

6000 level

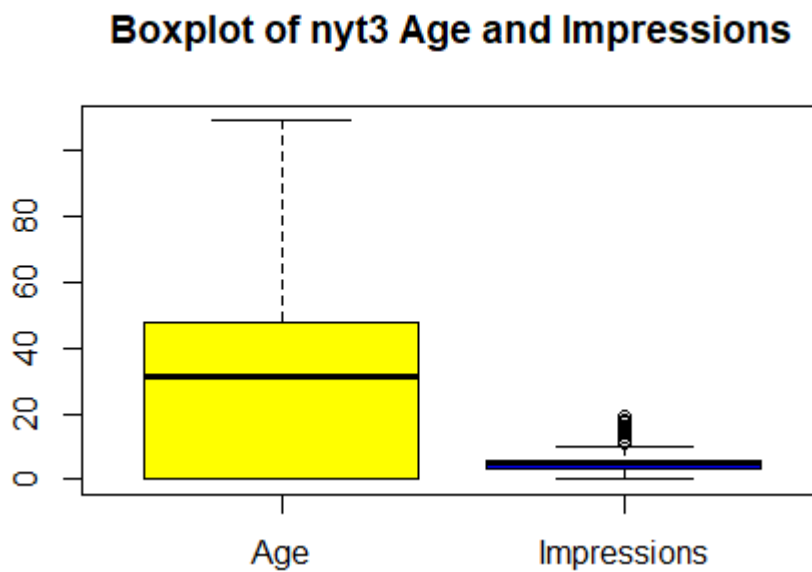
NAME: OGOCHUKWU NWACHUKWU GLORIA

RIN: 662058987

ASSIGNMENT 3

Question 1 a

Boxplot for nyt3, nyt4, nyt5, nyt6, nyt7, nyt8 and nyt9 Age and Impressions



Age:

- Minimum Age: 0.00
- 1st Quartile (25th percentile): 0.00
- Median (50th percentile): 31.00
- Mean: 29.47
- 3rd Quartile (75th percentile): 48.00
- Maximum Age: 109.00

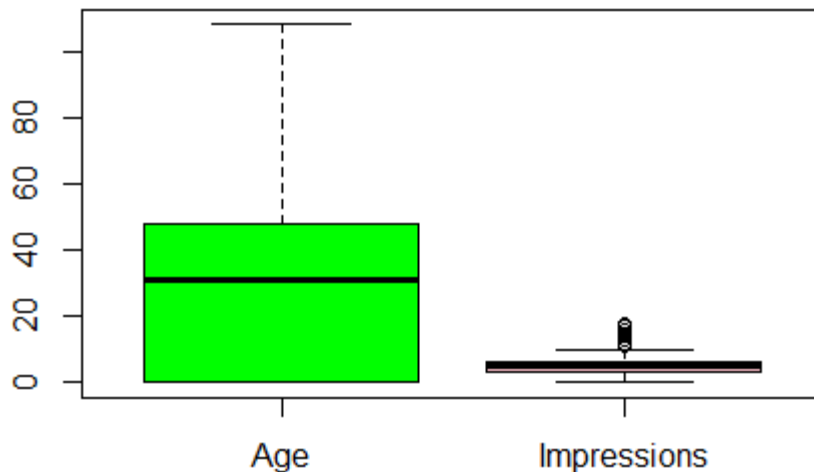
Impressions:

- Minimum Impressions: 0.000
- 1st Quartile (25th percentile): 3.000
- Median (50th percentile): 5.000
- Mean: 4.996
- 3rd Quartile (75th percentile): 6.000

- Maximum Impressions: 19.000

In Age, majority of the data falls within the range of 0 to 48, as indicated by the first and third quartiles. The median age is 31, and the mean age is approximately 29.47. In Impressions, most of the data is clustered between 3 and 6 impressions.

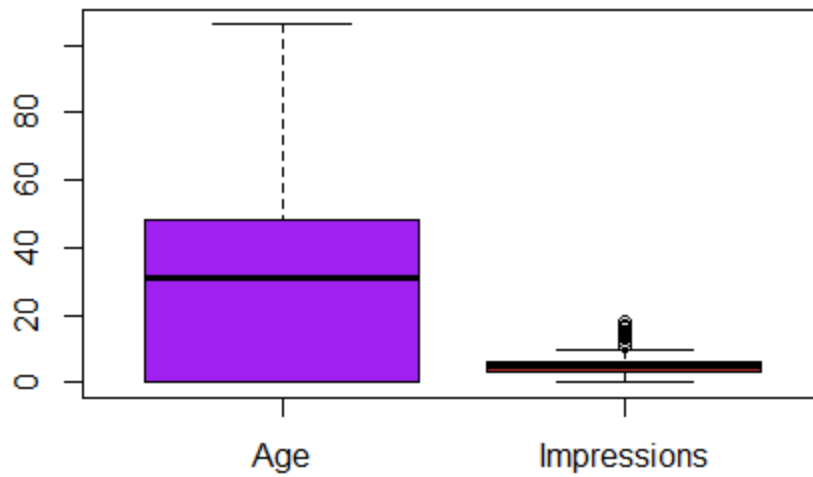
Boxplot of nyt4 Age and Impressions



- **Age:**
 - Min: 0.00
 - 1st Qu.: 0.00
 - Median: 31.00
 - Mean: 29.43
 - 3rd Qu.: 48.00
 - Max: 108.00
- **Impressions:**
 - Min: 0.000
 - 1st Qu.: 3.000
 - Median: 5.000
 - Mean: 5.004
 - 3rd Qu.: 6.000
 - Max: 18.000

The age distribution has a minimum of 0, a maximum of 108, and is centered around the median and mean of 31 and 29.43, respectively. Impressions range from 0 to 18, with a median of 5 and a slightly higher mean of 5.004 compared to nyt3.

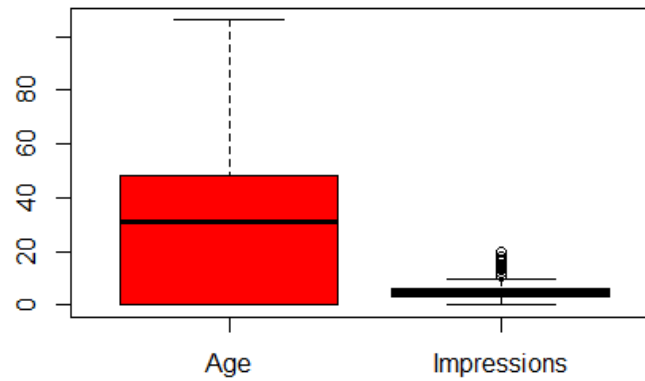
Boxplot of nyt5 Age and Impressions



- **Age:**
 - Min: 0.00
 - 1st Qu.: 0.00
 - Median: 31.00
 - Mean: 29.43
 - 3rd Qu.: 48.00
 - Max: 106.00
- **Impressions:**
 - Min: 0.000
 - 1st Qu.: 3.000
 - Median: 5.000
 - Mean: 4.999
 - 3rd Qu.: 6.000
 - Max: 18.000

Age has a median of 31, mean of 29.43, and a range from 0 to 106. Impressions also follow a similar pattern, with a median of 5.

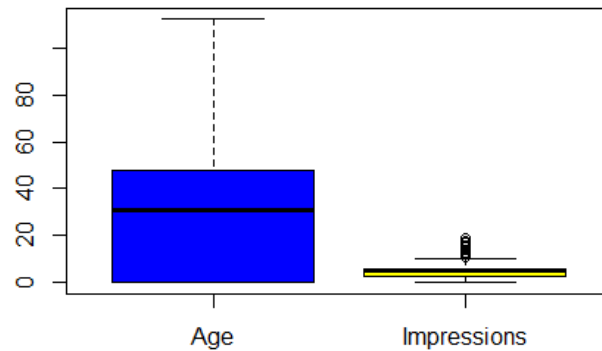
Boxplot of nyt6 Age and Impressions



- **Age:**
 - Min: 0.00
 - 1st Qu.: 0.00
 - Median: 31.00
 - Mean: 29.46
 - 3rd Qu.: 48.00
 - Max: 106.00
- **Impressions:**
 - Min: 0.000
 - 1st Qu.: 3.000
 - Median: 5.000
 - Mean: 4.995
 - 3rd Qu.: 6.000
 - Max: 20.000

Age has a median of 31, a mean of 29.46, and a range from 0 to 106. Impressions have a median of 5 and a mean of 4.995, similar to previous datasets.

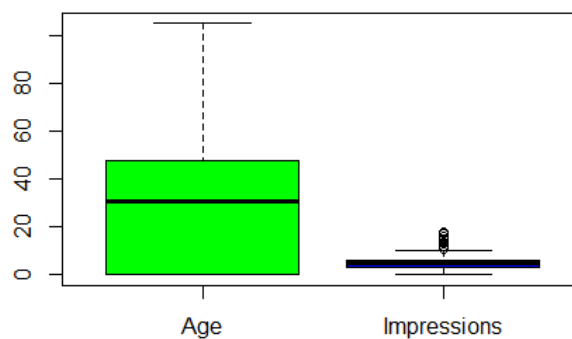
Boxplot of nyt7 Age and Impressions



- **Age:**
 - Min: 0.00
 - 1st Qu.: 0.00
 - Median: 31.00
 - Mean: 29.52
 - 3rd Qu.: 48.00
 - Max: 112.00
- **Impressions:**
 - Min: 0
 - 1st Qu.: 3
 - Median: 5
 - Mean: 5
 - 3rd Qu.: 6
 - Max: 19

Age has a median of 31, a mean of 29.52, and a range from 0 to 112. Impressions has a median of 5 and a mean of 5.

Boxplot of nyt8 Age and Impressions



Age:

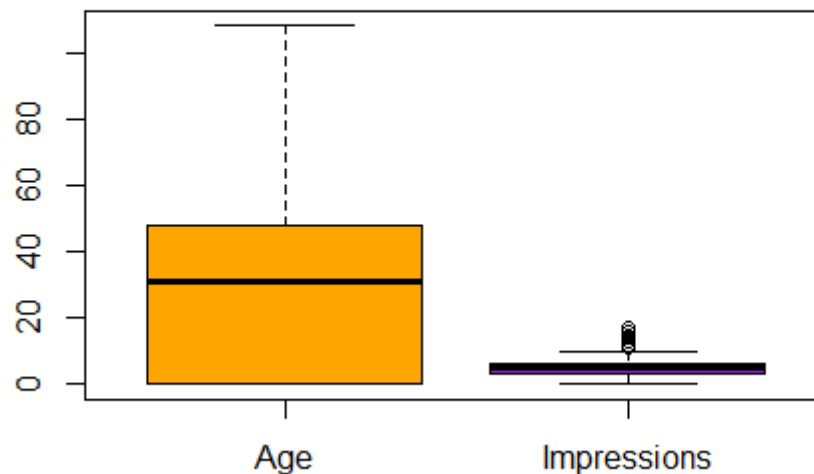
- Minimum (Min): 0.00
- 1st Quartile (1st Qu.): 0.00
- Median: 31.00
- Mean: 29.41
- 3rd Quartile (3rd Qu.): 48.00
- Maximum (Max): 105.00

Impressions:

- Minimum (Min): 0
- 1st Quartile (1st Qu.): 3
- Median: 5
- Mean: 5.001
- 3rd Quartile (3rd Qu.): 6
- Maximum (Max): 18

The Age variable the median (31.00) being slightly less than the mean (29.41). The 3rd quartile (48.00) is greater than the median, indicating the presence of potential outliers on the higher end of the data. In Impressions variable the median (5) is slightly less than the mean (5.001). The 3rd quartile (6) is greater than the median, suggesting possible outliers on the higher side of the data.

Boxplot of nyt9 Age and Impressions



- **Age:**
 - Min: 0.00
 - 1st Qu.: 0.00

- Median: 31.00
- Mean: 29.45
- 3rd Qu.: 48.00
- Max: 108.00
- **Impressions:**
 - Min: 0.000
 - 1st Qu.: 3.000
 - Median: 5.000
 - Mean: 4.998
 - 3rd Qu.: 6.000
 - Max: 17.000

Age statistics are in line with the previous datasets, with a median of 31, a mean of 29.45, and a range from 0 to 108. Impressions has a median of 5 and a mean just below 5. In the boxplot for nyt9, the distribution of age and impressions aligns with the previous datasets, with most data clustered at the lower end of the scale.

Conclusively, the datasets nyt3, nyt4, nyt5, nyt6, nyt7, nyt8 and nyt9 share similar statistical characteristics for age and impressions. They all have a concentration of data points with low age values and low impressions, while having some variations in means and max values. Most of the data is centered around the lower values. The boxplots visually represent these distributions, showing where the majority of the data falls and the presence of outliers. These visualizations help identify data patterns and potential outliers.

Question b

normality test using Anderson Darling test

The Anderson-Darling normality test assesses whether a given sample comes from a normally distributed population. In your analysis, you performed this test on both Age and Impressions for multiple datasets (nyt3, nyt4, nyt5, nyt6, nyt7, nyt8 and nyt9) and found that the p-values for all tests were less than $2.2e-16$. This extremely low p-value suggests that the data from all these datasets do not follow a normal distribution.

The results for nyt3, nyt4, nyt5, nyt6, nyt7, nyt8 and nyt9 Age and Impressions

For nyt3

data: nyt3\$Age

A = 12507, p-value < $2.2e-16$

data: nyt3\$Impressions
A = 4652.5, p-value < 2.2e-16

For nyt4

data: nyt4\$Age
A = 12645, p-value < 2.2e-16

data: nyt4\$Impressions
A = 4614.8, p-value < 2.2e-16

For nyt5

data: nyt5\$Age
A = 10541, p-value < 2.2e-16

data: nyt5\$Impressions
A = 3843.1, p-value < 2.2e-16

For nyt6

data: nyt6\$Age
A = 21752, p-value < 2.2e-16

data: nyt6\$Impressions
A = 8007.3, p-value < 2.2e-16

For nyt7

data: nyt7\$Age
A = 12842, p-value < 2.2e-16

data: nyt7\$Impressions
A = 4705, p-value < 2.2e-16

For nyt8

data: nyt8\$Age
A = 13297, p-value < 2.2e-16

data: nyt8\$Impressions
A = 4832.9, p-value < 2.2e-16

For nyt9

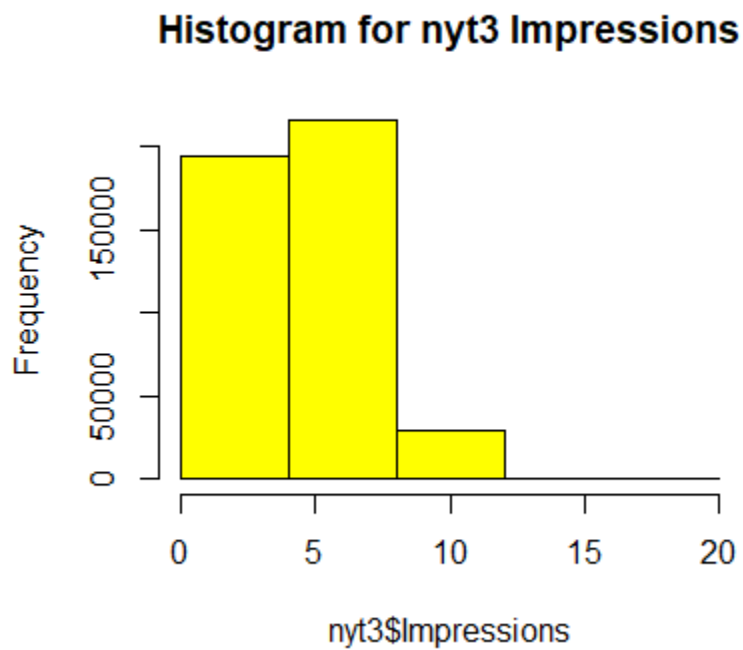
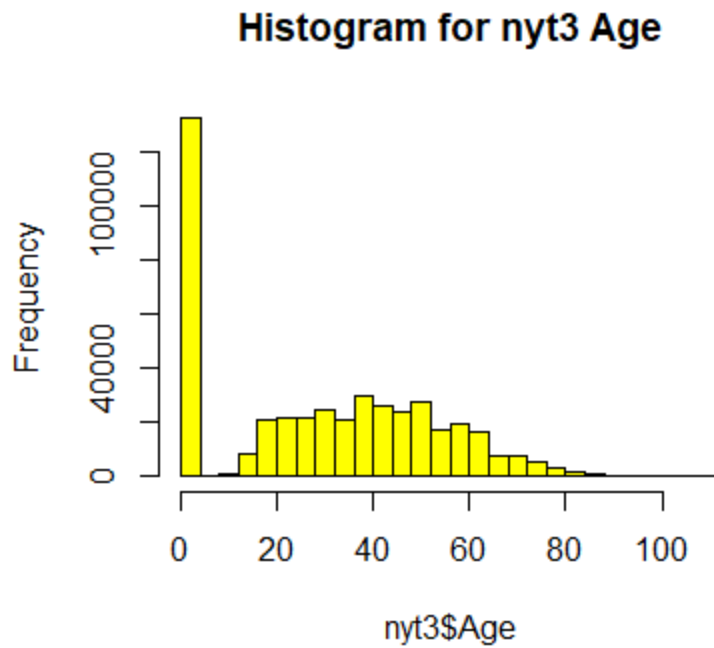
data: nyt9\$Age
A = 13033, p-value < 2.2e-16

data: nyt9\$Impressions
A = 4805.4, p-value < 2.2e-16

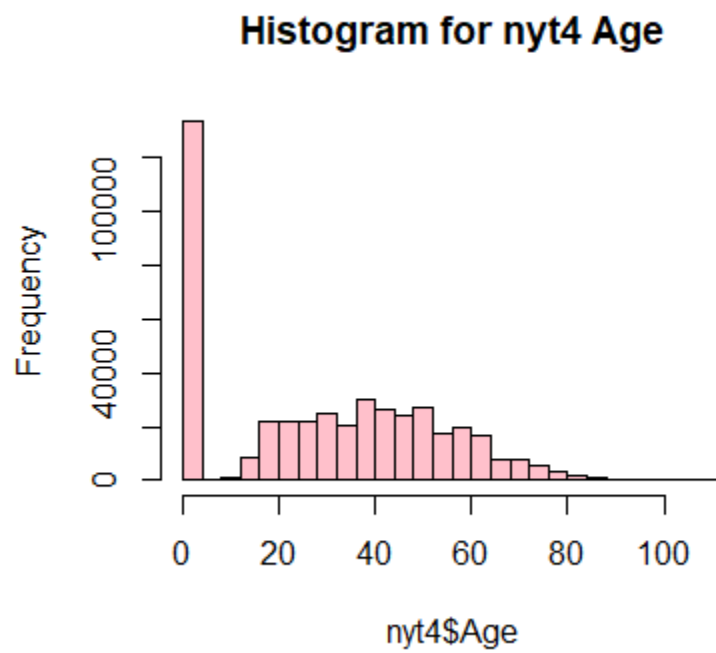
For both Age and Impressions in nyt3, nyt4, nyt5, nyt6, nyt7, nyt8, and nyt9, the data deviates significantly from a normal distribution as a result the null hypothesis (i.e., p-value > 0.05) was rejected as the data did not follow a normal distribution

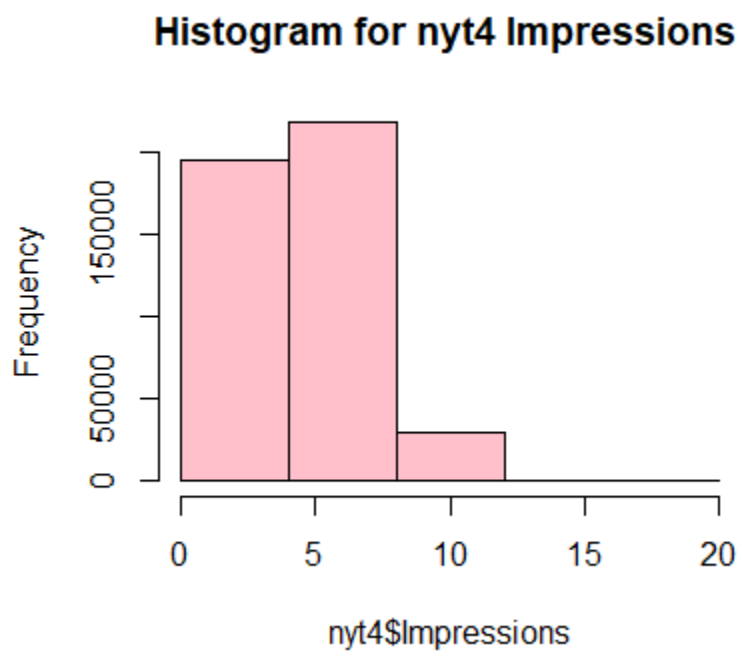
Histogram for nyt3, nyt4, nyt5, nyt6, nyt7, nyt8 and nyt9 "Age and Impressions"

for nyt3,



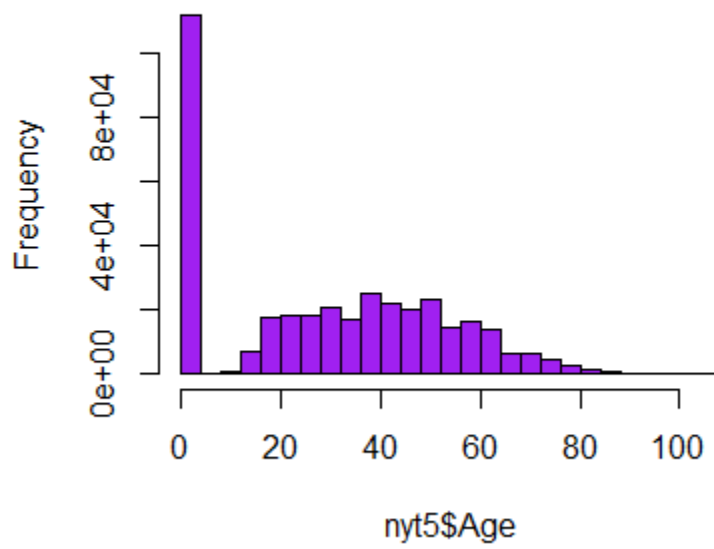
For nyt4,



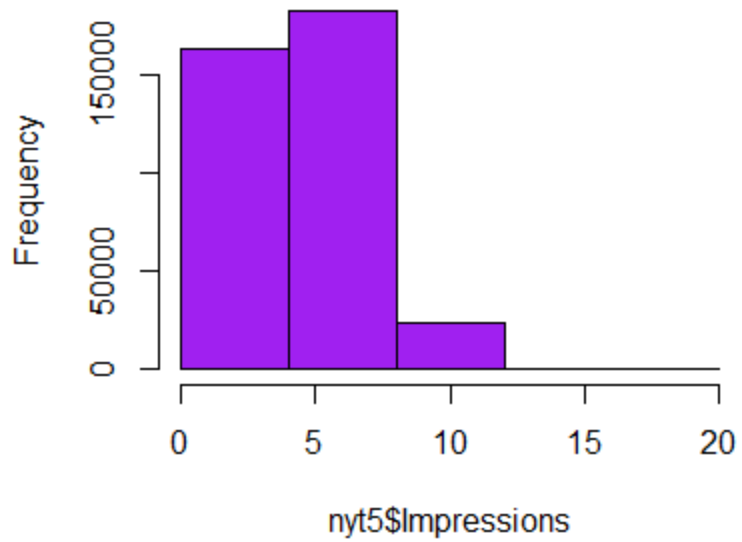


For nyt5,

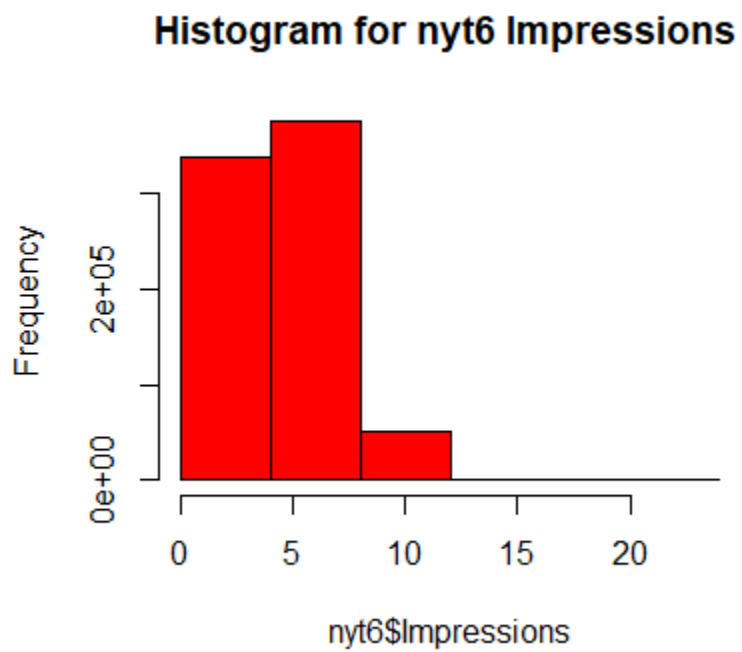
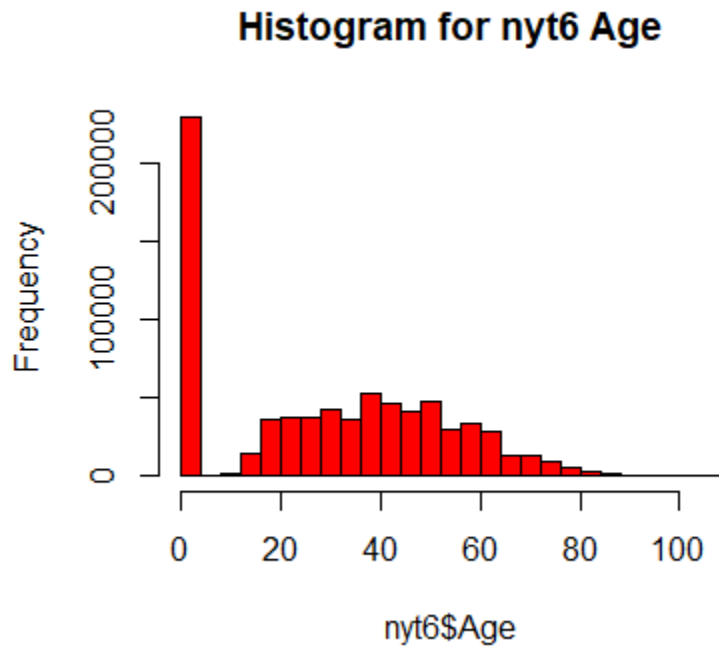
Histogram for nyt5 Age



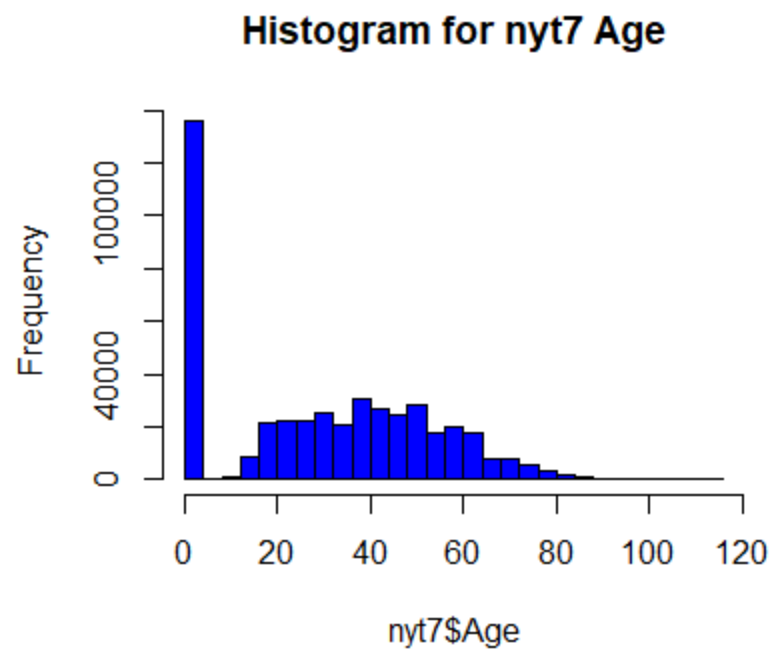
Histogram for nyt5 Impressions



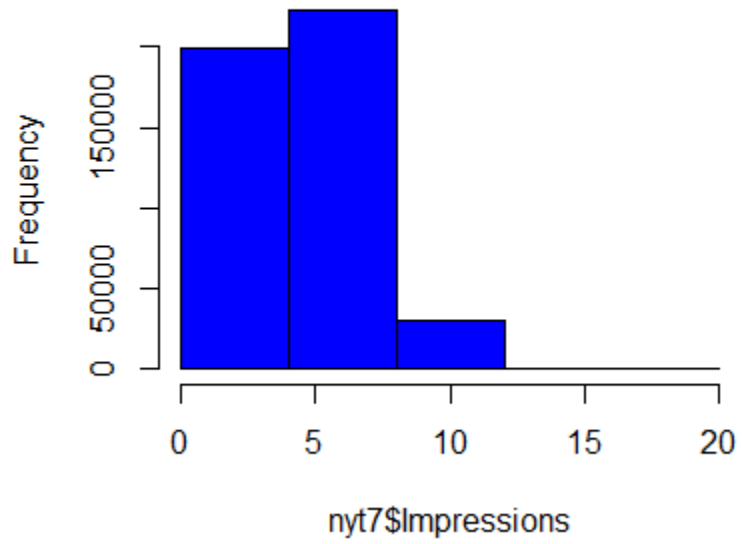
For nyt6,



For nyt7,

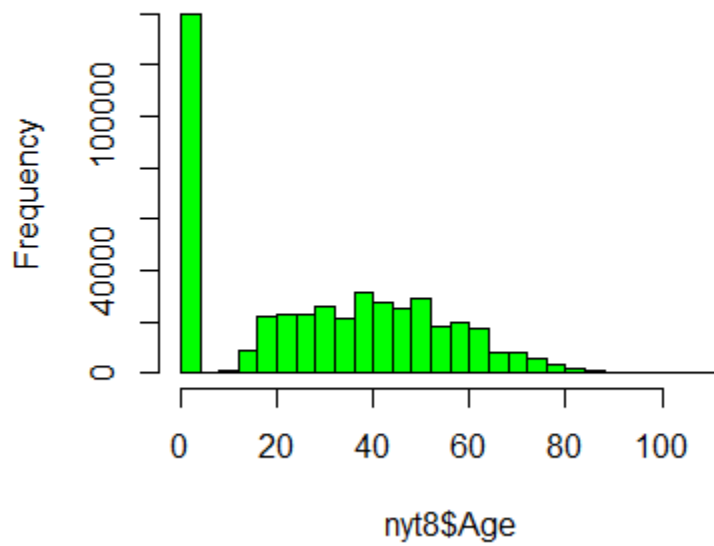


Histogram for nyt7 Impressions

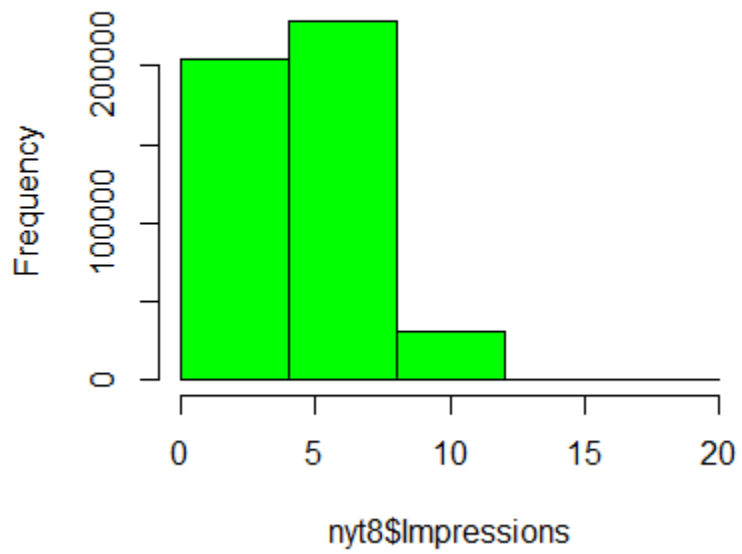


For nyt8,

Histogram for nyt8 Age

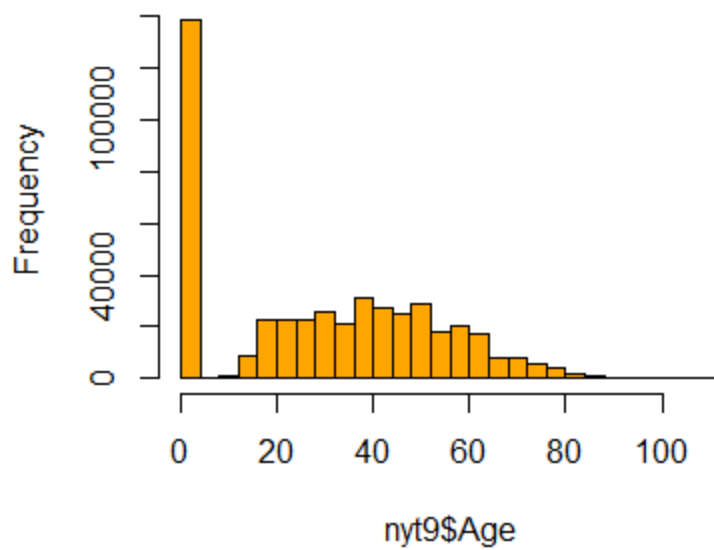


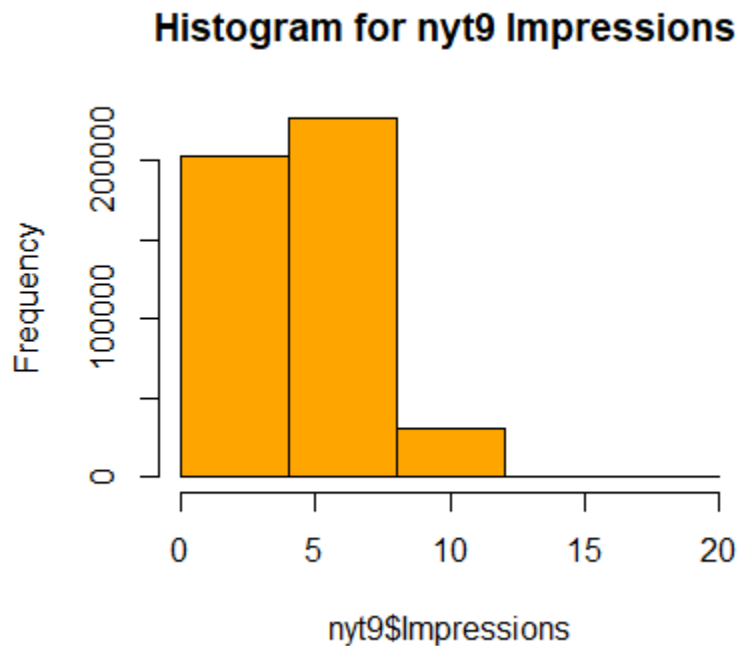
Histogram for nyt8 Impressions



For nyt9

Histogram for nyt9 Age





Age Distribution (for all datasets):

- The "Age" variable for all datasets appears to be positively skewed, as the mean is less than the median, and the distribution extends to higher ages.
- The minimum age is 0, and the maximum age is around 109, with a median age around 31.
- The distribution is right-skewed, indicating that the majority of observations are concentrated towards the younger age group.
- The boxplots show similar trends across datasets, with some variations in the range and quartiles.

Impressions Distribution (for all datasets):

- The "Impressions" variable also exhibits positive skewness in all datasets. The mean is very close to the median, indicating a relatively symmetric distribution.
- Impressions have a minimum value of 0, with the maximum value around 19, and a median value of 5.
- The distribution of Impressions shows fewer extreme values and is more concentrated around the median compared to the Age distribution.

Differences Among Datasets:

- In terms of Age and Impressions, the seven datasets (nyt3, nyt4, nyt5, nyt6, nyt7, nyt8, nyt9) exhibit similar patterns and have nearly identical summary statistics.

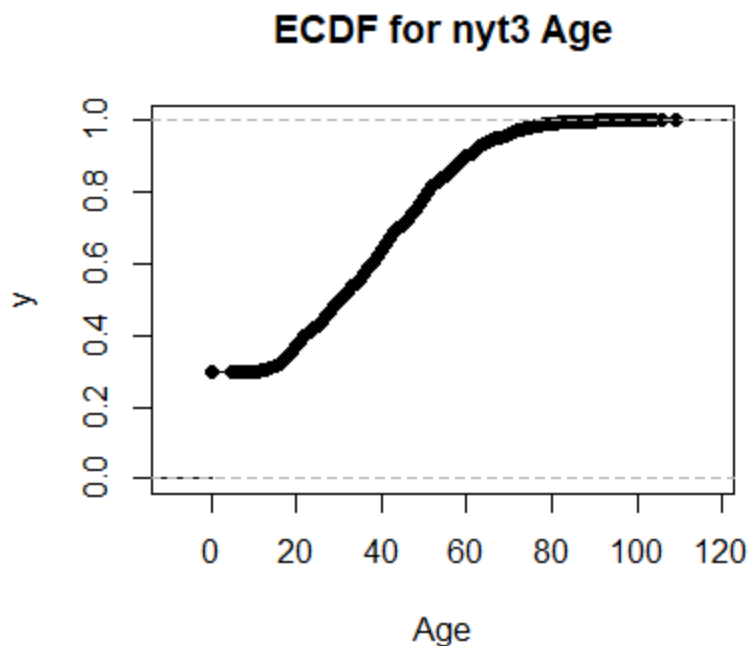
- The distributions are all right-skewed for Age and have a more symmetric distribution for Impressions.
- The distributions are consistent across datasets, with similar central tendencies and ranges.
- The boxplots show consistent characteristics across datasets for both Age and Impressions.

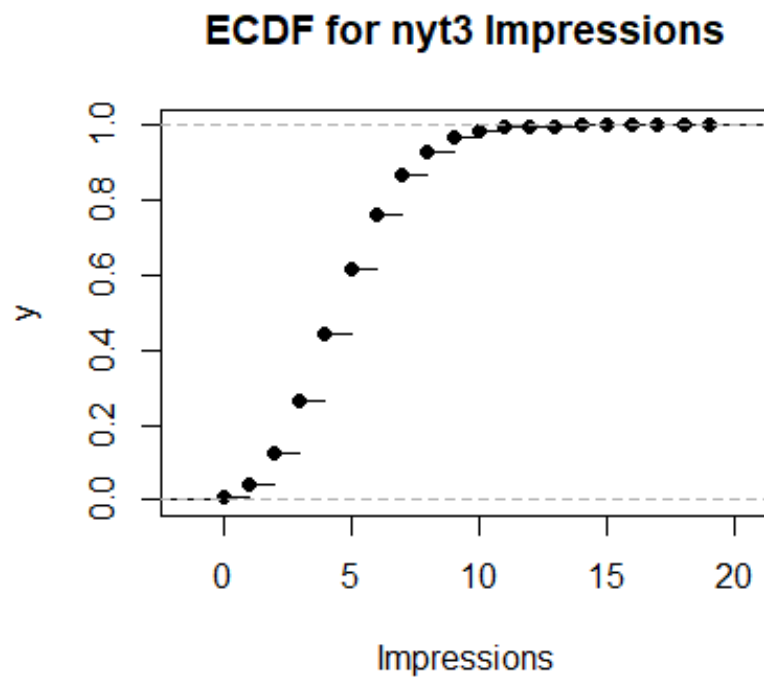
Conclusively, it can be seen that the Age and Impressions variables in these datasets exhibit similar statistical properties and right-skewed, indicating that the distributions are largely consistent among these datasets.

Question c

Empirical Cumulative Distribution Function (ECDFs) for nyt3, nyt4, nyt5, nyt6, nyt7, nyt8 and nyt9 Age and Impressions

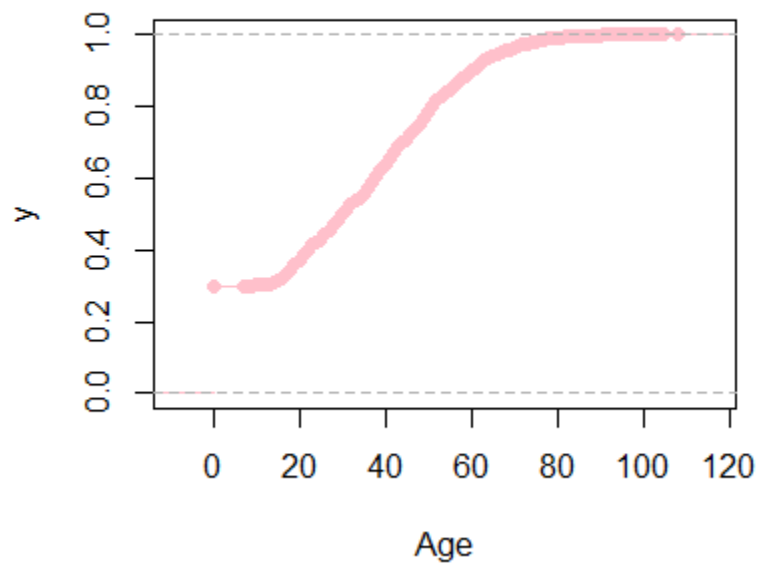
for nyt3



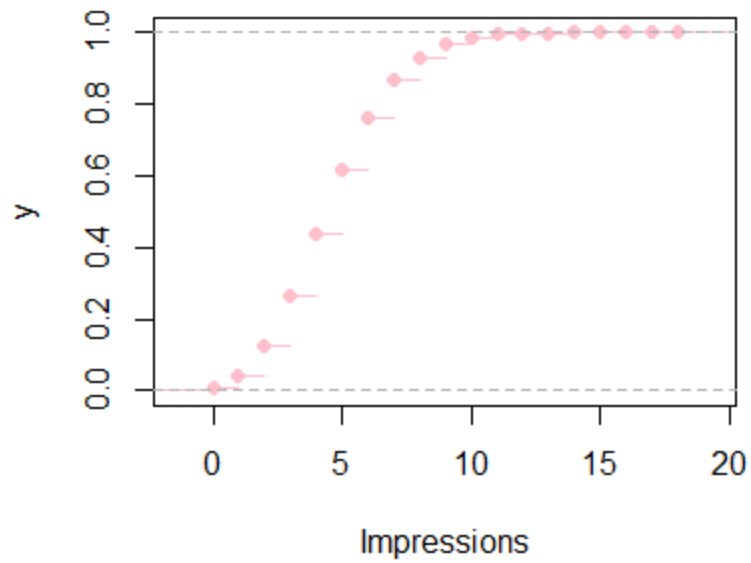


for nyt4

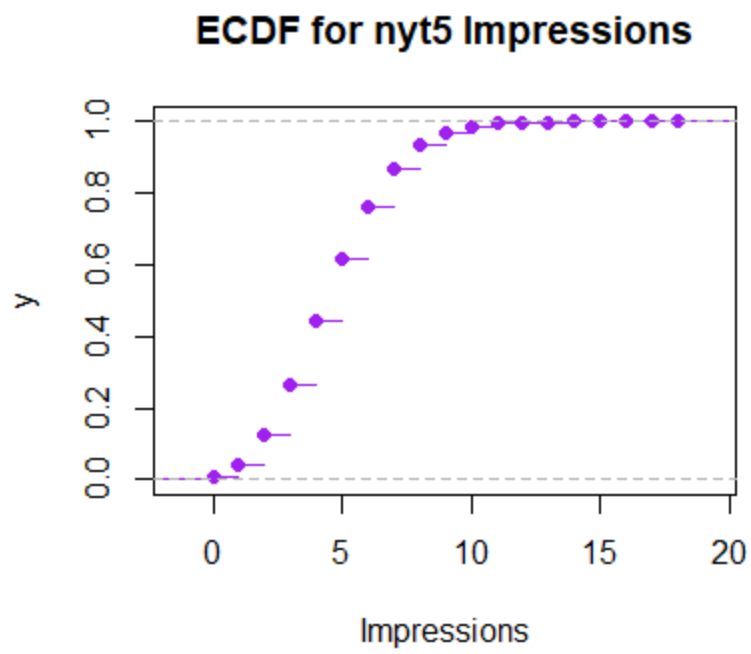
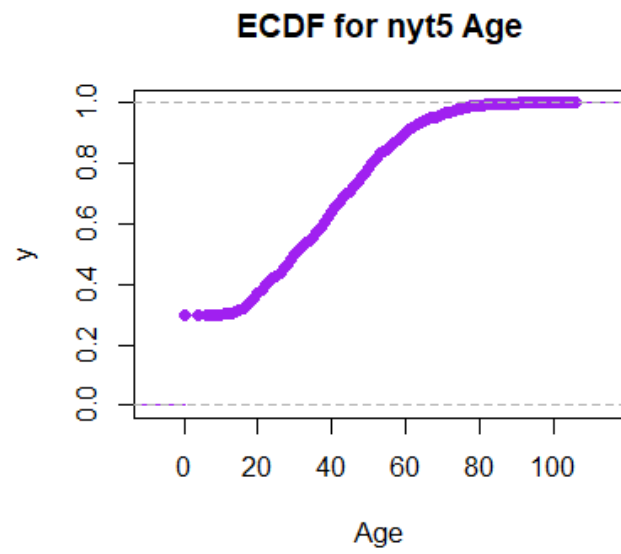
ECDF for nyt4 Age



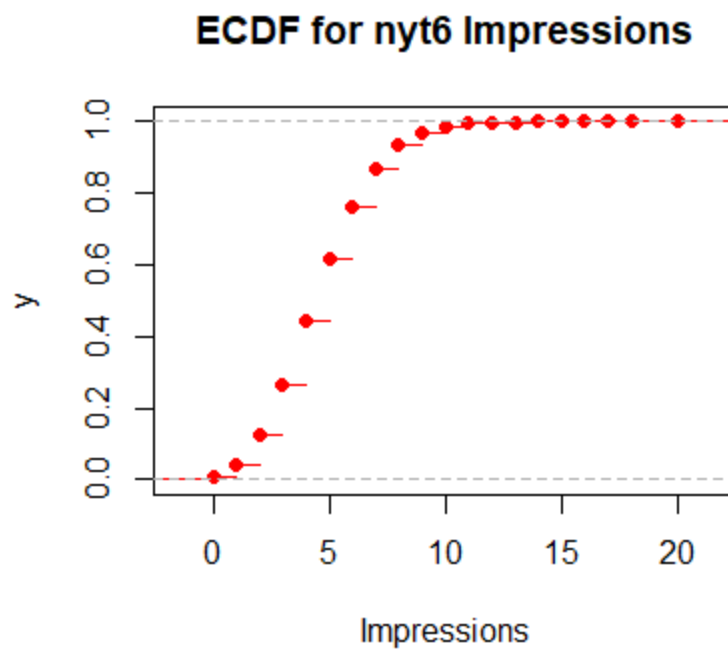
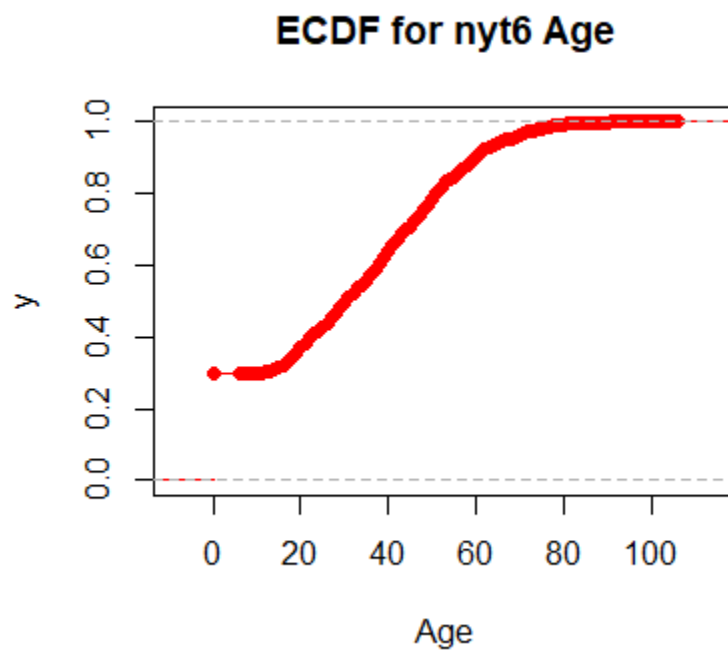
ECDF for nyt4 Impressions



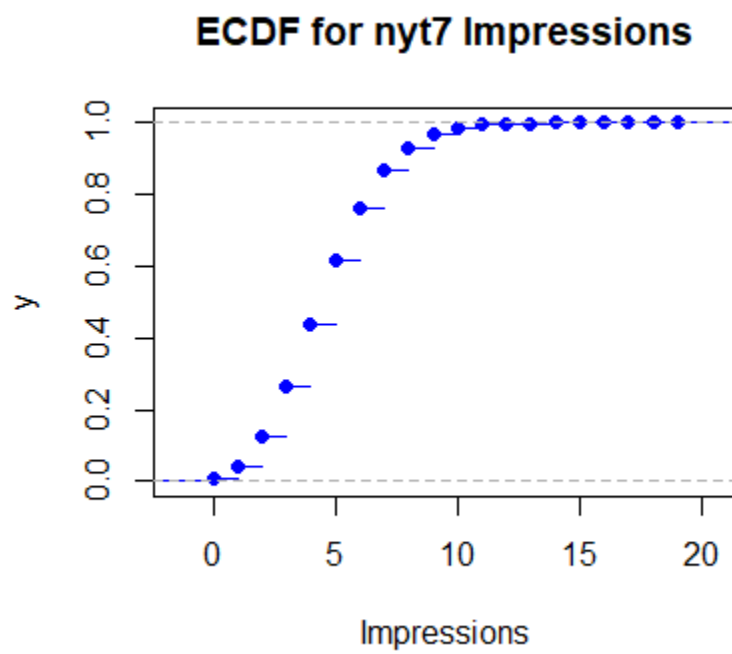
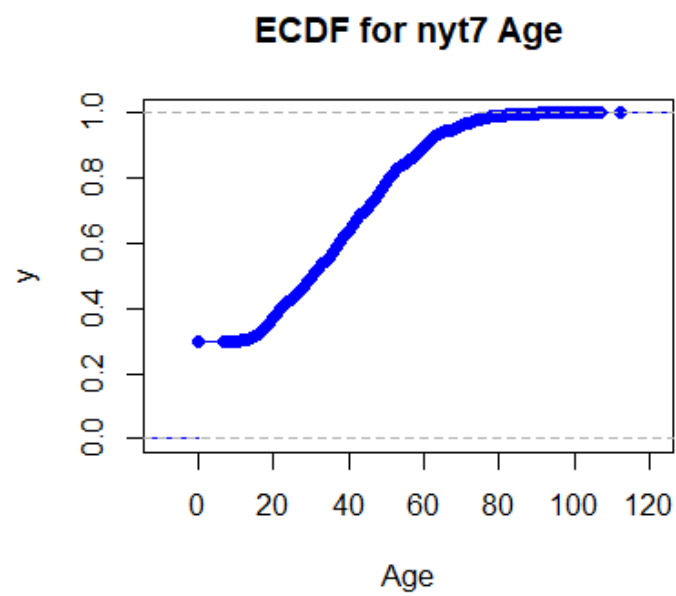
for nyt5



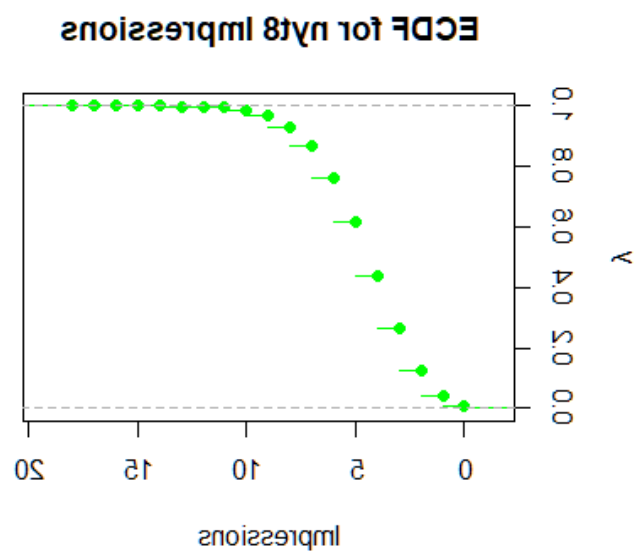
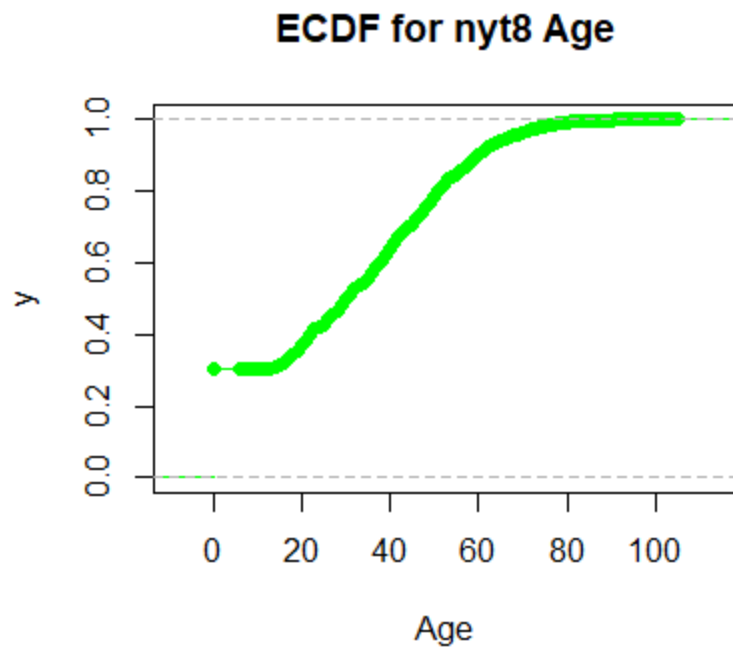
for nyt6



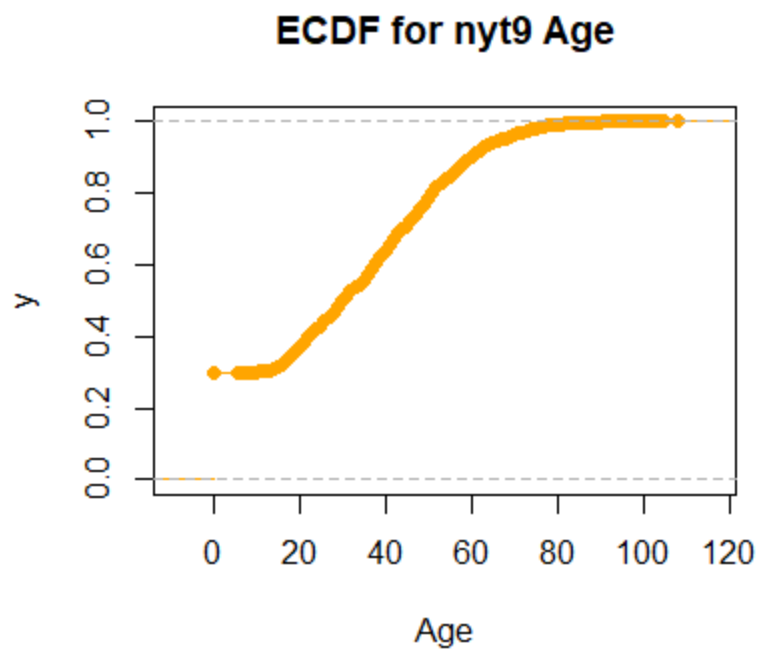
for nyt7

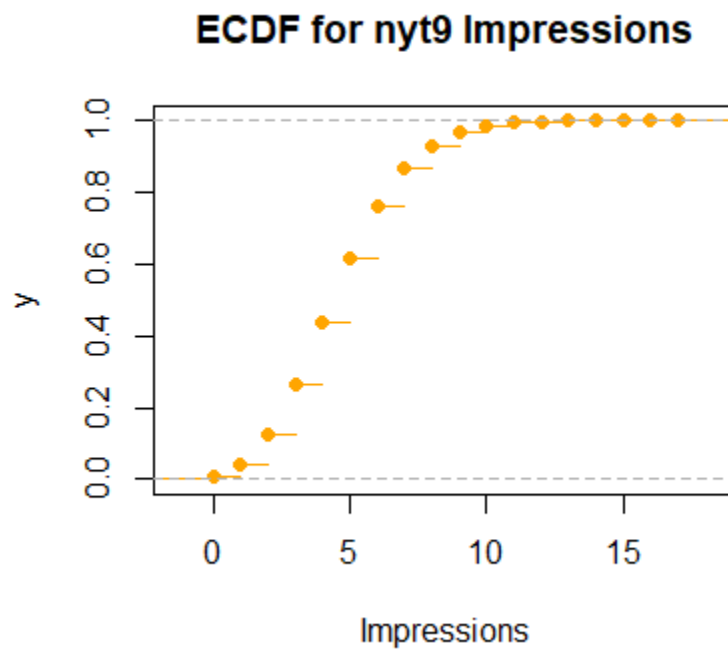


for nyt8



for nyt9

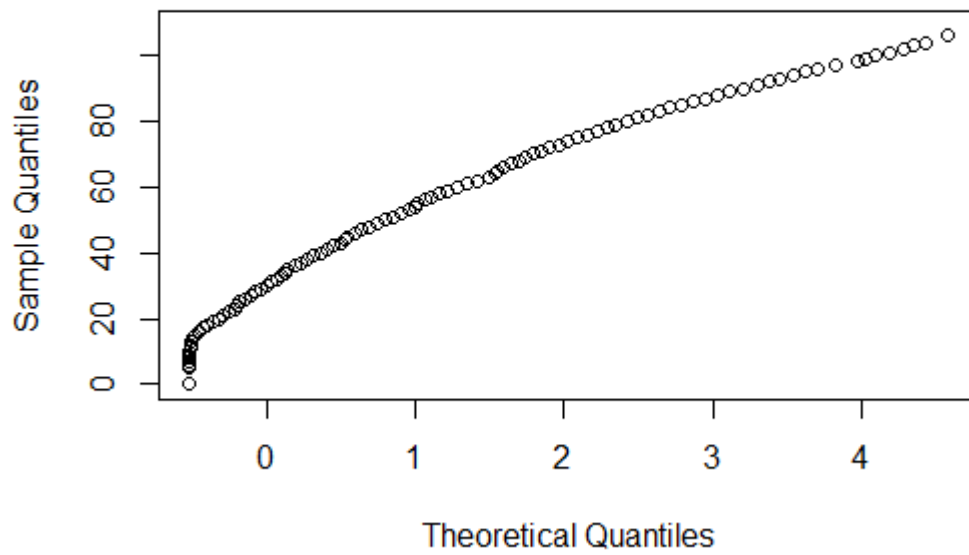




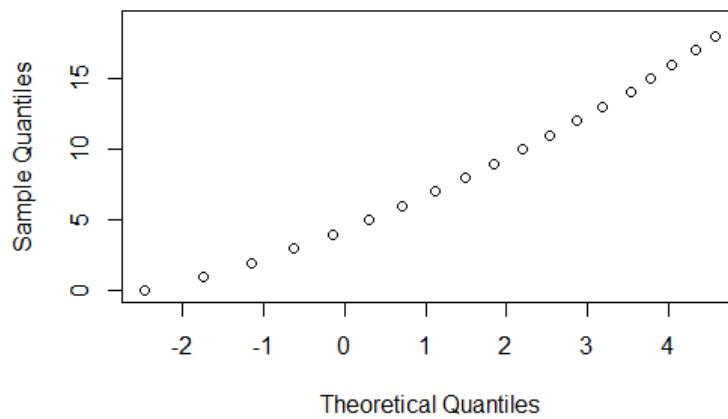
Q-Q plot for nyt3, nyt4, nyt5, nyt6, nyt7, nyt8 and nyt9 Age and Impressions

For nyt3

Q-Q Plot for nyt3 Age

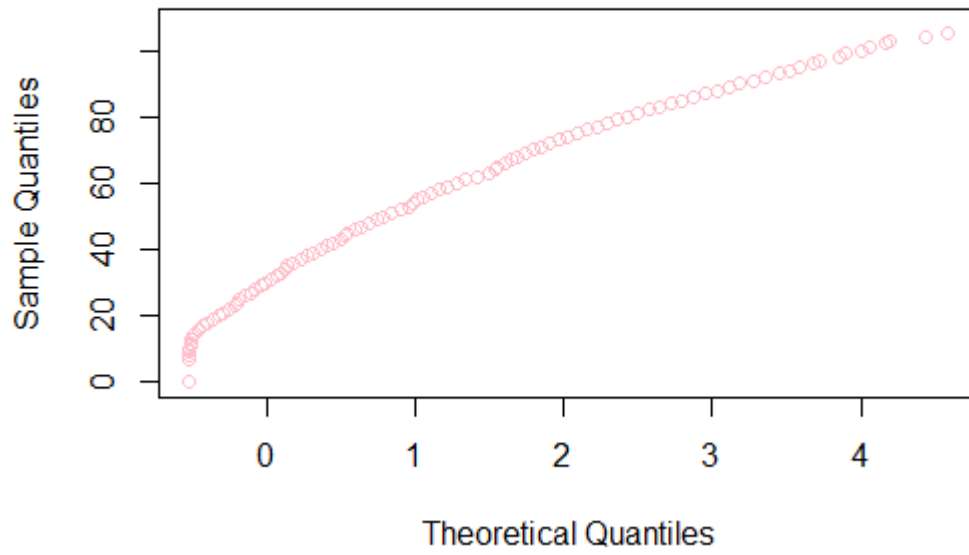


Q-Q Plot for nyt3 Impressions

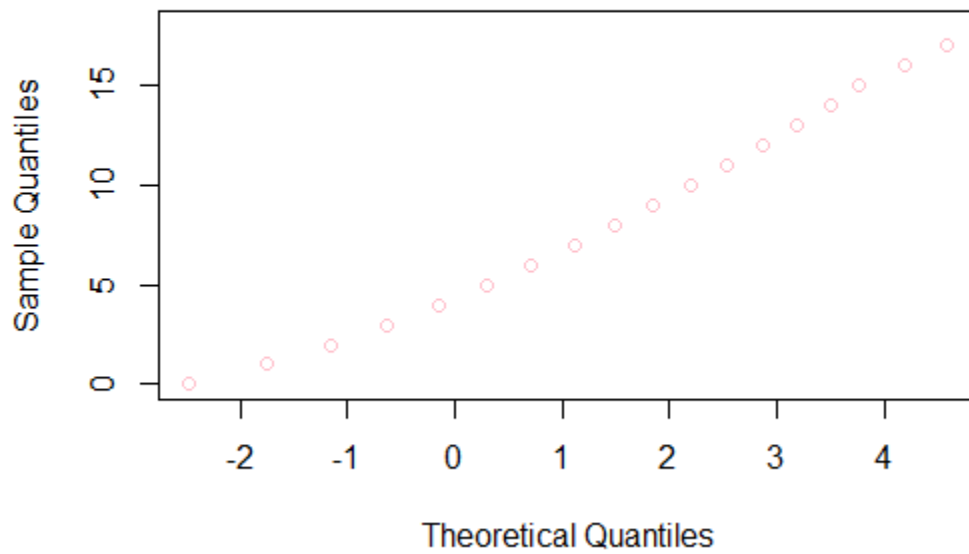


For nyt4

Q-Q Plot for nyt4 Age

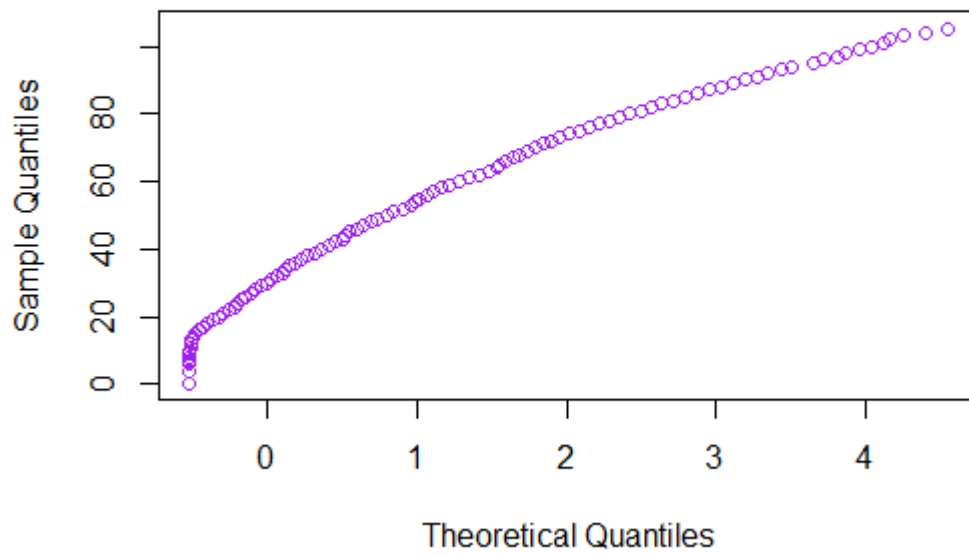


Q-Q Plot for nyt4 Impressions

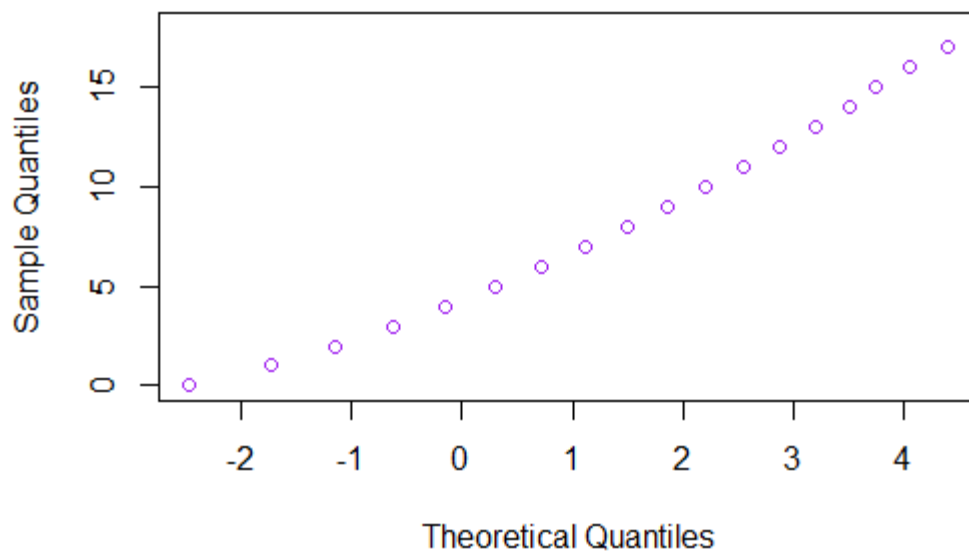


For nyt5

Q-Q Plot for nyt5 Age

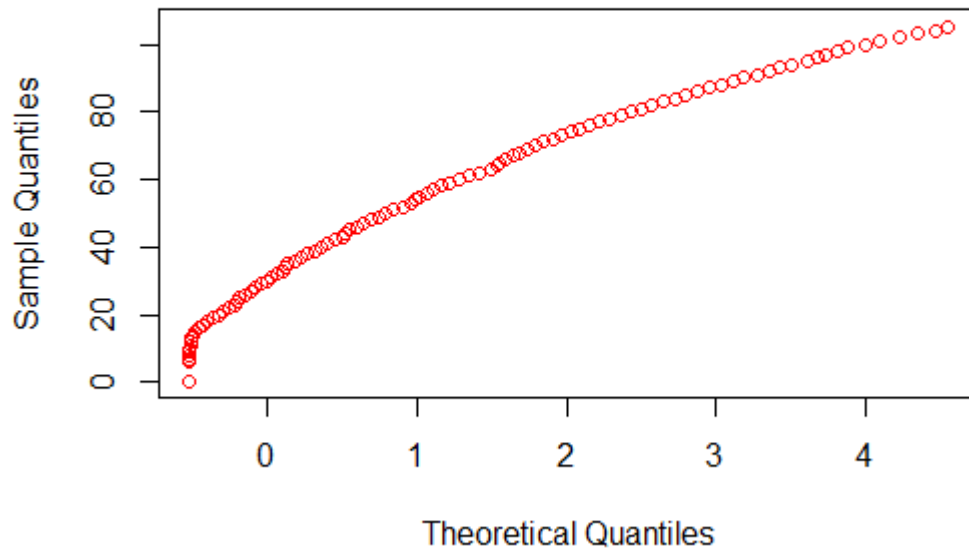


Q-Q Plot for nyt5 Impressions

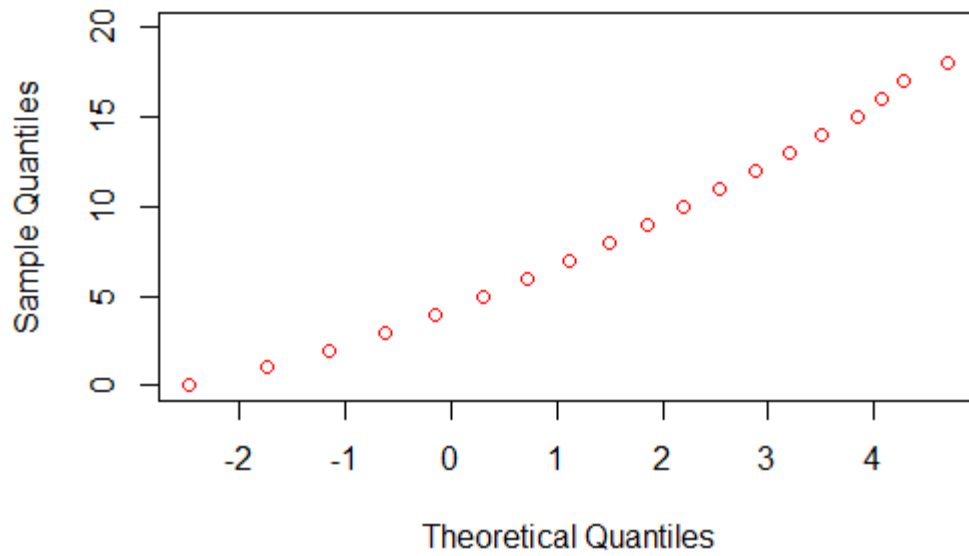


For nyt6

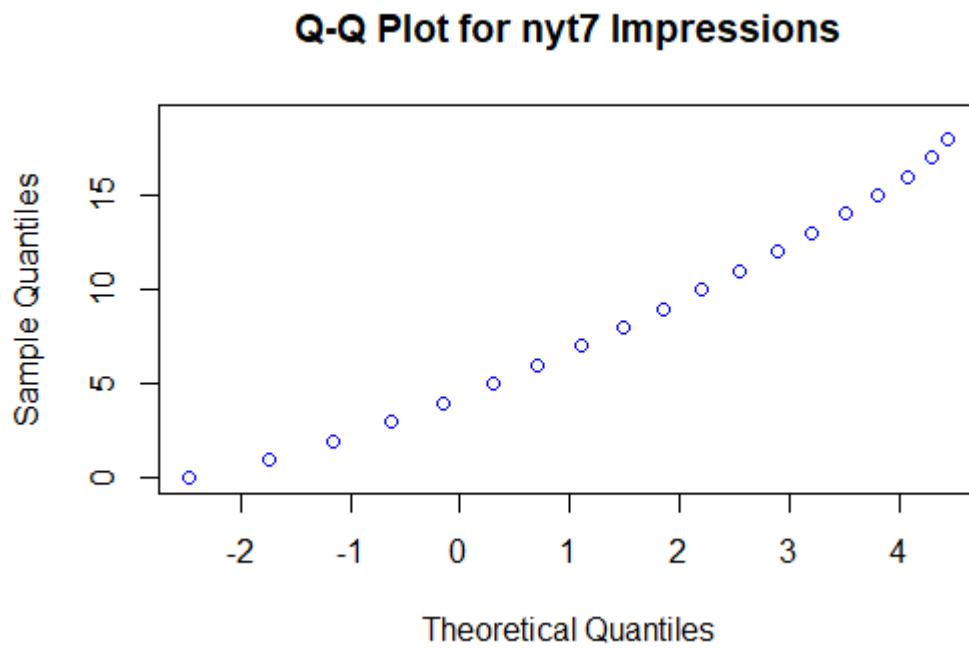
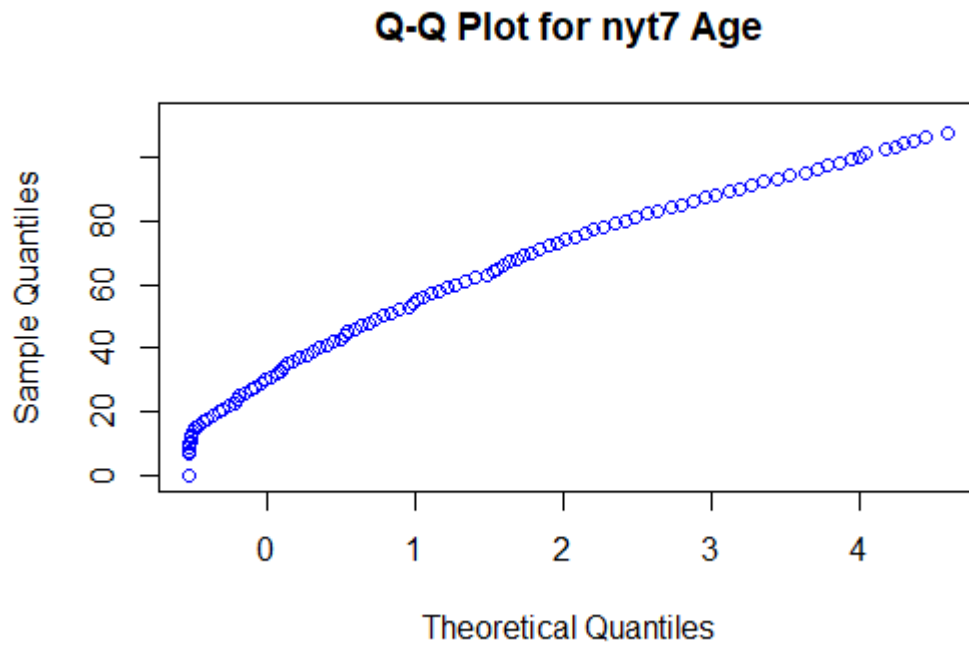
Q-Q Plot for nyt6 Age



Q-Q Plot for nyt6 Impressions

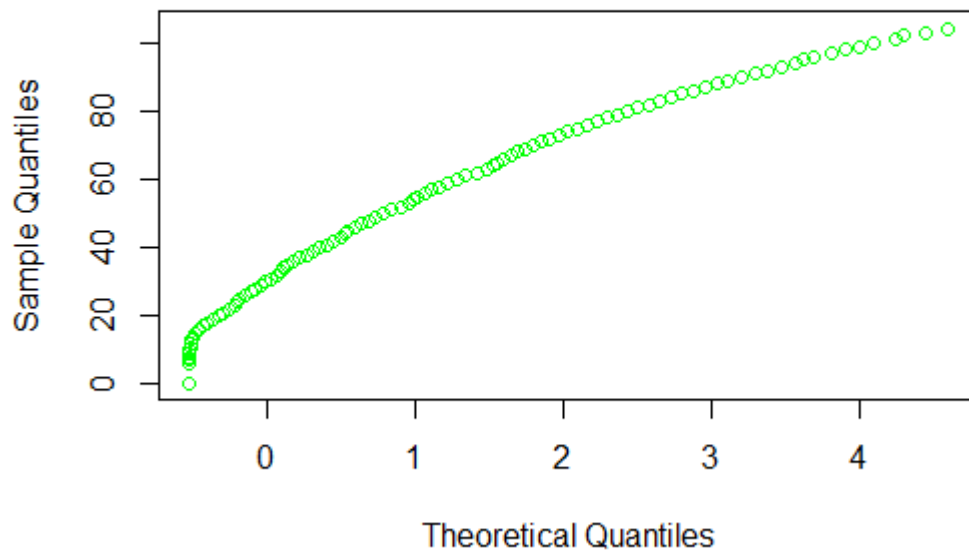


For nyt7

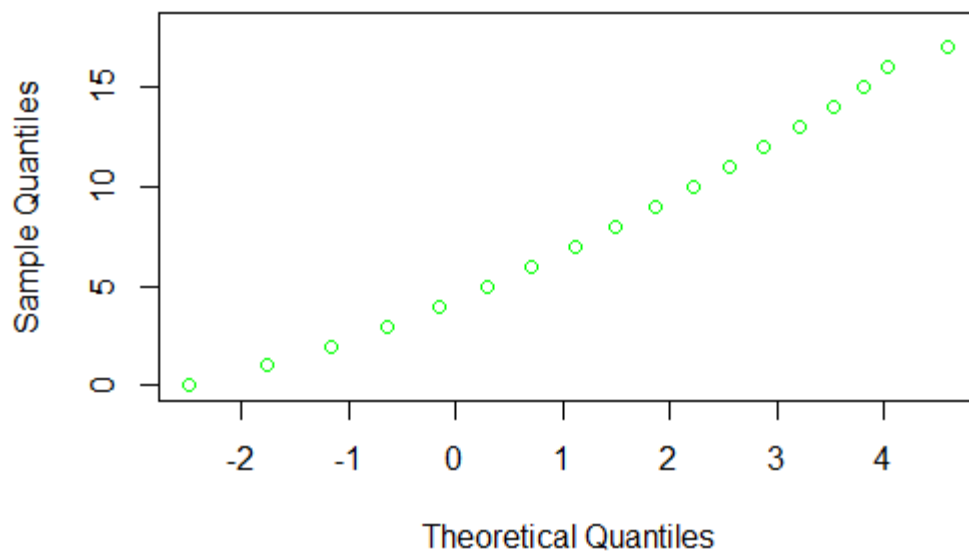


For nyt8

Q-Q Plot for nyt8 Age

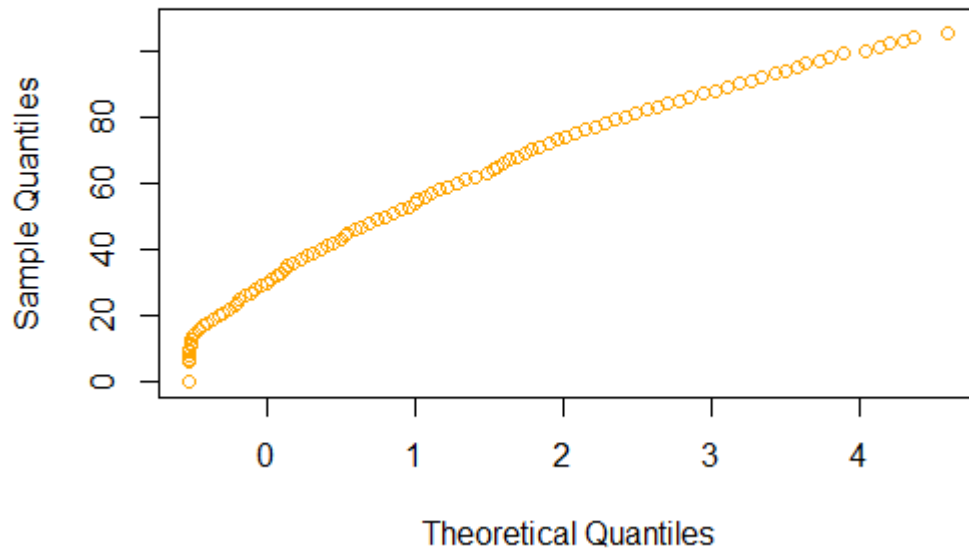


Q-Q Plot for nyt8 Impressions

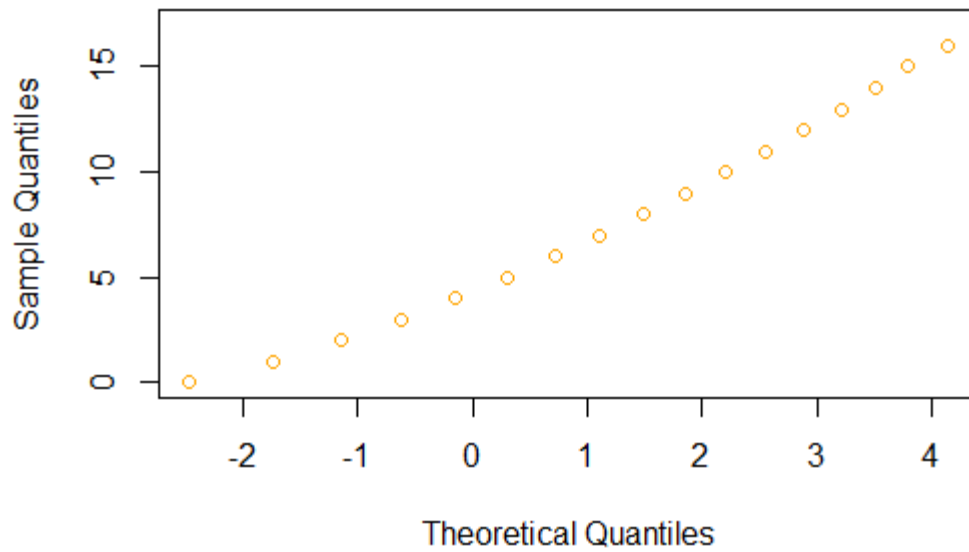


For nyt9

Q-Q Plot for nyt9 Age



Q-Q Plot for nyt9 Impressions



From the ECDFs of Impressions all the values are centered around 5 which is the expected the median. The quantile-quantile plots also follow a near straight line and hence we can conclude that the distribution is normal between 0 - 10. The same can be said about the Age variable, which follows an almost normal distribution between 20 to 60.

Question d

Significance Testing for nyt3, nyt4, nyt5, nyt6, nyt7, nyt8 and nyt9 Age and Impressions

I used T- test for my Significance Testing. T-test, is a statistical hypothesis test used to determine whether there is a significant difference between the means of two groups or conditions.

T-Test Results for Age and Impressions Data:

nyt3:

- t-value: 683.93
- Degrees of Freedom: 448,241
- p-value: $< 2.2e-16$ (essentially zero)

A highly significant difference between the means of "nyt3\$Age" and "nyt3\$Impressions" was observed. The 95% confidence interval for the difference in means (24.40155, 24.54181) does not contain zero, indicating a substantial difference.

nyt4:

- t-value: 684.26
- Degrees of Freedom: 450,778
- p-value: $< 2.2e-16$ (essentially zero)

A highly significant difference between the means of "nyt4\$Age" and "nyt4\$Impressions" was observed. The 95% confidence interval for the difference in means (24.35846, 24.49841) does not contain zero, indicating a substantial difference.

nyt5:

- t-value: 626.03
- Degrees of Freedom: 376,957
- p-value: $< 2.2e-16$ (essentially zero)

A highly significant difference between the means of "nyt5\$Age" and "nyt5\$Impressions" was observed. The 95% confidence interval for the difference in means (24.35811, 24.51111) does not contain zero, indicating a substantial difference.

nyt6:

- t-value: 901.74

- Degrees of Freedom: 778,197
- p-value: $< 2.2e-16$ (essentially zero)

A very substantial difference between the means of "nyt6\$Age" and "nyt6\$Impressions" was observed. The 95% confidence interval for the difference in means (24.40973, 24.51607) does not contain zero, indicating a substantial difference.

nyt7:

- t-value: 693.47
- Degrees of Freedom: 460,565
- p-value: $< 2.2e-16$ (essentially zero)

A substantial difference between the means of "nyt7\$Age" and "nyt7\$Impressions" was observed. The 95% confidence interval for the difference in means (24.4504, 24.5890) does not contain zero, indicating a substantial difference.

nyt8:

- t-value: 699.93
- Degrees of Freedom: 471,444
- p-value: $< 2.2e-16$ (essentially zero)

A substantial difference between the means of "nyt8\$Age" and "nyt8\$Impressions" was observed. The 95% confidence interval for the difference in means (24.34483, 24.48156) does not contain zero, indicating a substantial difference.

nyt9:

- t-value: 697.77
- Degrees of Freedom: 467,674
- p-value: $< 2.2e-16$ (essentially zero)

A substantial difference between the means of "nyt9\$Age" and "nyt9\$Impressions" was observed. The 95% confidence interval for the difference in means (24.38298, 24.52034) does not contain zero, indicating a substantial difference.

Conclusively, in all datasets, the extremely low p-values and non-overlapping confidence intervals provide strong evidence against the null hypothesis.

Question e

observations about the datasets/ variables

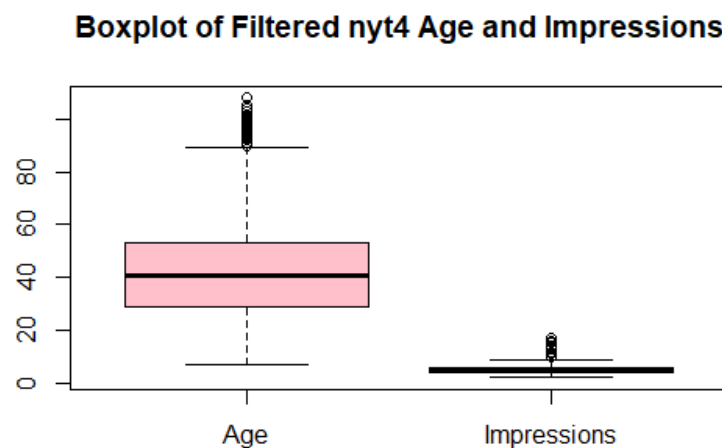
I observed the graphical representations of the data for both boxplots and histograms, when generated, exhibited strikingly similar patterns across all four datasets (nyt3, nyt4, nyt5, and nyt6). Also, the right-skewed pattern was consistent among the age and impressions distributions graphical representations. I noticed that that the datasets contain a significant number of zero values

In addition to this, the results of the Anderson-Darling normality tests; obtained p-values were consistently less than 0.05 for both age and impressions data, signaling a clear departure from normal distribution. This reinforces the understanding that these datasets do not conform to a typical bell-shaped curve. These findings collectively underscore the uniqueness of the datasets and their non-normally distributed nature.

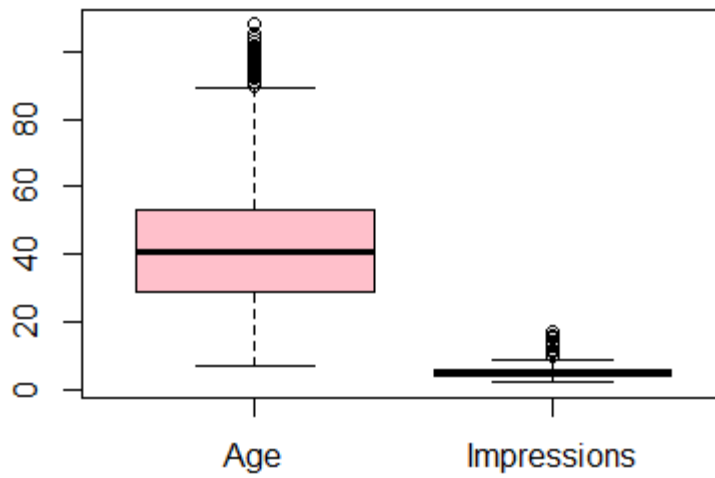
Question 2

Filtered the distributions for Age and Impressions for nyt3, nyt4, nyt5 and nyt6

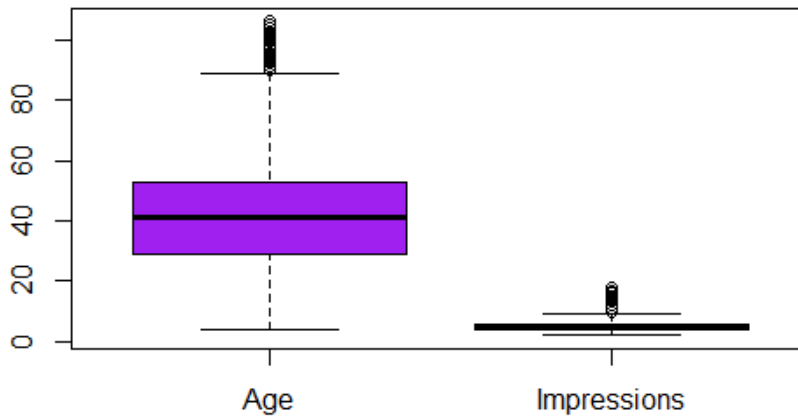
Box plot



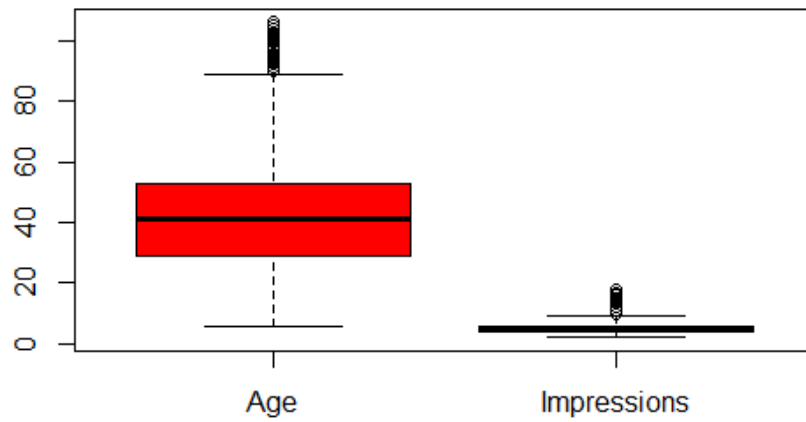
Boxplot of Filtered nyt4 Age and Impressions:



Boxplot of Filtered nyt5 Age and Impressions

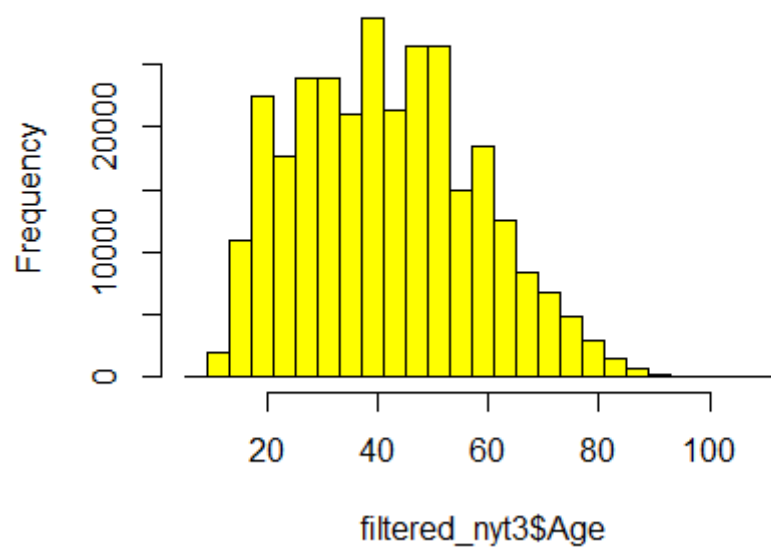


Boxplot of Filtered nyt6 Age and Impressions

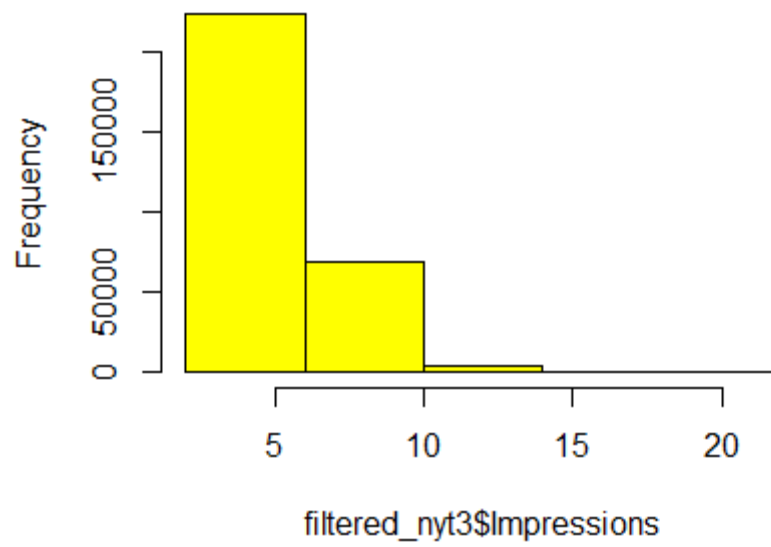


Histogram

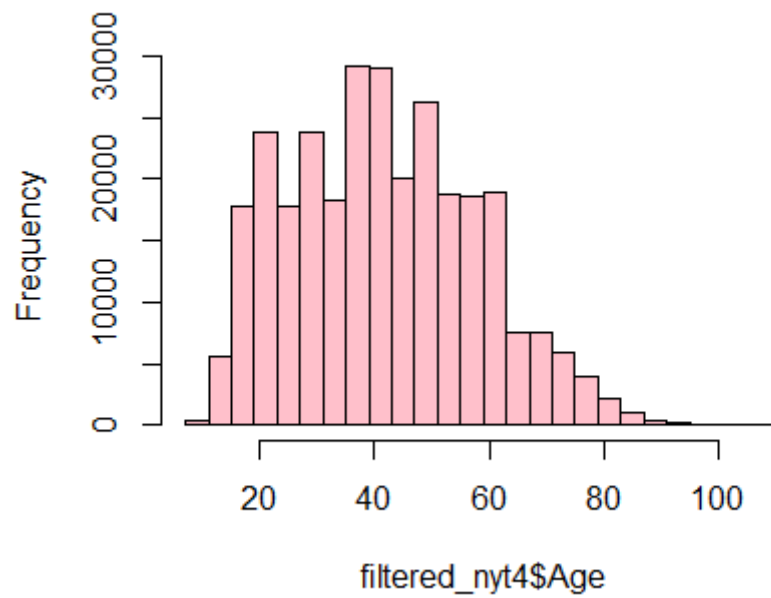
Histogram for filtered_nyt3 Age



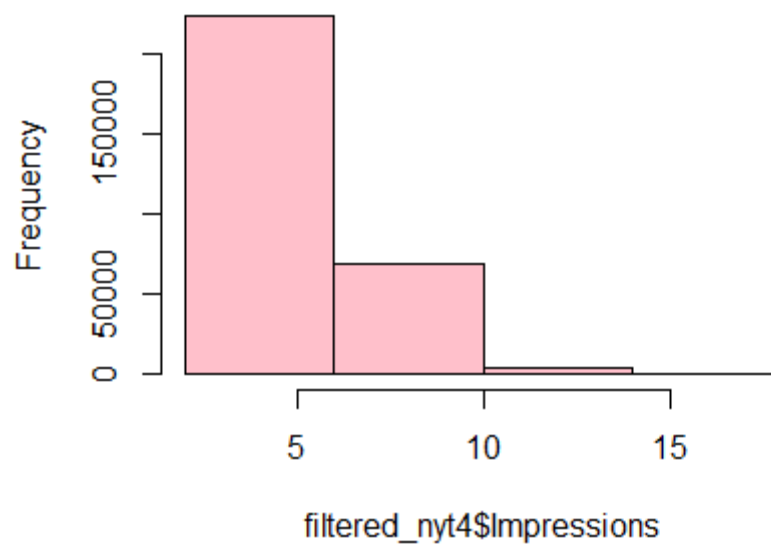
Histogram for filtered_nyt3 Impressions



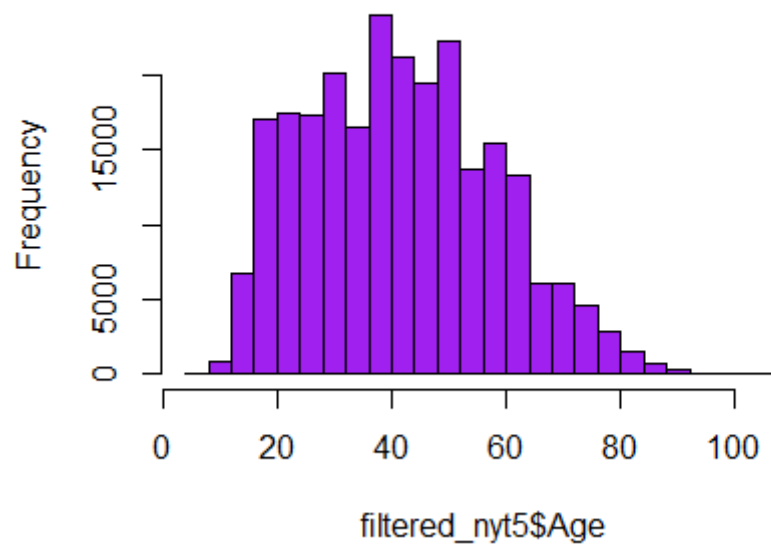
Histogram for filtered_nyt4 Age



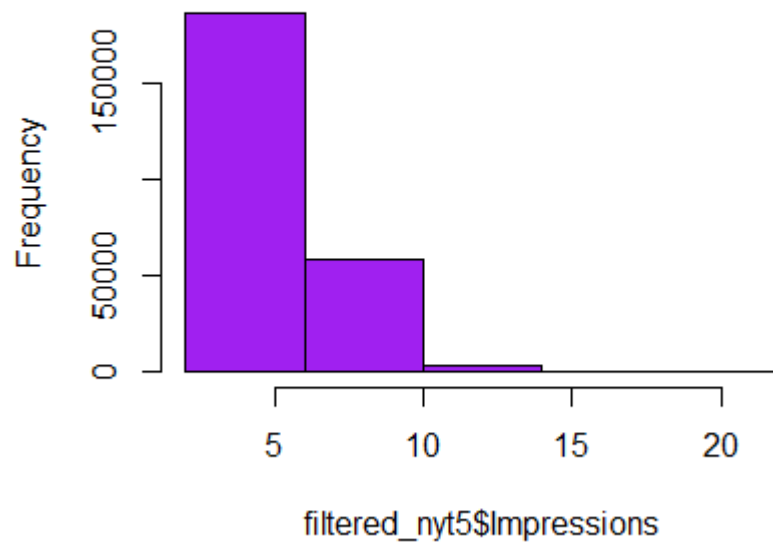
Histogram for filtered_nyt4 Impressions



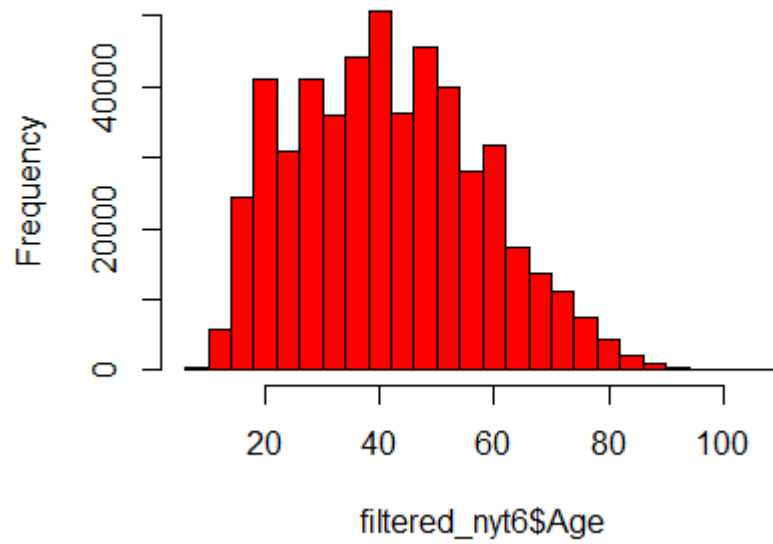
Histogram for filtered_nyt5 Age

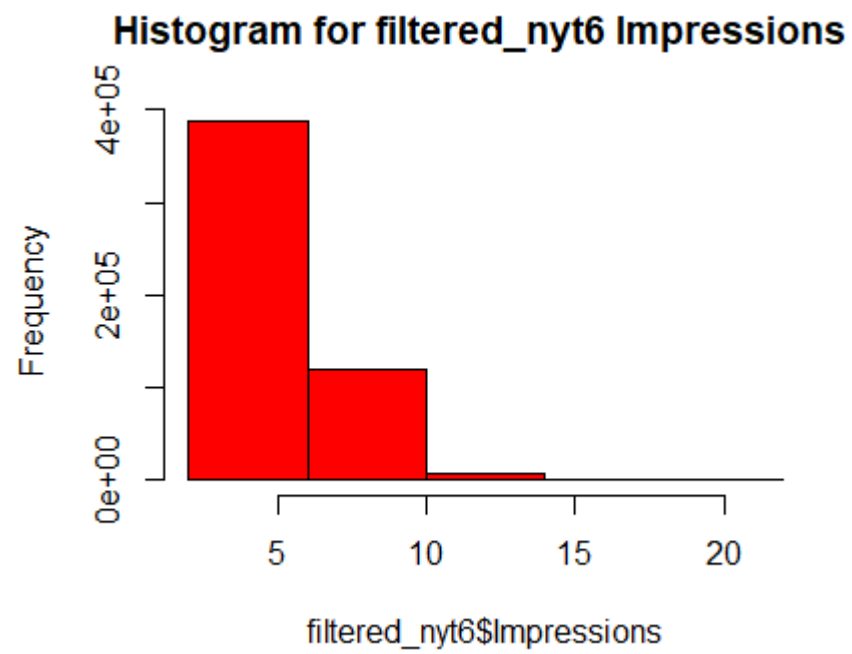


Histogram for nyt5 Impressions

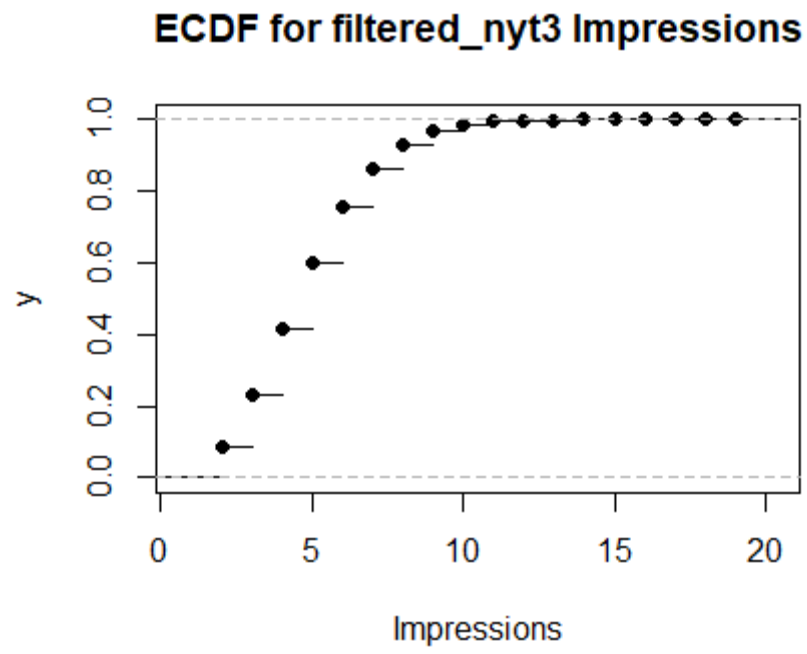


Histogram for filtered_nyt6 Age



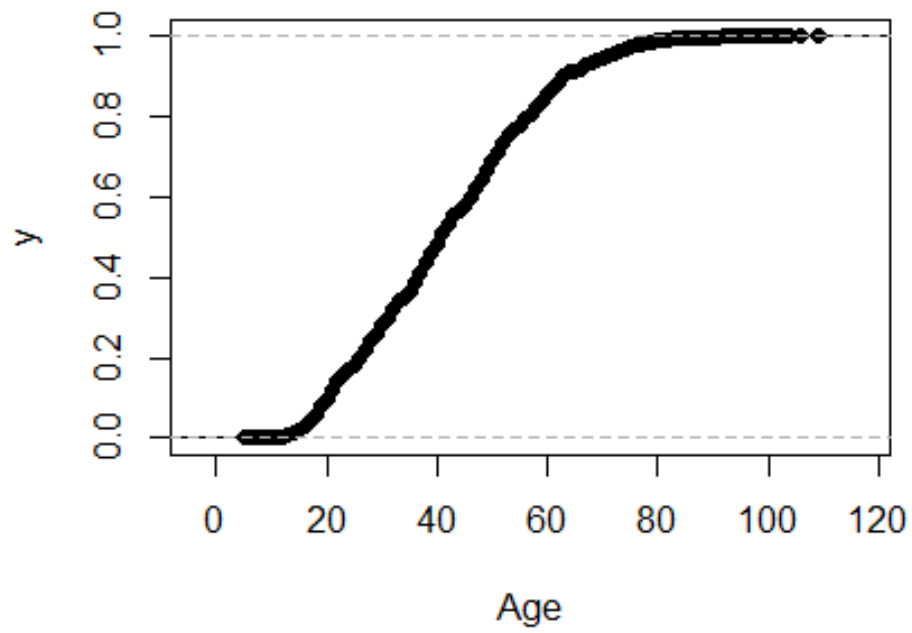


ECDF

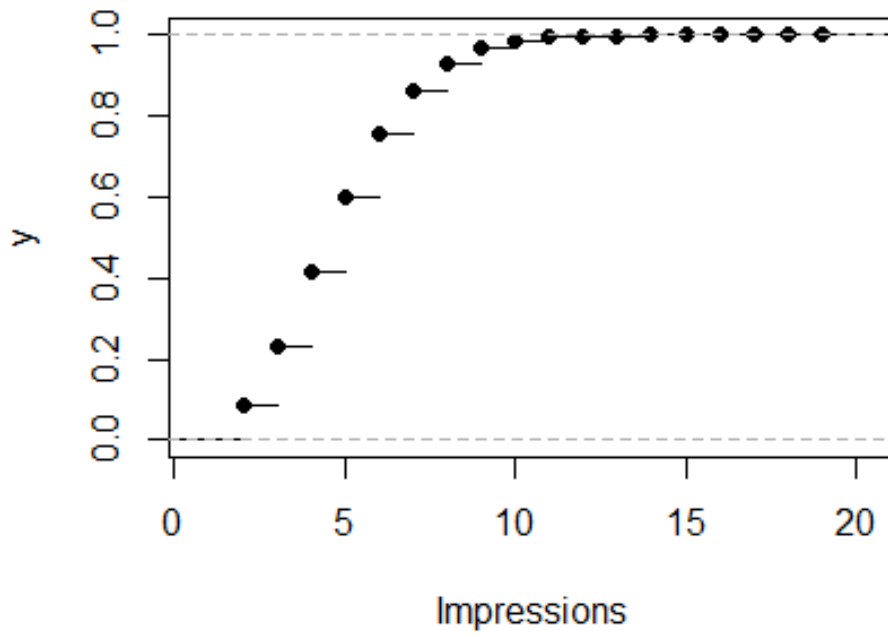


ECDF

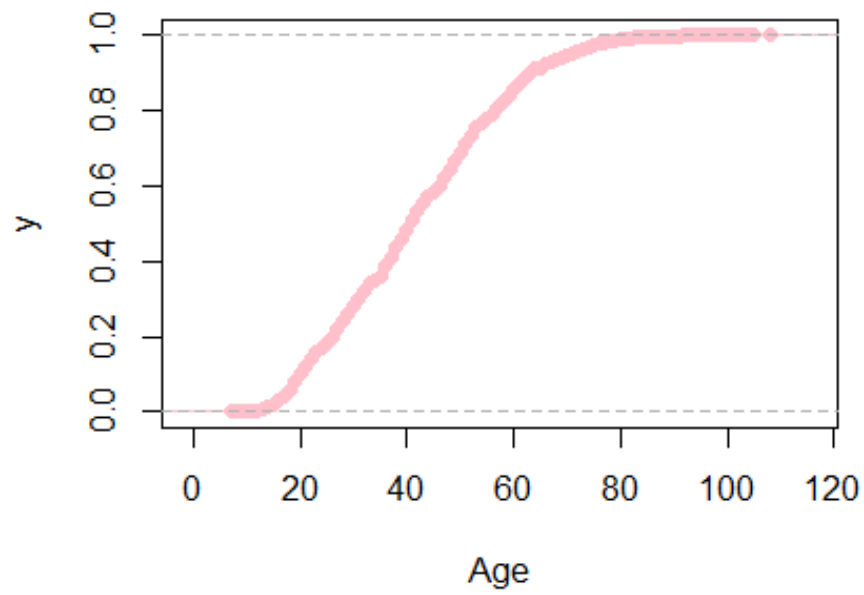
ECDF for filtered_nyt3 Age



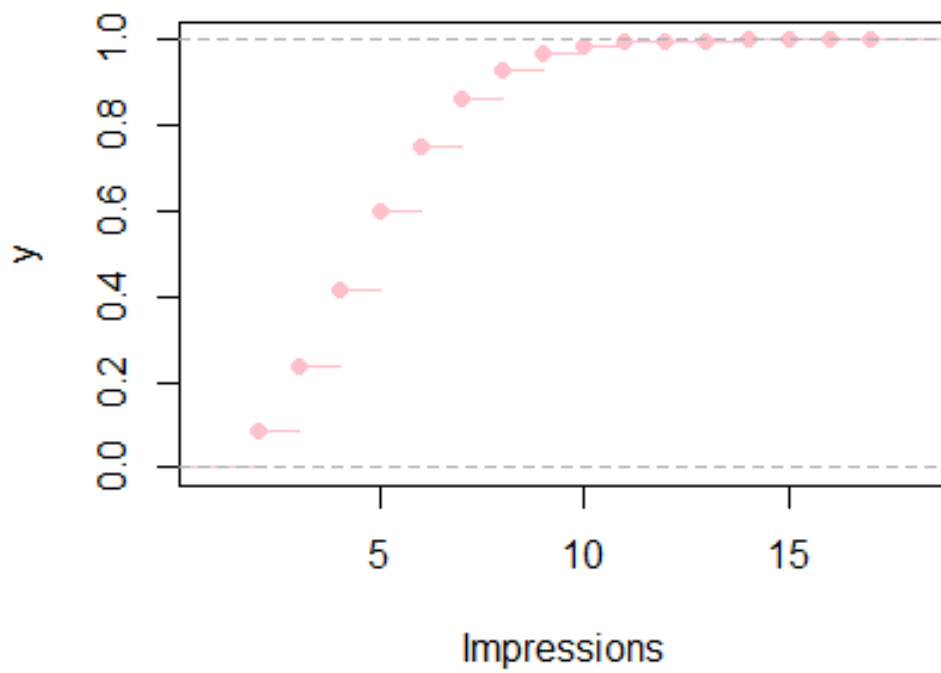
ECDF for filtered_nyt3 Impressions



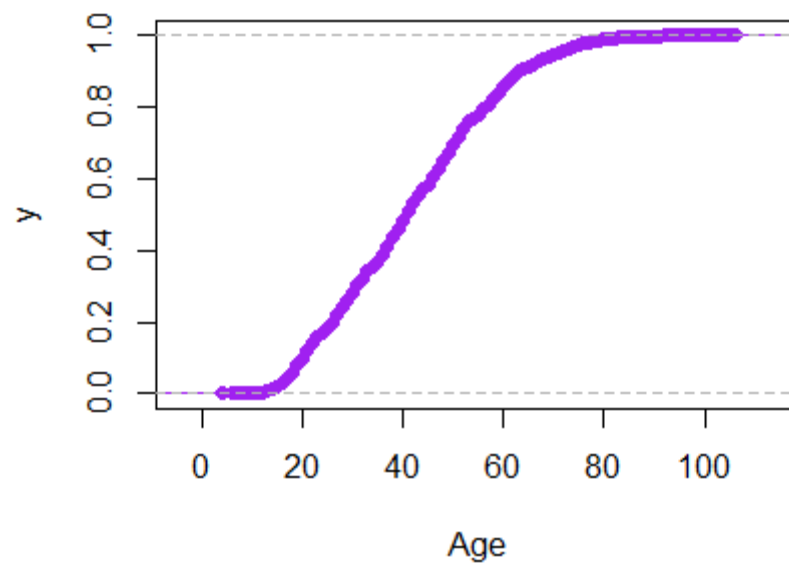
ECDF for filtered_nyt4 Age



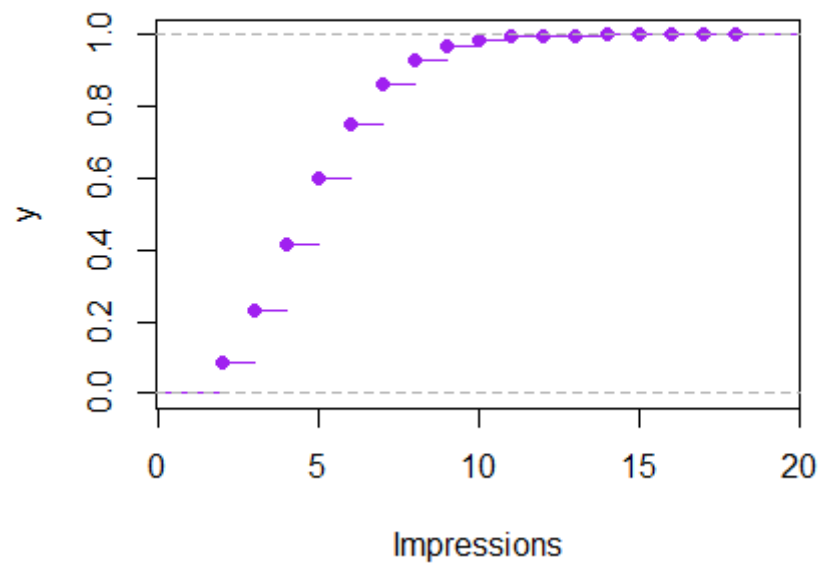
ECDF for filtered_nyt4 Impressions



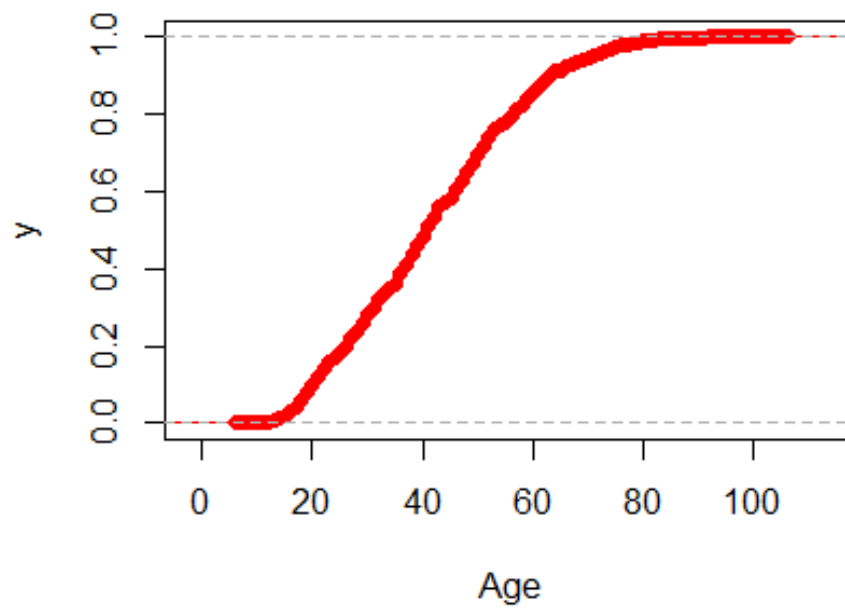
ECDF for filtered_nyt5 Age



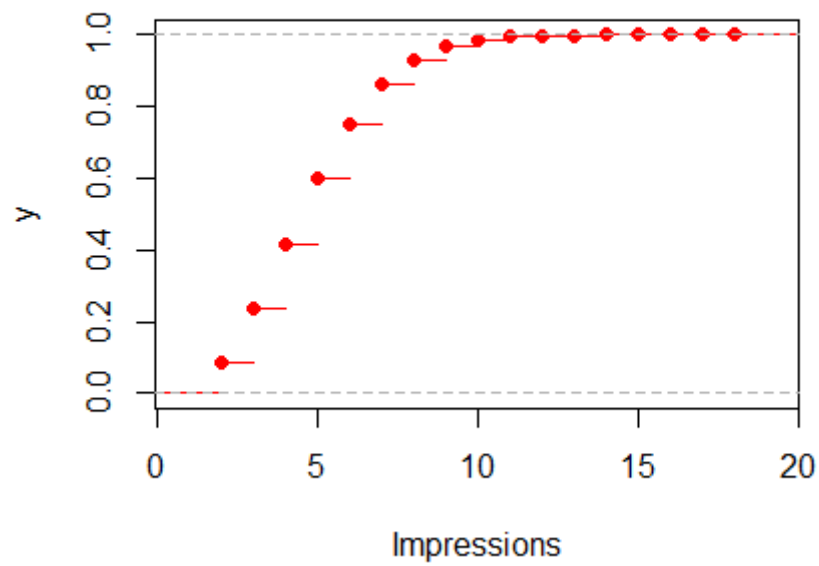
ECDF for filtered_nyt5 Impressions



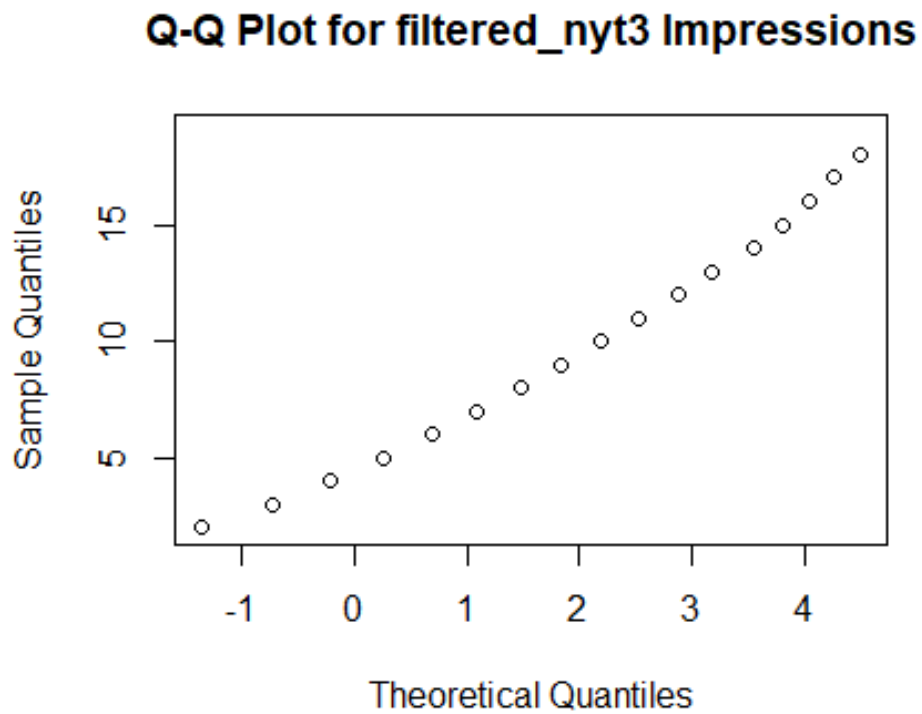
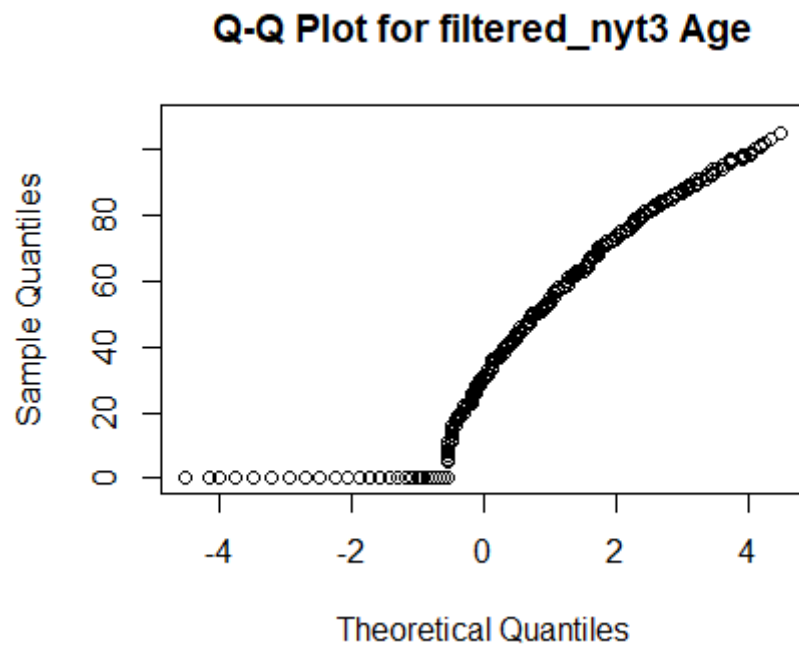
ECDF for filtered_nyt6 Age



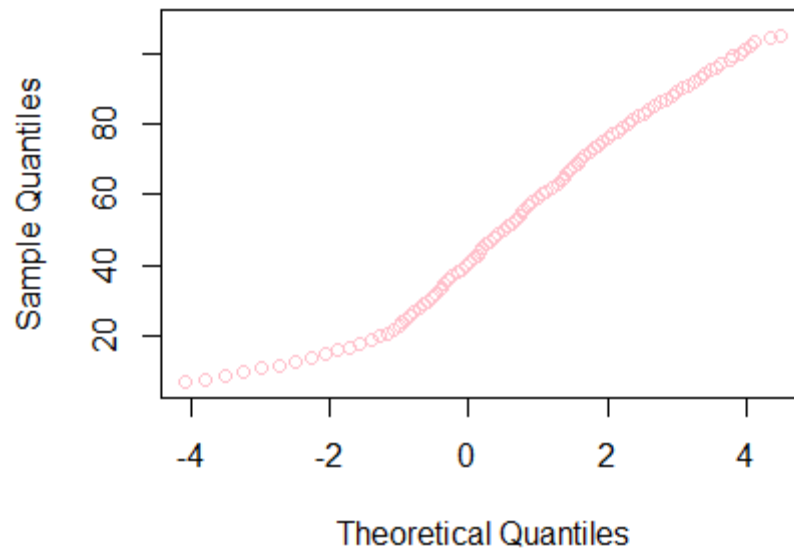
ECDF for filtered_nyt6 Impressions



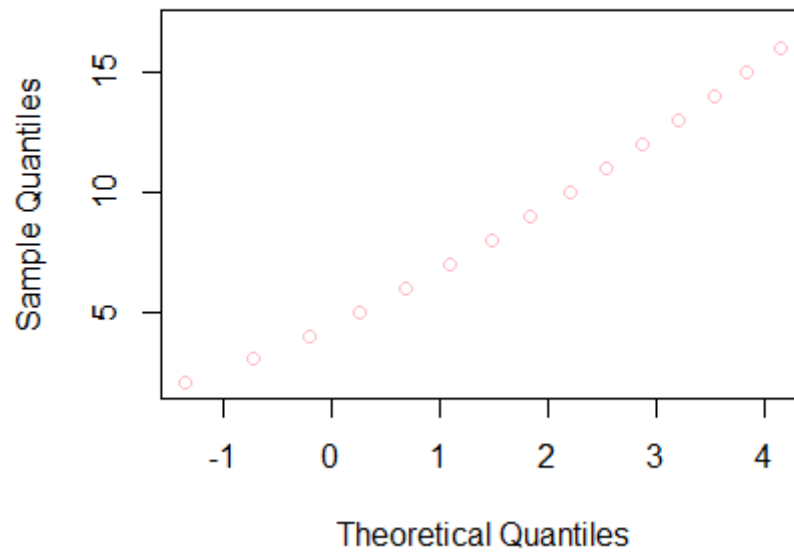
Q-Q of filtered nyt3, nyt4, nyt5, nyt6



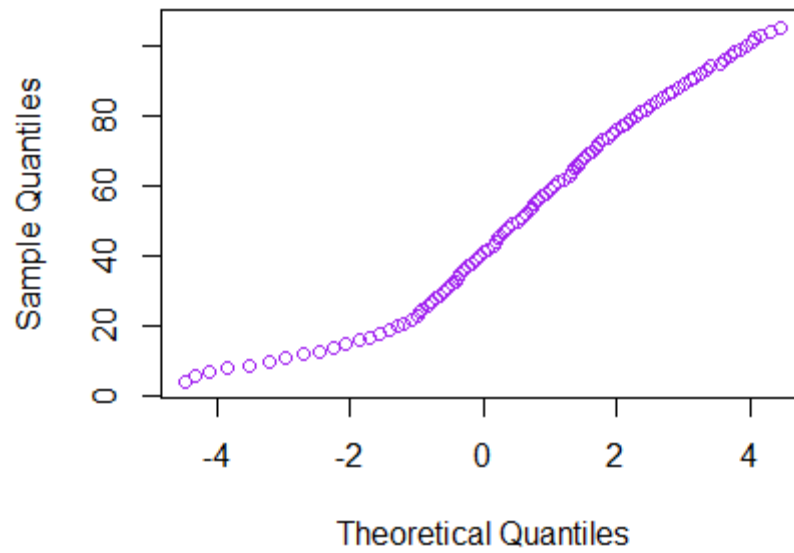
Q-Q Plot for filtered_nyt4 Age



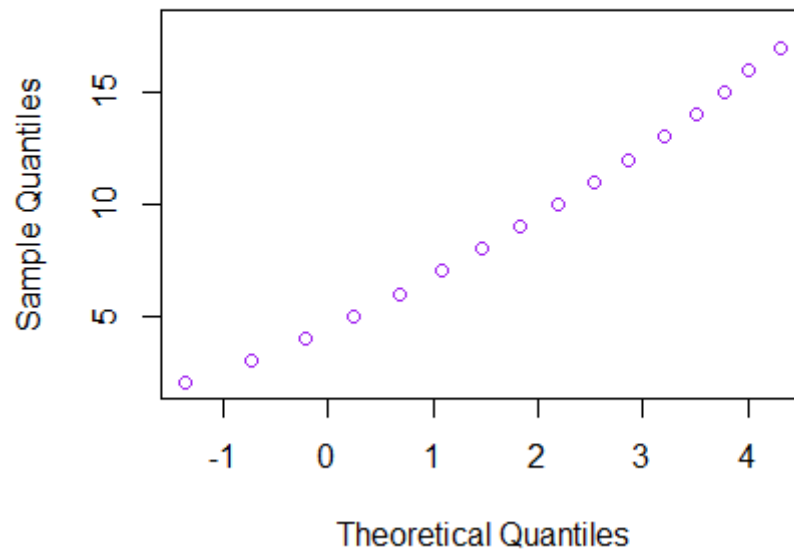
Q-Q Plot for filtered_nyt4 Impressions



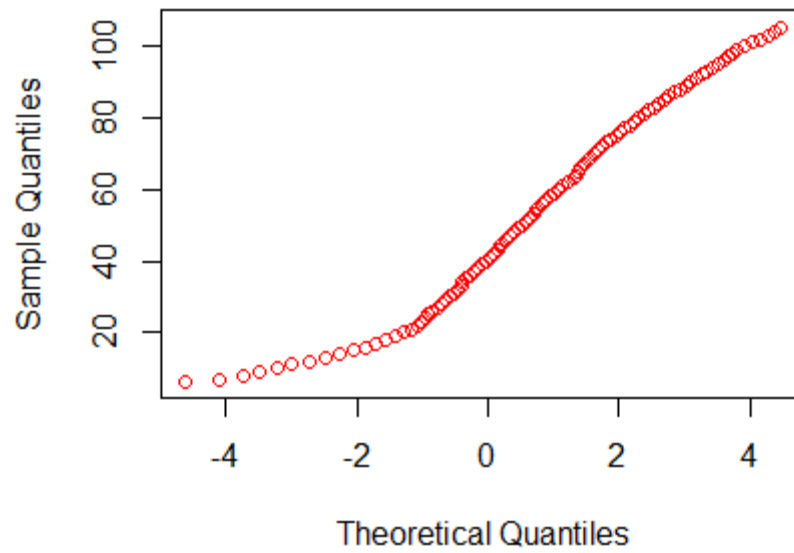
Q-Q Plot for filtered_nyt5 Age



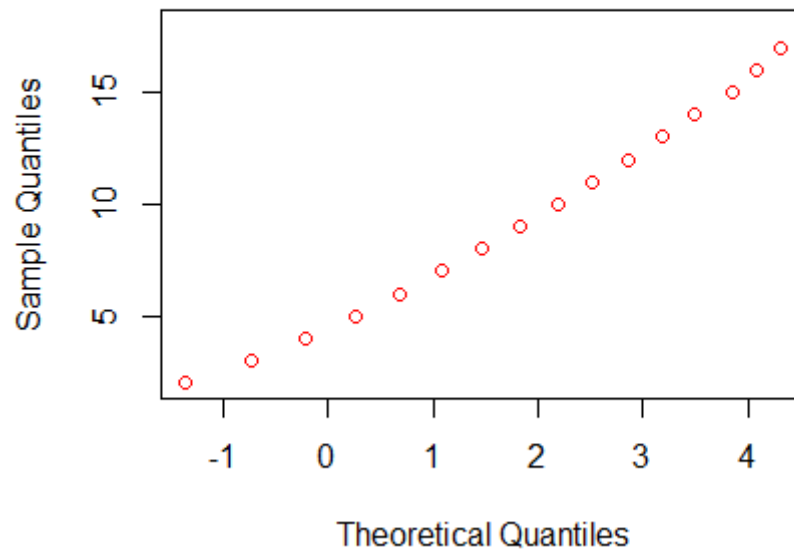
Q-Q Plot for filtered_nyt5 Impressions



Q-Q Plot for filtered_nyt6 Age



Q-Q Plot for filtered_nyt6 Impressions



In summary, the age distributions within all four datasets (nyt3, nyt4, nyt5, and nyt6) exhibit right-skew, with a central tendency indicated by a mean age of 42-43 years. Similarly, the impressions distributions in these datasets display a right-skewed pattern, with a median value hovering around 5 across the board.

Also, observed from the histograms, most of the individuals around the age of forty had greater number of impressions compared to individuals of other age groups

Furthermore, the results of Anderson-Darling normality tests yield low p-values (less than 0.05) for both age and impressions, providing evidence that these datasets do not adhere to a normal distribution. Consequently, I confidently reject the null hypothesis.