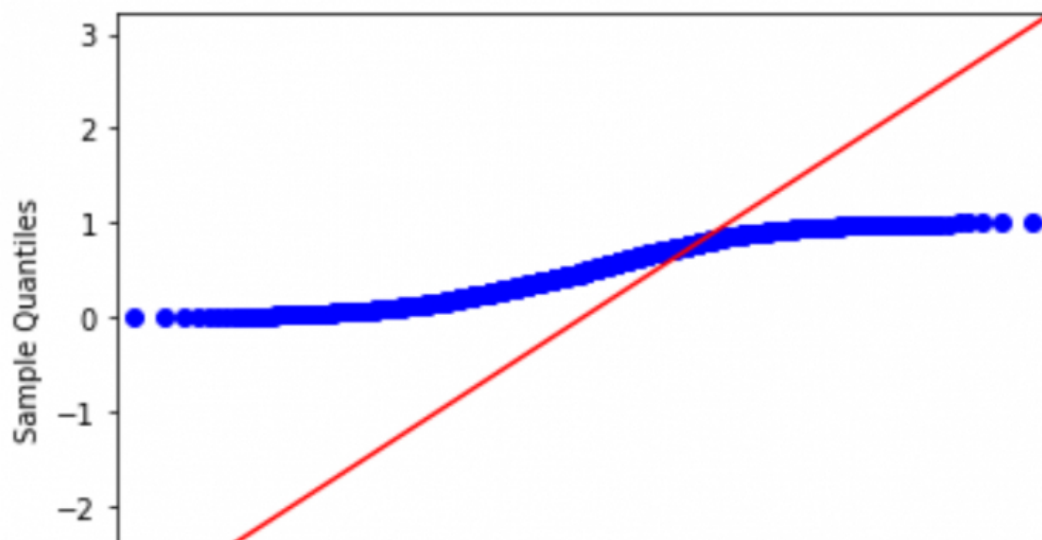# QQ Plot (Quantile- Quantile Plot)

## 1.Introduction

- A Q-Q plot is a scatterplot created by plotting two sets of quantiles against one another.
- If both sets of quantiles came from the same distribution, we should see the points forming a line that's roughly straight.
- x axis displays the theoretical quantiles.
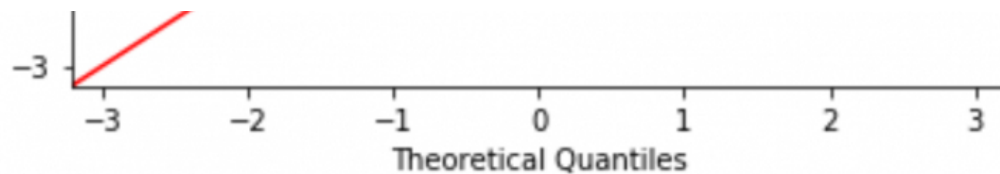- y axis displays the actual data

**QQ Plot for a normally distributed sample**



**QQ Plot of samples which are not normally distributed**

$-3$ axis with values $-3, -2, -1, 0, 1, 2, 3$

**Theoretical Quantiles**

## 2.Where the visualization can be used

- Some machine learning models like linear and logestic regression assume that the variables are normally distributed.
- The normal distributed varaiables may boost the machine learning algorithm performance.
- So we can use QQ plot to check a set of observations are normally distributed.

## 3.Python Code for QQ Plot

### Libraries used

- pandas: It offers data structures and operations for manipulating numerical tables.
- numpy: Python library used for working with arrays
- matplotlib: Used for visualization.
- scipy: Is a library that uses NumPy for more mathematical functions.

In [43]:
```python
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
%matplotlib inline
import scipy.stats as stats
```

### Dataset Used

- We are using titanic data set.
- For our purpose we are only taking 3 columns Age,Fare and surivive.
- Top 5 datas are shown below

In [44]:
```python
data=pd.read_csv('titanic.csv',usecols=['Age','Fare','Survived'])
data.head()
```

Out[44]:

|   | Survived | Age  | Fare    |
|---|----------|------|---------|
| 0 | 0        | 22.0 | 7.2500  |
| 1 | 1        | 38.0 | 71.2833 |
| 2 | 1        | 26.0 | 7.9250  |
| 3 | 1        | 35.0 | 53.1000 |
| 4 | 0        | 35.0 | 8.0500  |

### Data preprocessing

- When we check for null values, we can find that the age contains 177 null values.
- We are removing rows with null value.

In [45]:
```python
data.isnull().sum()
```

Out[45]:
```
Survived      0
Age         177
Fare          0
dtype: int64
```

In [46]:
```python
df = data[pd.notnull(data['Age'])]
df.head()
```

Out[46]:

|   | Survived | Age  | Fare    |
|---|----------|------|---------|
| 0 | 0        | 22.0 | 7.2500  |
| 1 | 1        | 38.0 | 71.2833 |
| 2 | 1        | 26.0 | 7.9250  |
| 3 | 1        | 35.0 | 53.1000 |

|   | 4 | 0 | 35.0 | 8.0500 |

- We can see that all the null values are removed.

In [47]: 
```python
1 df.isnull().sum()
```

Out[47]: 
```
Survived    0
Age         0
Fare        0
dtype: int64
```
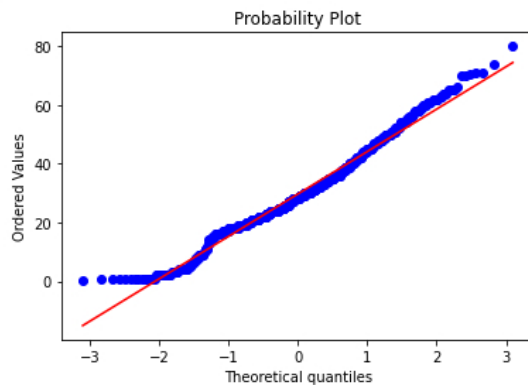
## QQ Plot

- A function is defined to draw the QQ plot.
- Two parameters are accepted the dataset and the varaiable.
- QQ plot is drawn using calling stats.probplot()

In [48]: 
```python
1 def diagnostic_plots(df, variable):
2
3     plt.subplot(1, 1, 1)
4     stats.probplot(df[variable], dist="norm", plot=plt)
5
6     plt.show()
```

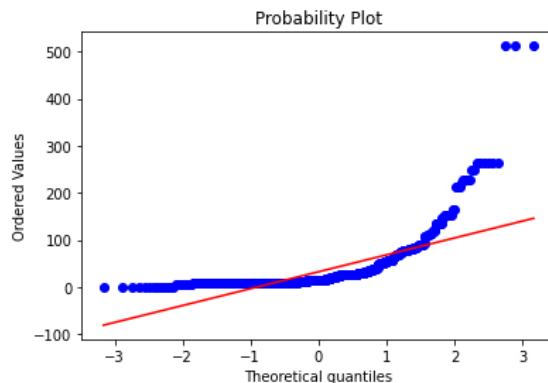- QQ plot for Age column.
- We can see that most of the points are near the line.

In [49]: 
```python
1 diagnostic_plots(df, 'Age')
```



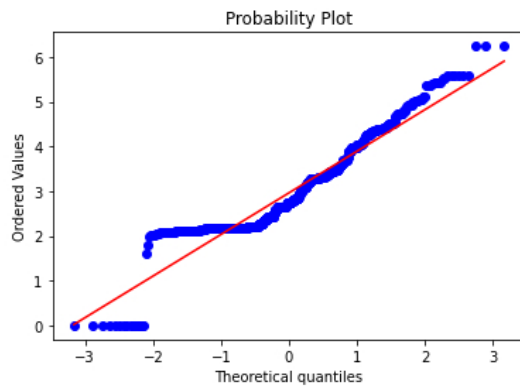- QQ plot for Fare column.
- We can see that most of the points are not on the red line, so we can say it is not a linear distributed.

In [50]: 
```python
1 diagnostic_plots(data, 'Fare')
```



- If a variable is not normally distributed, sometimes it is possible to find a mathematical transformation.
- One of such transformation is Logarithmic transformation.
- Here we are doing this logarithmic transformation to the 'fare'.
- After the transformation we can see that it is better than the last one.

```
1  data['Log_Fare']=np.log(data['Fare']+1)
2  diagnostic_plots(data,'Log_Fare')
3
```



Probability Plot

## 4. R code for QQ plot

- Here we are using wine classification dataset.
- We are mainly considering hue and malicAcid columns for ploting.



- qqnorm() is used to draw qq plot on R
- qqline() is used to draw the 45 degree line.
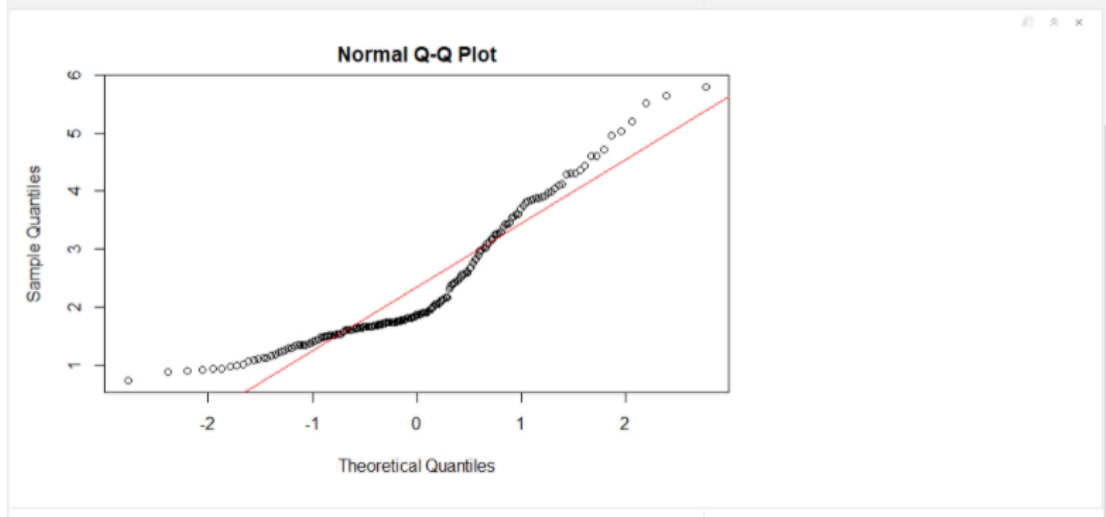- We can see that hue varaiable is almost linearly distributed.

```{r}
qqnorm(data$hue)
qqline(data$hue, col = 'red')
```



Normal Q-Q Plot

- We can see that malicAcid varaiable is not linearly distributed. As most of the points doesnt pass through the line.

```{r}
qqnorm(data$malicAcid)
```

```
qqline(data$malicAcid, col = 'red')
```

**Normal Q-Q Plot**



## 5. Purpose of the visualization

- In most cases, this plot is used to determine whether or not a set of data follows a normal distribution.
- All point lie on or close to straight line at an angle of 45 degree from x – axis. It indicates that the samples have similar distributions.

# Thank You

In [ ]:  1