

Data Assignment - 1

Group - 19

-SHEFALI MITTAL 2019389
HARSHIT CHOUDHARY 2021254
KANISHK KUKREJA 2021393
YASHILA ARORA 2021436

5) Tables:

Amount.of.Arsenic

Var 1	Fre q	
1	0	7785
2	0.167	1
3	0.5	2
4	1.5	3
5	1.595	1
6	1.8	1
7	2.189	1

8	2.31	1
9	2.5	4
10	2.548	1
11	2.786	2

Amount.Of.carbonate

Am ou nt. of. car bo nat e	Fre q
---------------------------------------------------------------	------------------

1	0	74992
2	0.002	5
3	0.01	10
4	0.02	1

5	0.1	5
6	0.2	1
7	0.3	1
8	0.6	1
9	1	14
10	1.2	4
11	2	4

Amount .of. Calcium

Var	Fre
1	q

1	0	35
2	0.01	5
3	0.7	1

4	0.8	34
5	1.2	11
6	1.4	1
7	1.5	2
8	1.6	150
9	1.7	1
10	1.8	1
11	1.9	1

Amount.of.Chloride

Am ou nt. Of. chl ori de	Fre q
---------------------------------------------------------	------------------

--	--	--	--

1	0	26766 1
2	0.2	1
3	0.3	4
4	0.4	12
5	0.5	14
6	0.56	1
7	0.6	11
8	0.7	9
9	0.8	14
10	0.9	8
11	1	7

Amount.of.Electrical.Conductivity

Amou nt.of. Electri cal.Co nducti vity	Fre q
-------------------------------------------------------	----------

1	0	2
2	2	2
3	4	1
4	10	1
5	13	1
6	21	1
7	22	3
8	23	1
9	24	4

10	25	5
11	26	5

Amount.of.Fluorine

Amou nt.of. Flouri ne	Fre q
--------------------------------	----------

1	0	2224
2	1e-04	6
3	2e-04	10
4	6e-04	3
5	7e-04	4
6	8e-04	4
7	0.001	9

8	0.002	15
9	0.003	15
10	0.003707 824	2
11	0.0039	6

Amount.of.Iron

Var 1	Fre q
----------	----------

1	0	8310
2	0.001	48
3	0.0016	1
4	0.0017	1
5	0.002	29
6	0.0021	1

7	0.0023	1
8	0.003	16
9	0.0034	1
10	0.0035	1
11	0.004	33

Amount.of.HydrogenCarbonate

Am ou nt. of. Hy dro ge nc arb on ate	Fre q
---------------------------------------------------------------------	----------

1	0	708
---	---	-----

2	0.01	1
3	0.5	1
4	1.59	1
5	1.7	1
6	1.9	1
7	2.07	1
8	2.1	1
9	2.2	1
10	2.3	1
11	2.32	1

Amount.of.Pottasium

Am ou nt. of. Pot	Fre q
------------------------------------------	------------------

tas iu m		
----------------	--	--

1	0	185
2	0.001	2
3	0.002	1
4	0.01	60
5	0.0102	1
6	0.02	60
7	0.03	57
8	0.04	56
9	0.05	61
10	0.051	1

11	0.06	73
----	------	----

Amount.Of.Magnesium

Amount.of. Magnesium	Frequency
-------------------------	-----------

1	-67	1
2	-22	1
3	-4.8608	1
4	-2.4304	3
5	0	69
6	8.64e-14	1
7	0.01	42
8	0.02	24

9	0.03	7
10	0.04	13
11	0.05	6

Amount.of.Nitrate

Am ou nt. of. Nit rat e	Fre q
-------------------------------------------	----------

1	0	5181
2	0.001	1
3	0.002	1
4	0.003	1
5	0.005	1

6	0.01	76
7	0.011	1
8	0.014	1
9	0.0145	1
10	0.015	1
11	0.017	1

Amount.of.Sodium

Am ou nt. of. So diu m	Fre q
------------------------------------------	----------

1	0	33
2	0.01	1

3	0.03	1
4	0.05	2
5	0.09	1
6	0.1	5
7	0.104	1
8	0.13	1
9	0.18	1
10	0.2	12
11	0.21	1

Amount.of.Phosphatelon

Amou nt.of. Phos phate lon	Fre q
----------------------------------------	----------

--	--	--	--

1	0	7400
2	0.001	102
3	0.0012	2
4	0.0014	1
5	0.002	18
6	0.003	13
7	0.004	32
8	0.005	4
9	0.006	13
10	0.006497 726	15
11	0.007	19

Amount.of.Sodium

Amou nt.of.S odium	Fre q
-----------------------------------	------------------

1	-230.78	1
2	-196.609 7471	1
3	-193.22	1
4	-192.8	1
5	-180.4	1
6	-177.2	1
7	-171.465 3191	1
8	-156.4	1
9	-156.263 5943	1
10	-156	1

11	-152.469 1737	1
----	------------------	---

Amount.of.sulfate

Amou nt.of.S ulfate	Fre q
---------------------------	----------

1	-0.46620 6751	1
2	-0.24149 3149	1
3	0	3845
4	0.001	3
5	0.01	20
6	0.014327 283	1
7	0.02	2

8	0.023838 602	1
9	0.03	3
10	0.04	10
11	0.05	21

Amount.of.Hardness.Total

Am ou nt. of. Ha rdn es s.T ota l	Fre q	
--------------------------------------------------------------	----------	--

1	0	151
2	0.72	1
3	1	1

4	1.14	1
5	2	2
6	3	2
7	3.6194 17	1
8	4	19
9	5	12
10	6	81
11	8	179

Amount.of.Alkalinity.Total

Amou nt.of. Alkali nity.T otal	Fre q
--------------------------------------------	----------

1	0	4205
---	---	------

2	0.008196 721	1
3	0.4	1
4	1.3	1
5	1.4	1
6	1.6	1
7	1.7	2
8	1.9	2
9	1.967213 115	5
10	1.97	12
11	2	28

Amount.of.Total.Dissolved.Solids

Amou nt.of.T otal.D	Fre q
------------------------------------	------------------

issolv ed.Sol ids		
----------------------------------	--	--

1	0	6059
2	8.28	1
3	8.74	1
4	12	1
5	16	1
6	17	1
7	18	1
8	18.1	1
9	19	3
10	20	3

11	20.7	1
----	------	---

Amount.of.Potential.of.Hydrogen

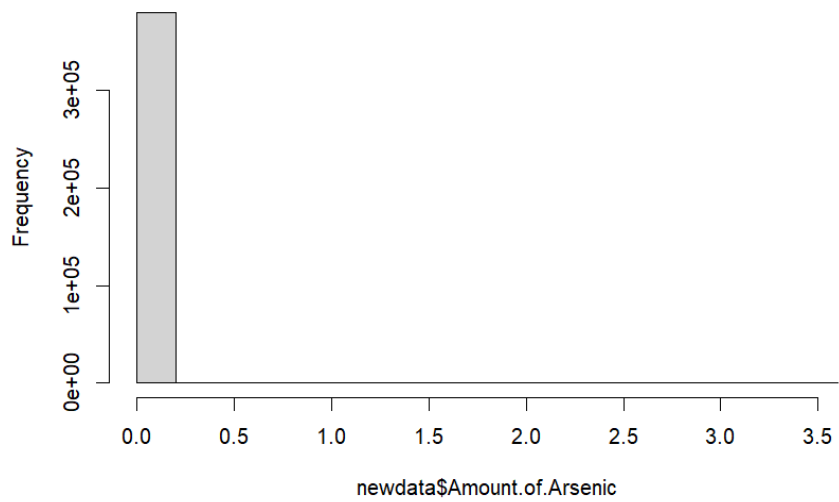
Amou nt.of. Poten tial.of. Hydro gen	Fre q
-----------------------------------------------------	----------

1	0	1
2	2	1
3	3	1
4	3.11	1
5	3.29	1
6	3.36	1
7	3.39	1

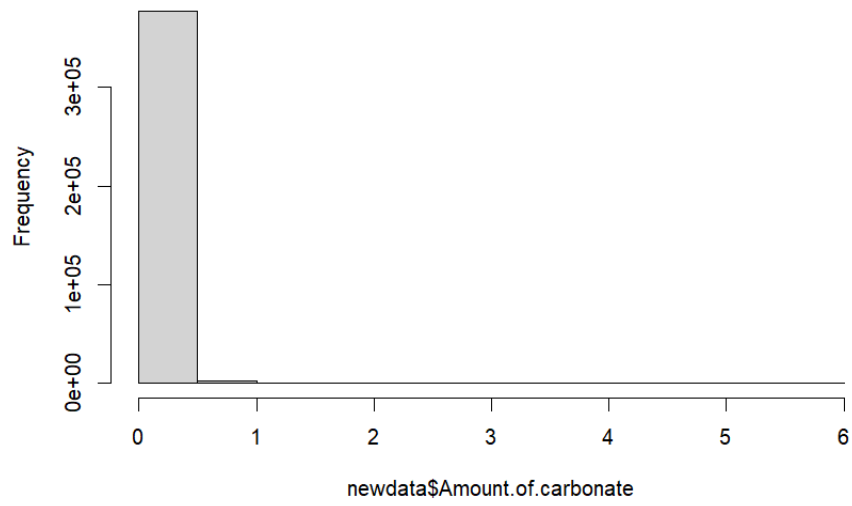
8	3.4	1
9	3.42	1
10	3.46	1
11	3.5	1

Histograms:

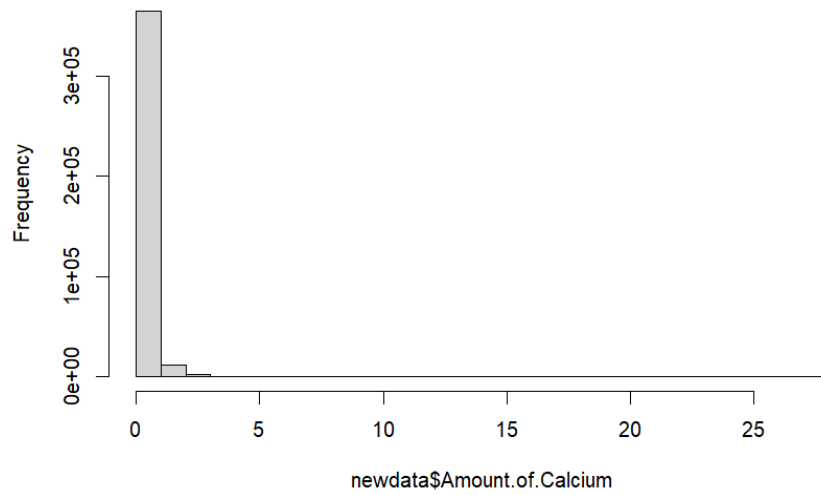
Residuals Histogram



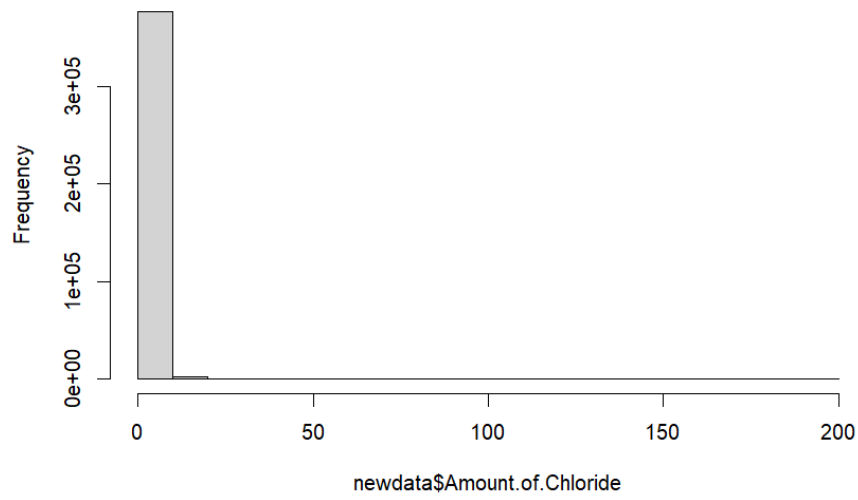
Residuals Histogram



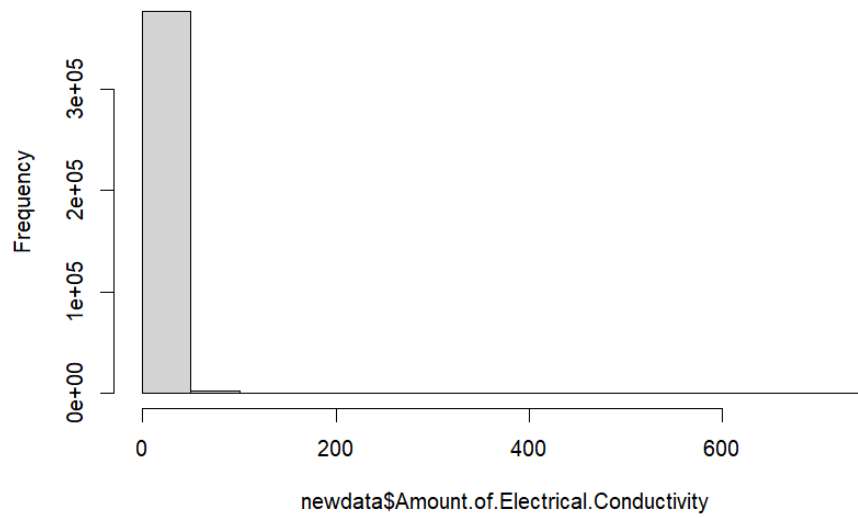
Residuals Histogram



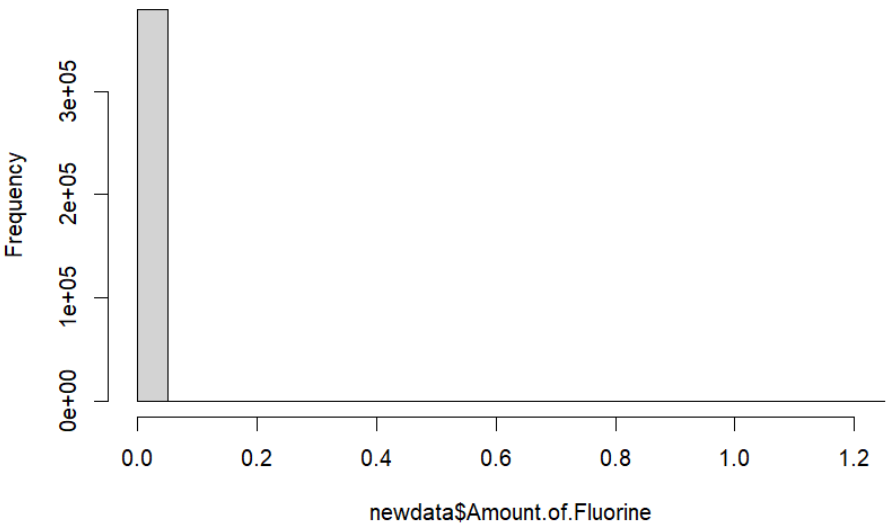
Residuals Histogram



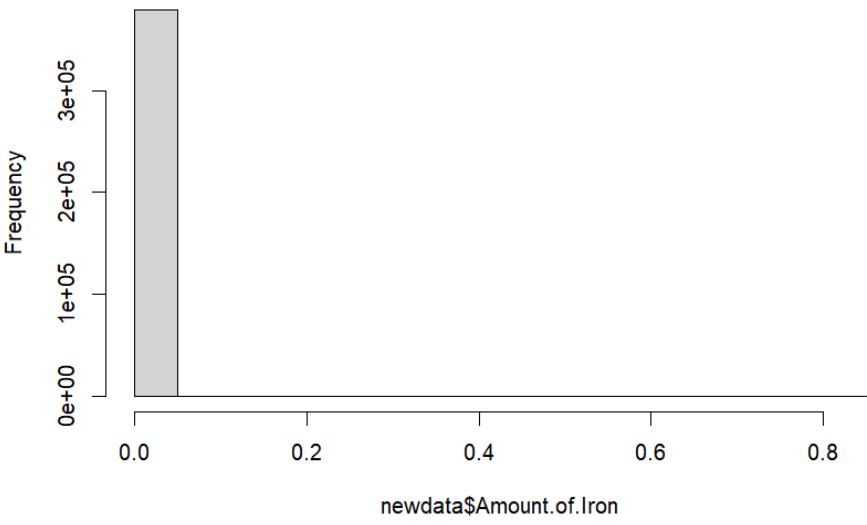
Residuals Histogram



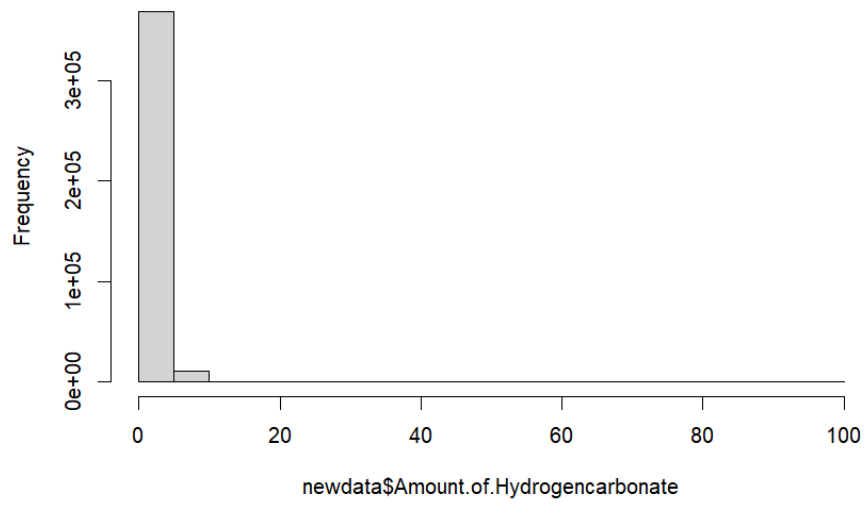
Residuals Histogram



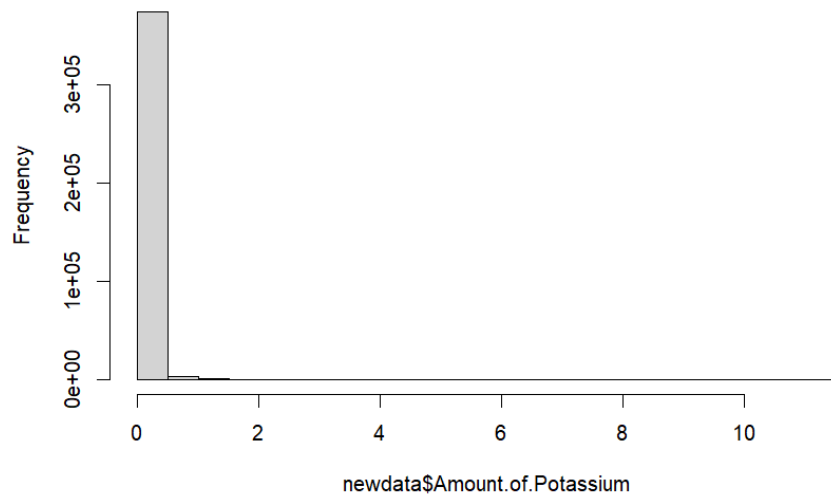
Residuals Histogram



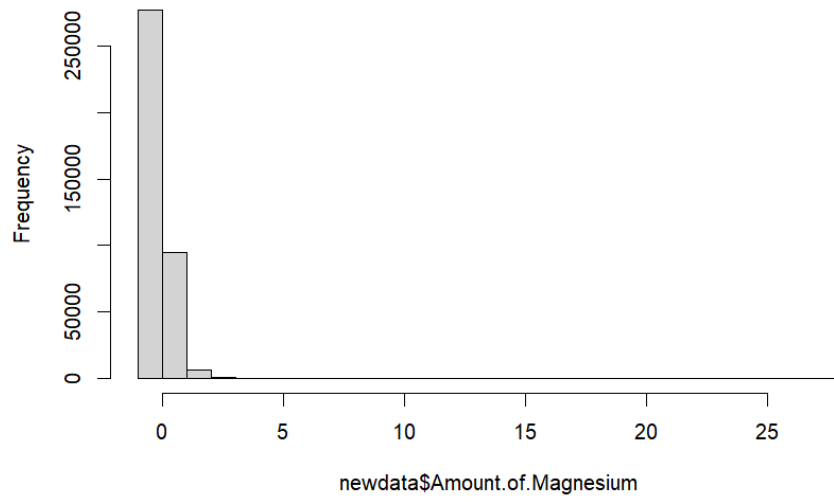
Residuals Histogram



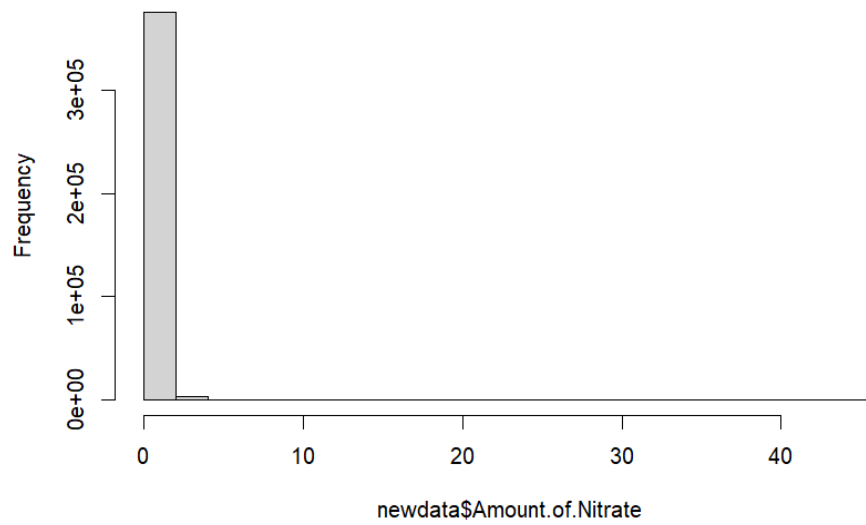
Residuals Histogram



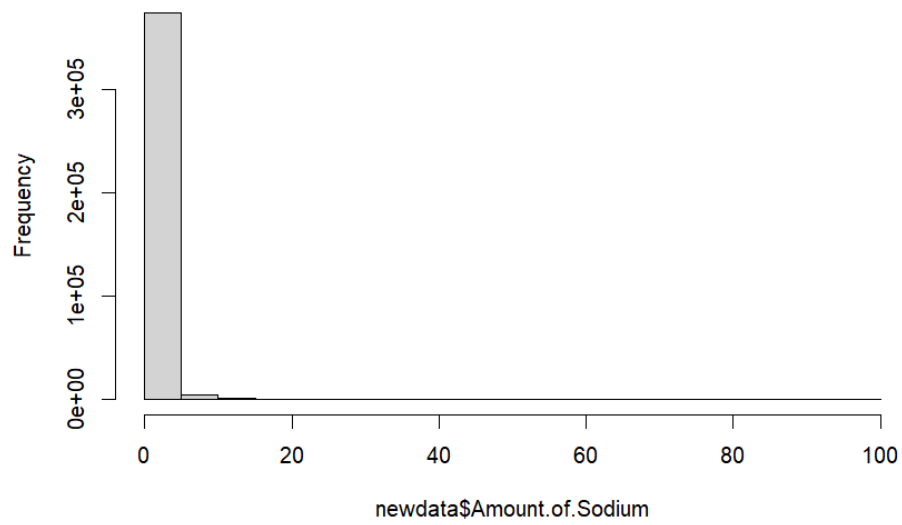
Residuals Histogram



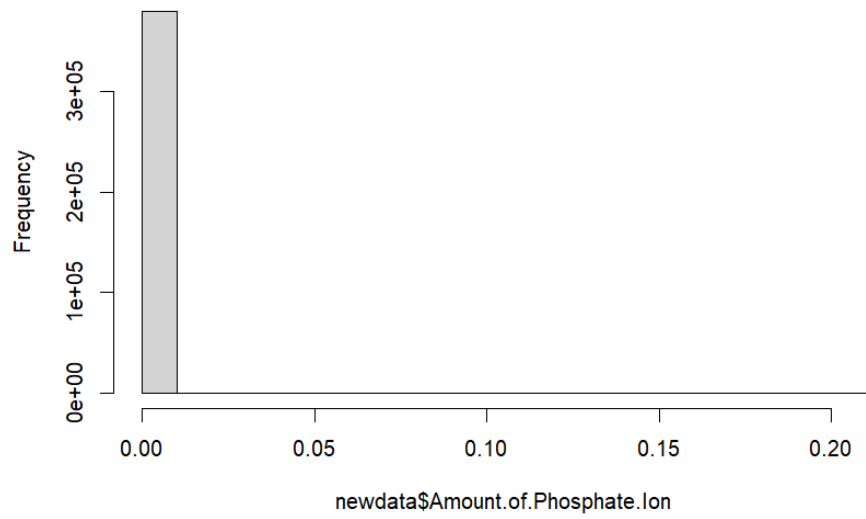
Residuals Histogram

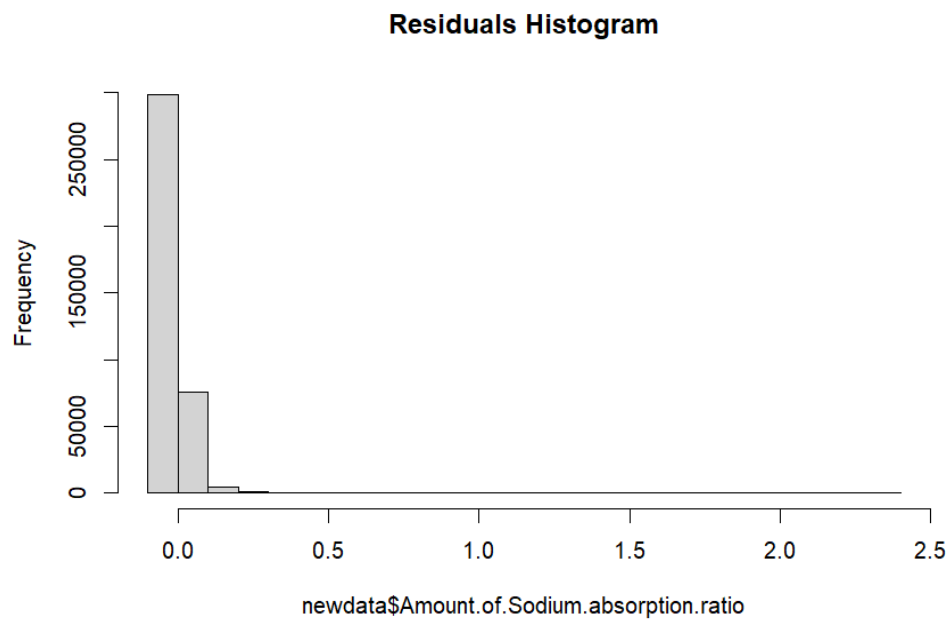
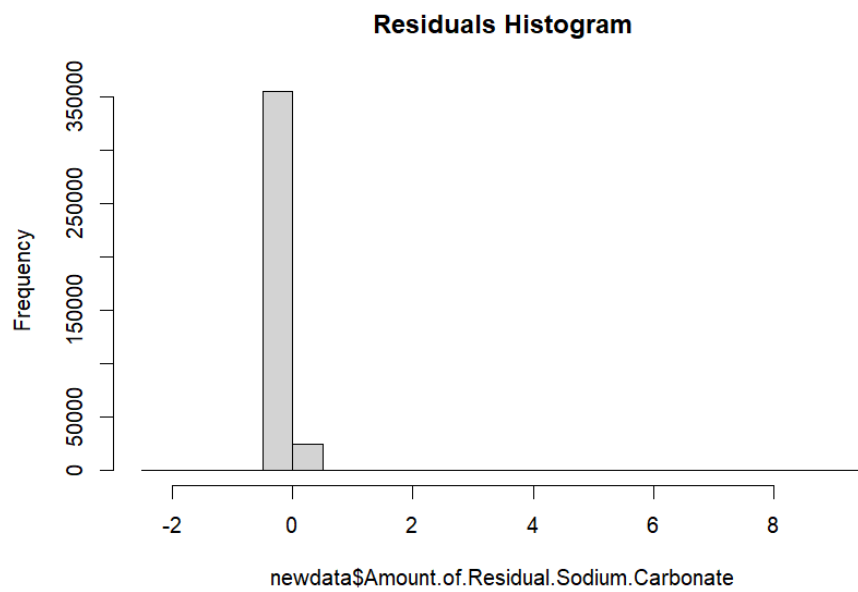


Residuals Histogram

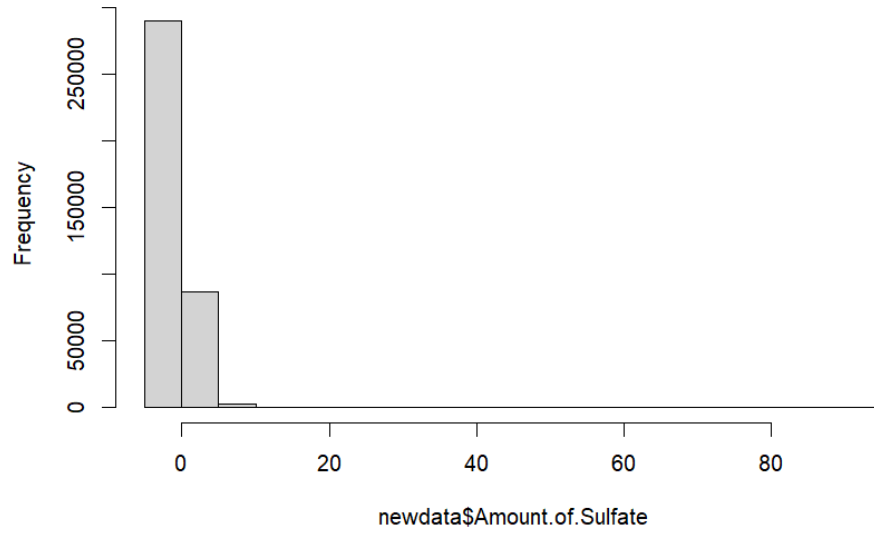


Residuals Histogram

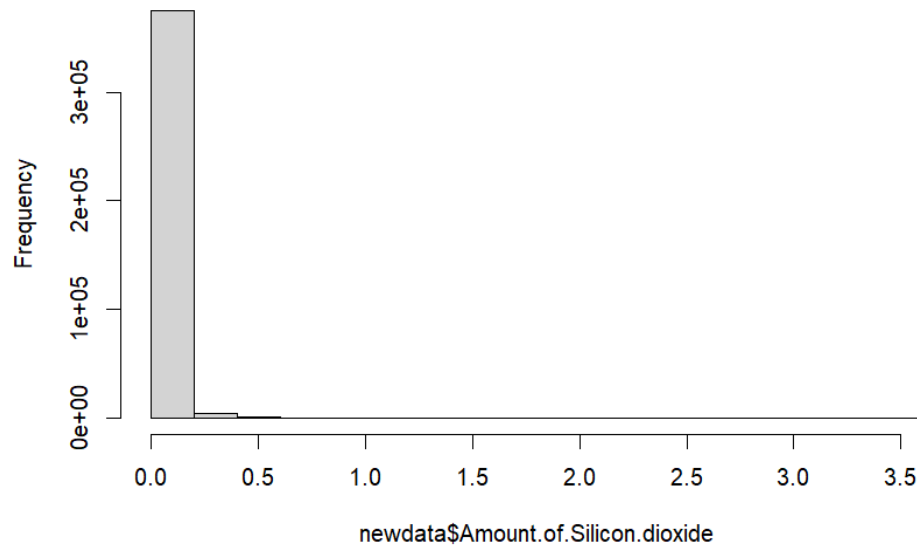




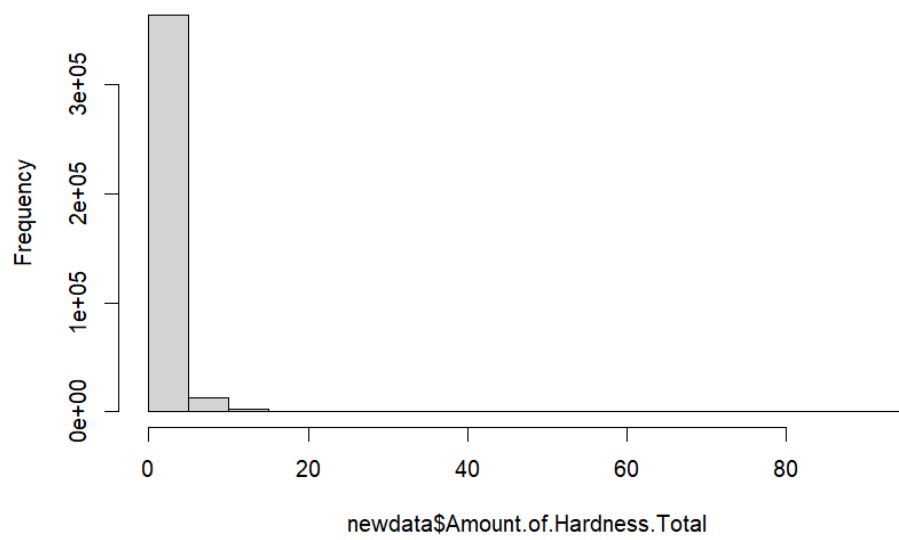
Residuals Histogram



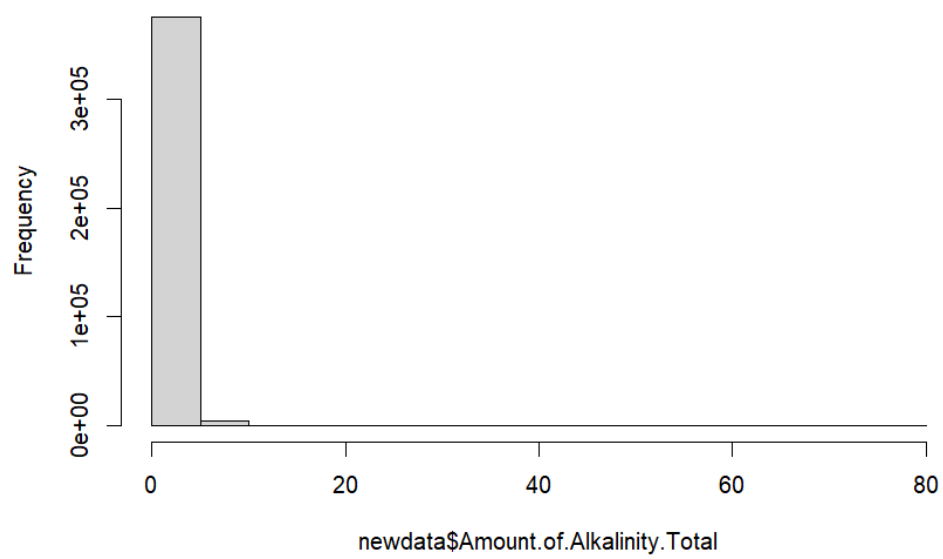
Residuals Histogram



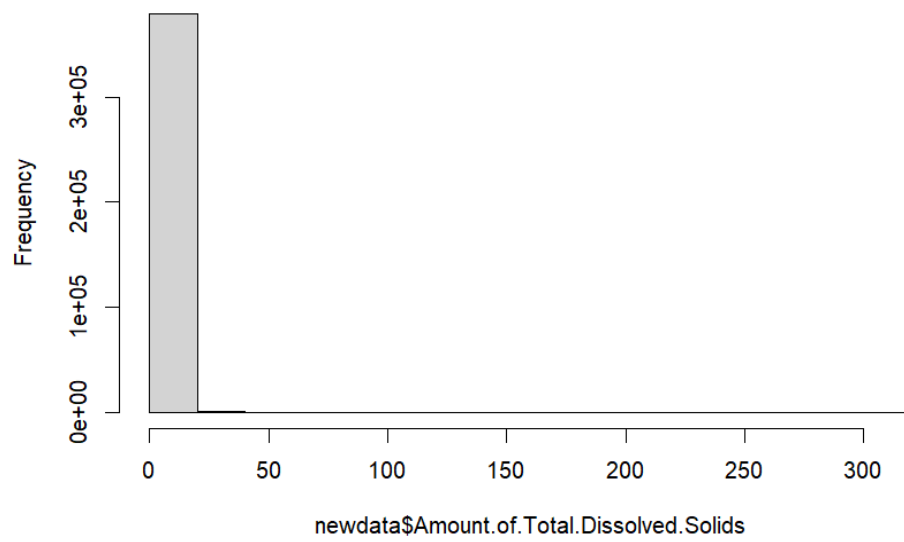
Residuals Histogram



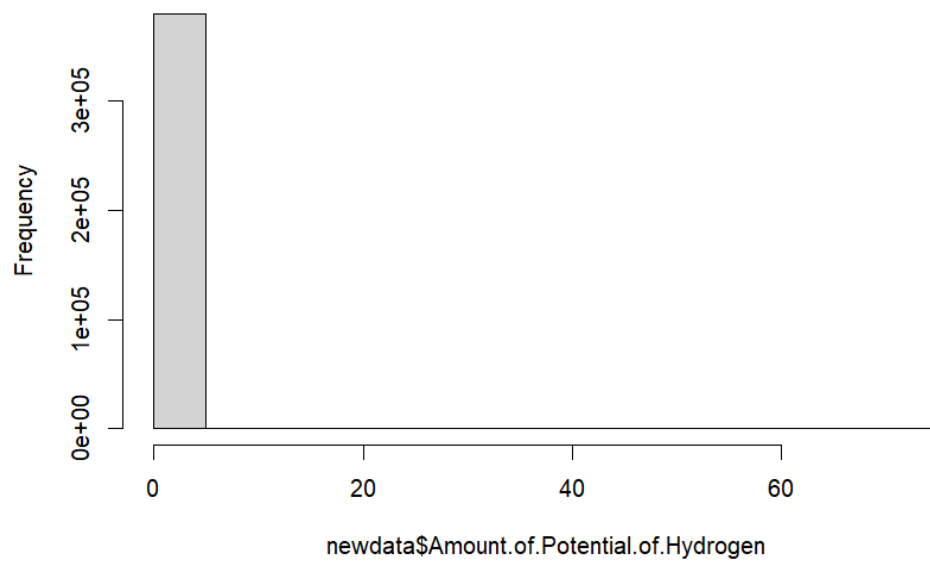
Residuals Histogram



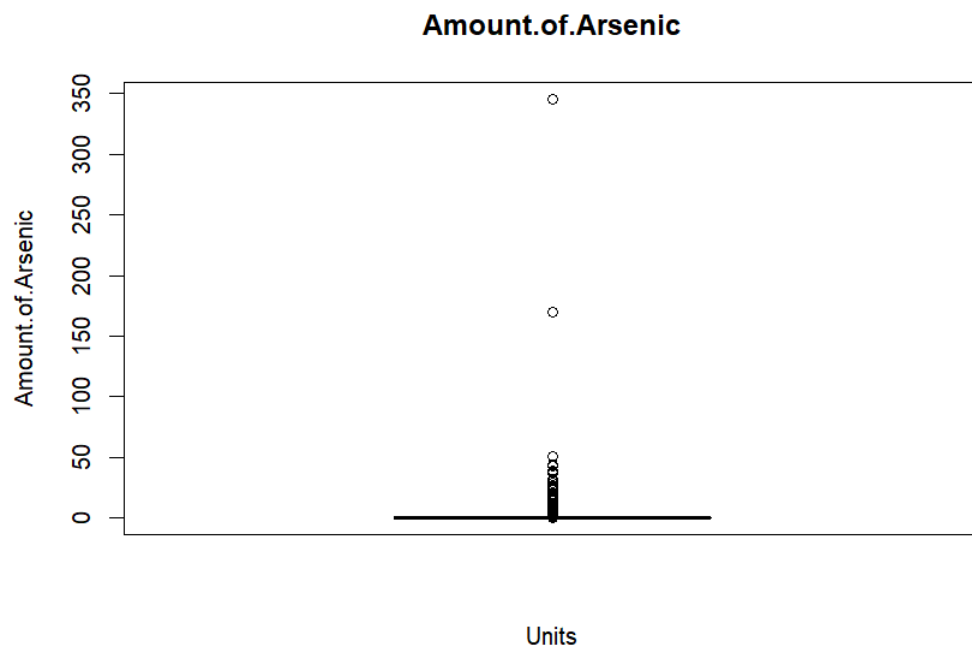
Residuals Histogram



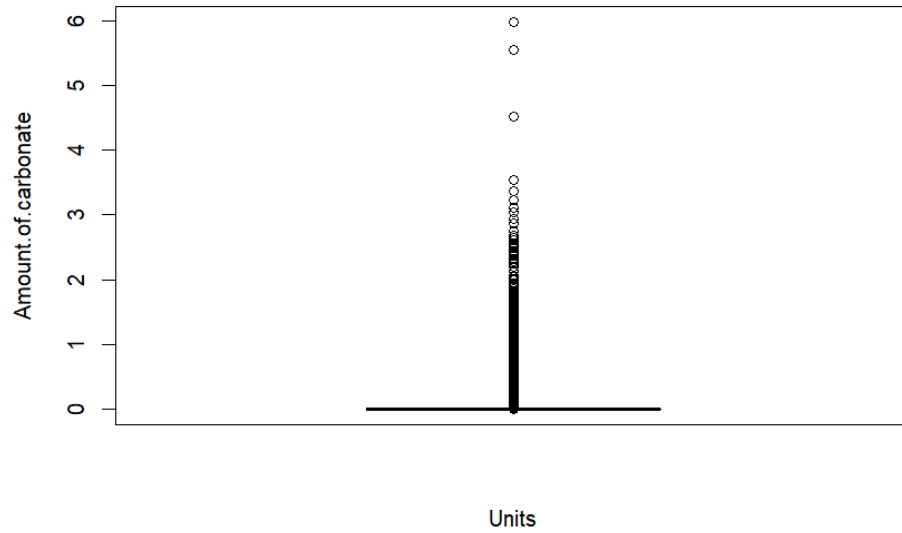
Residuals Histogram



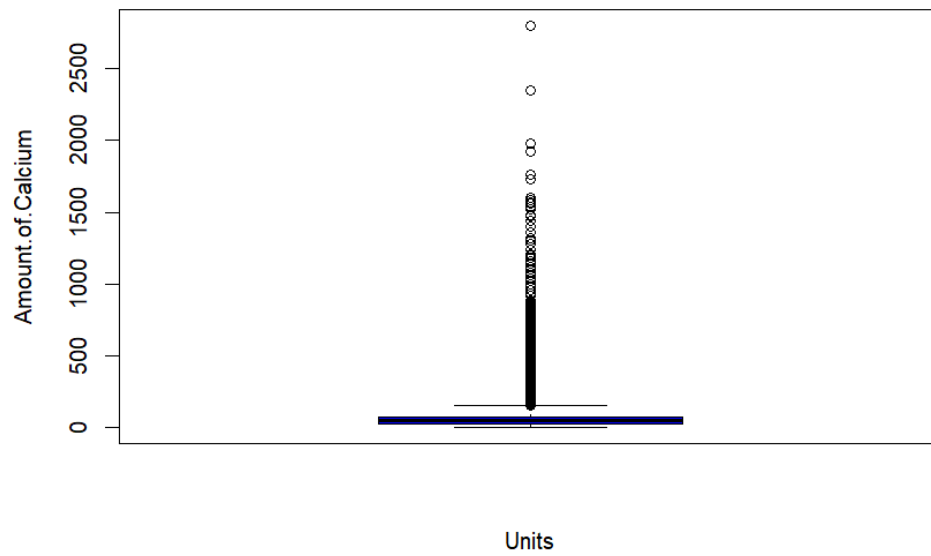
Box-Plots:

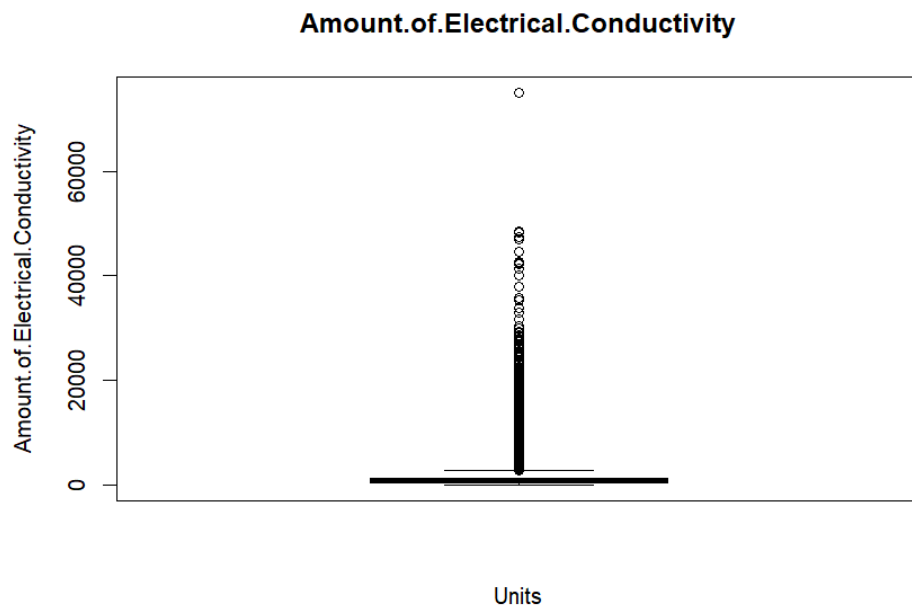
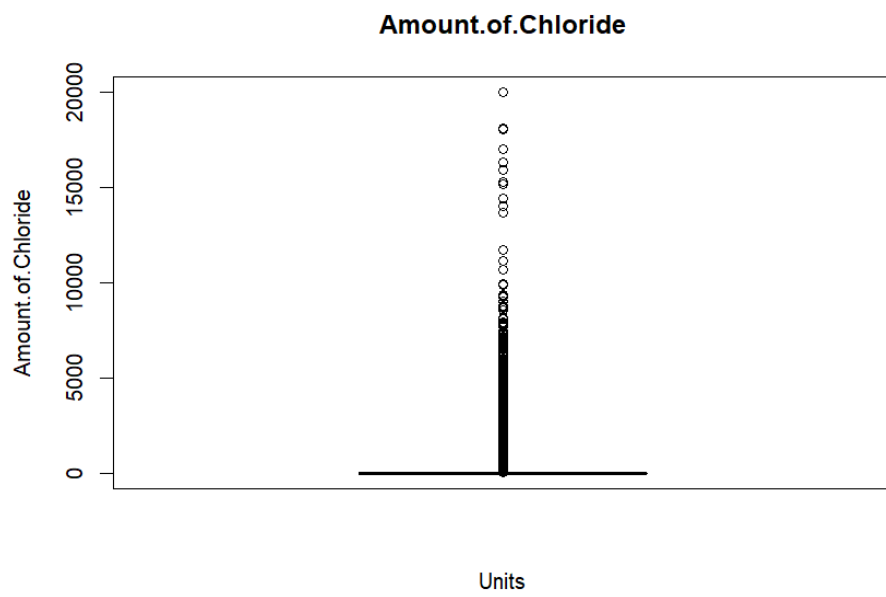


Amount.of.carbonate

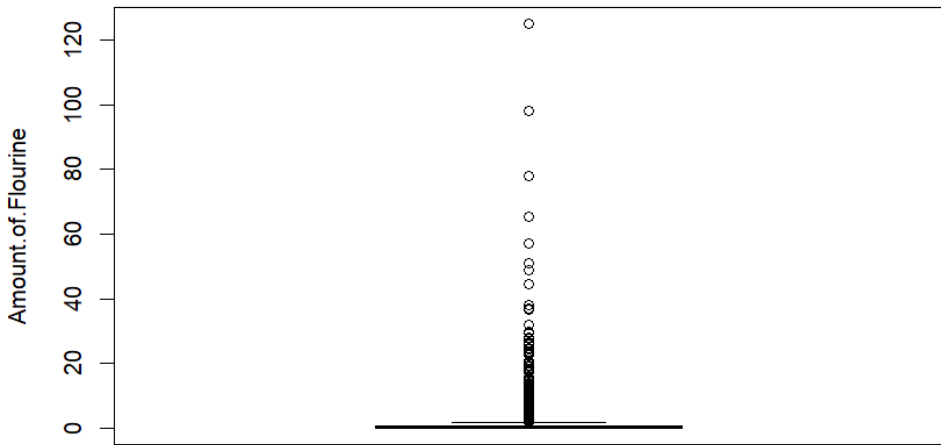


Amount.of.Calcium



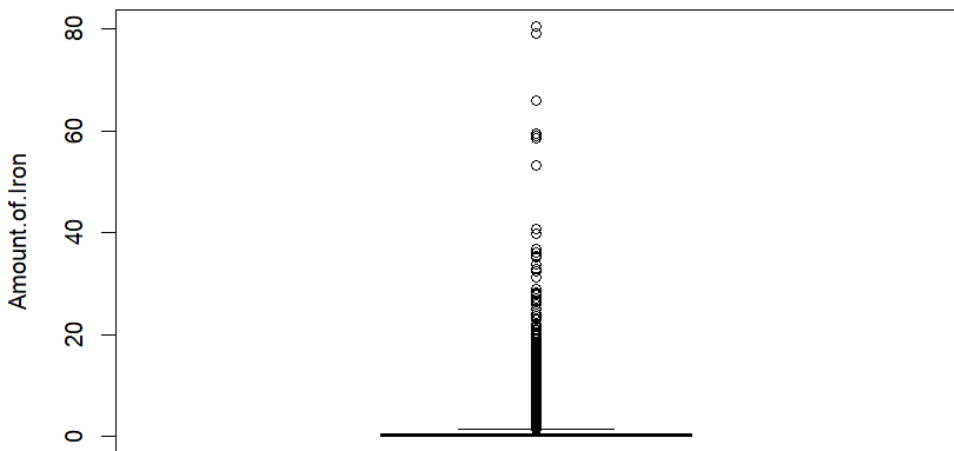


Amount.of.Flourine



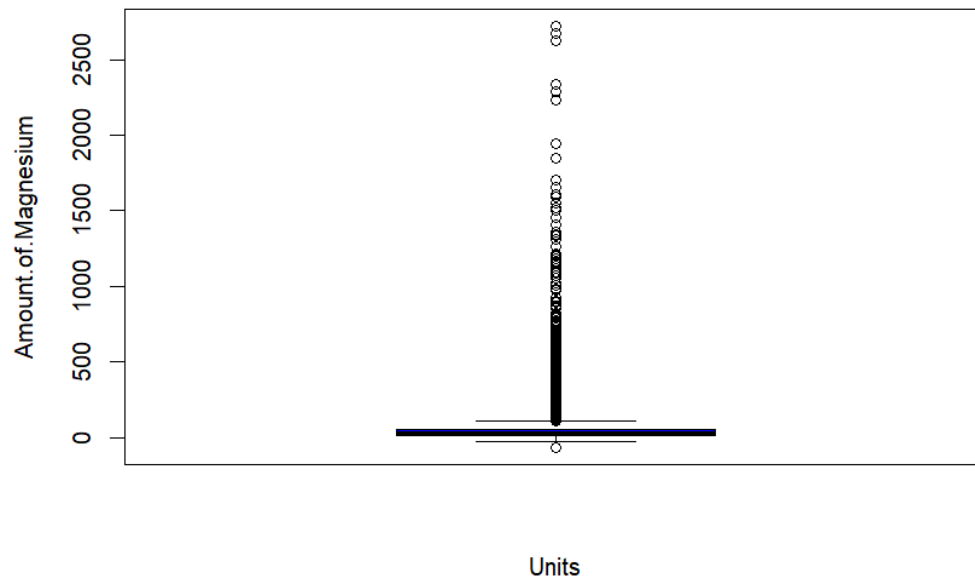
Units

Amount.of.Iron

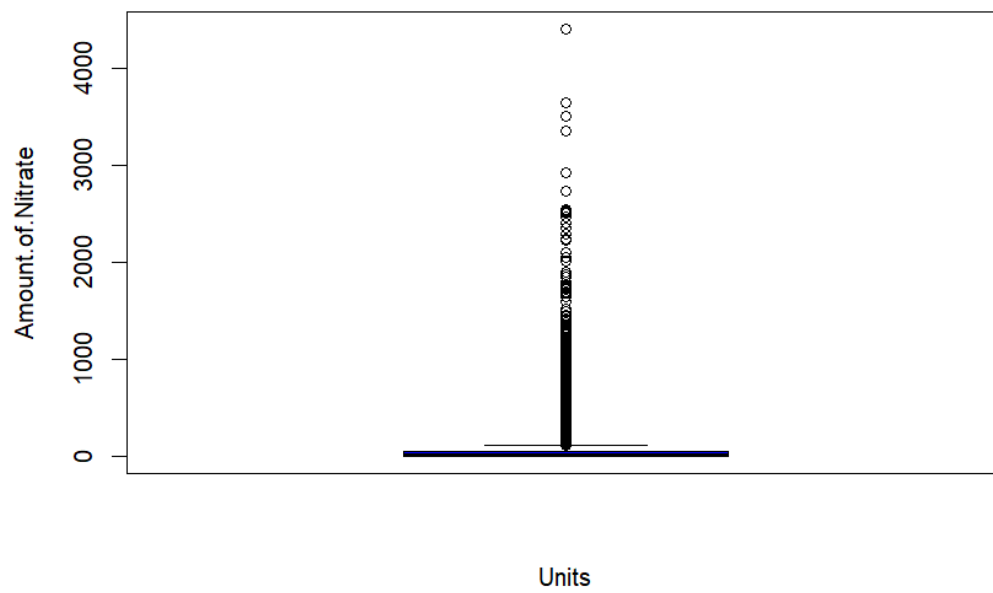


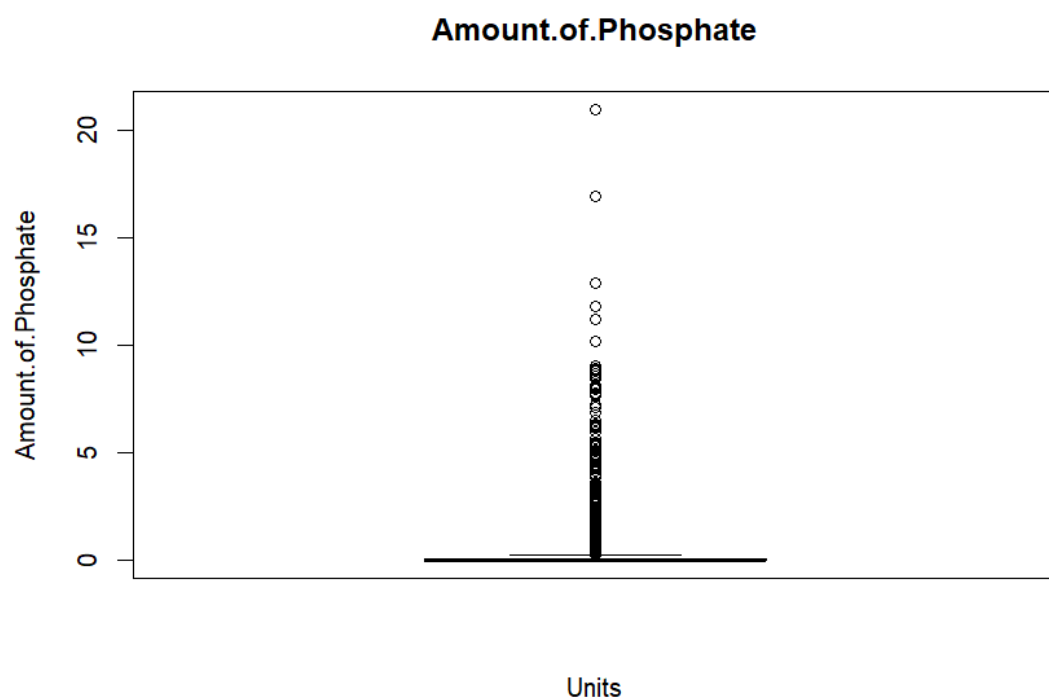
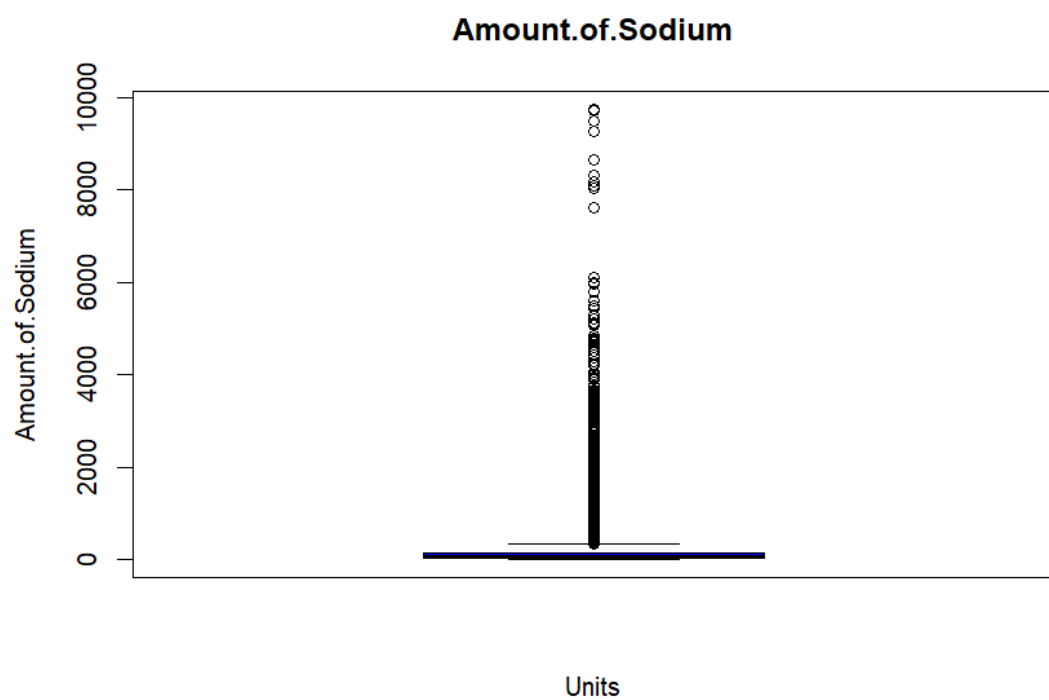
Units

Amount.of.Magnesium

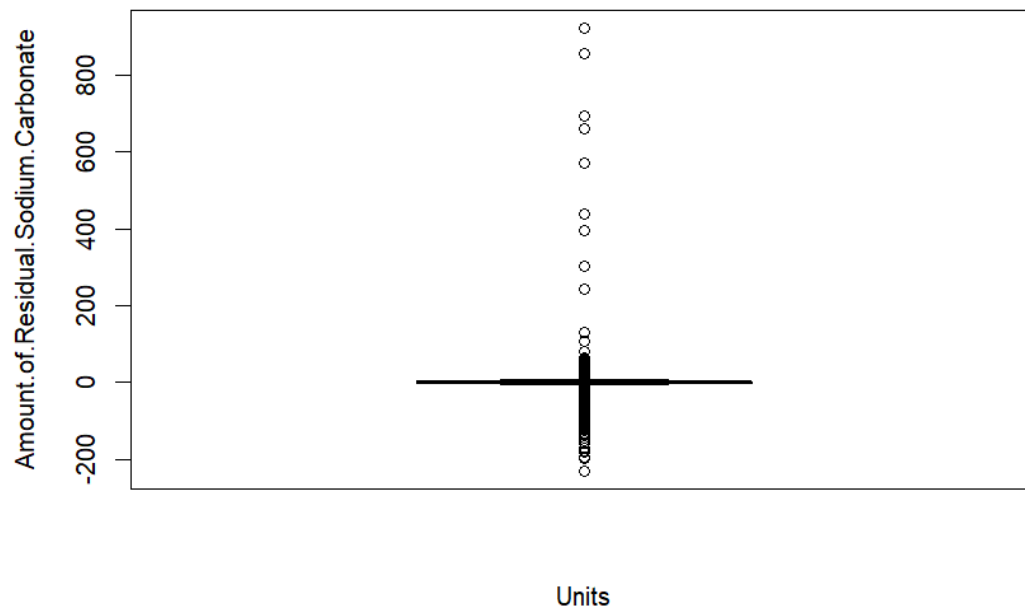


Amount.of.Nitrate

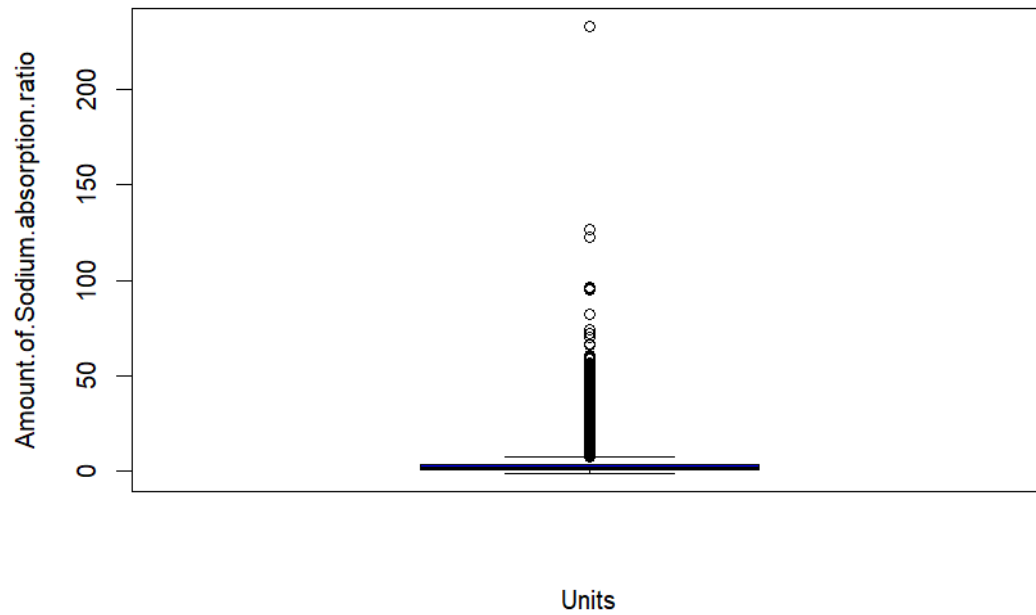




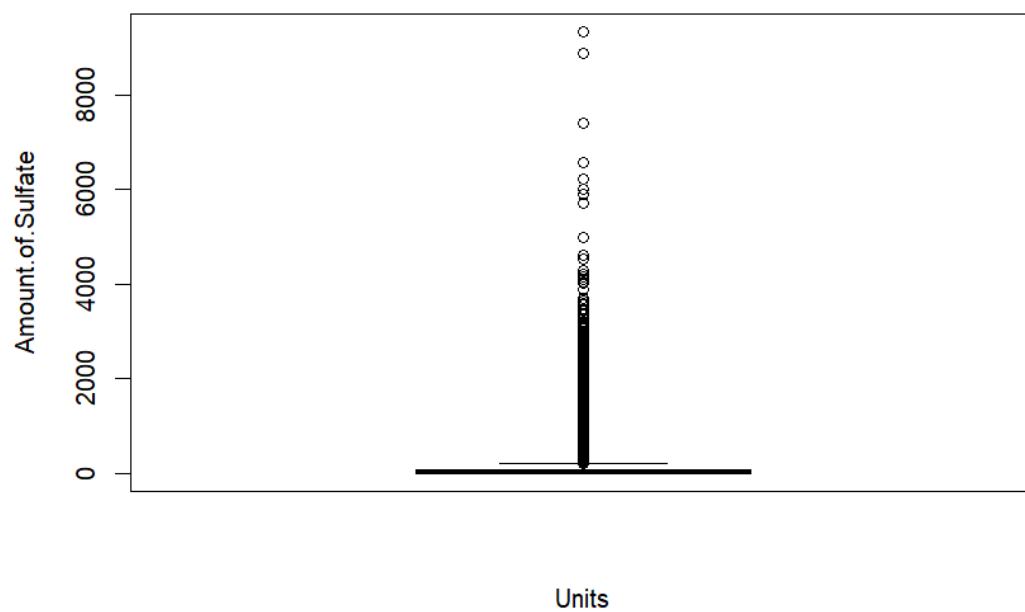
Amount.of.Residual.Sodium.Carbonate



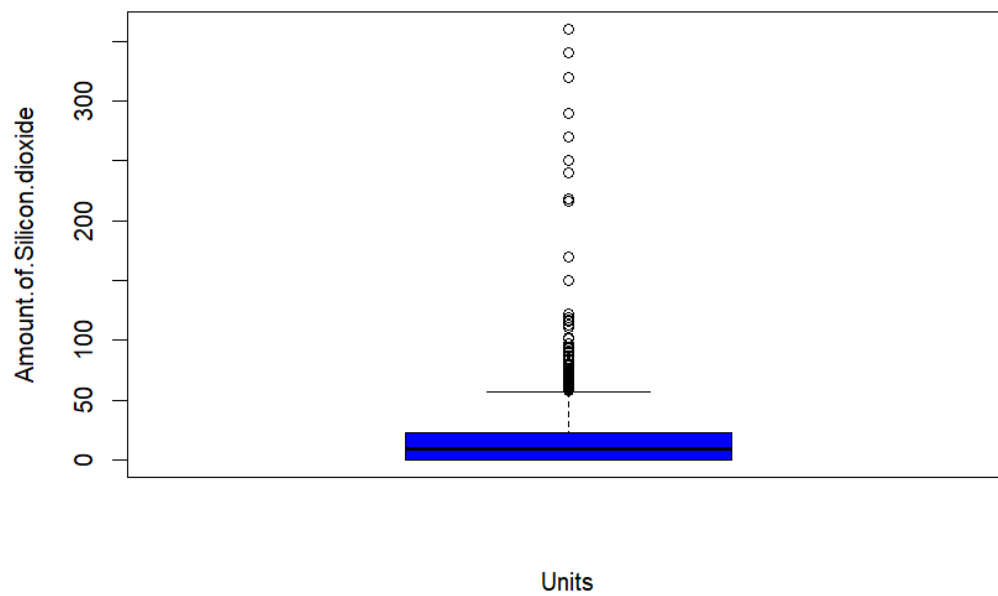
Amount.of.Sodium.absorption.ratio



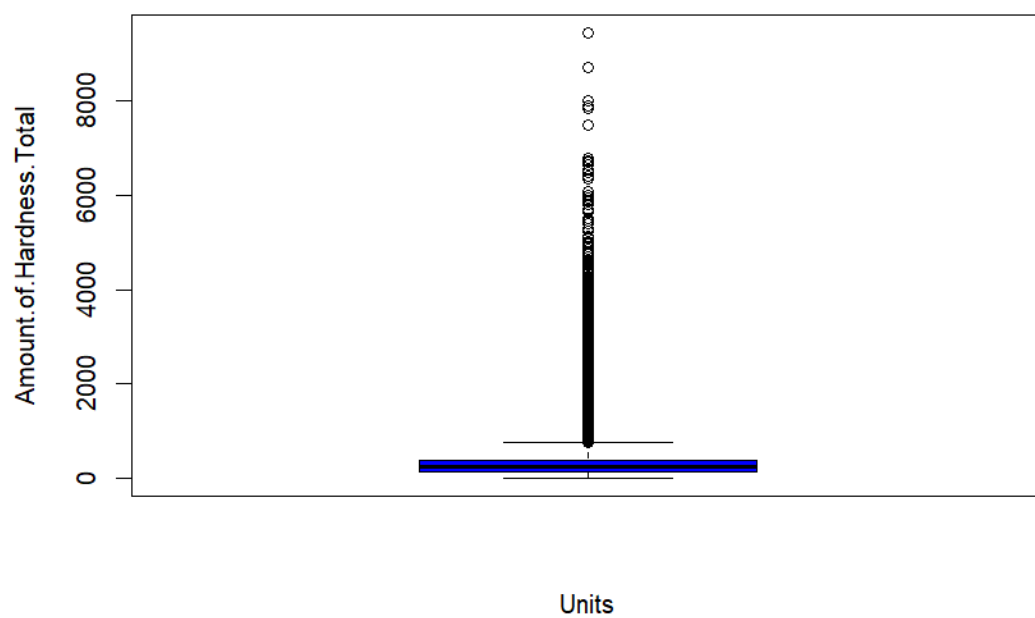
Amount.of.Sulfate



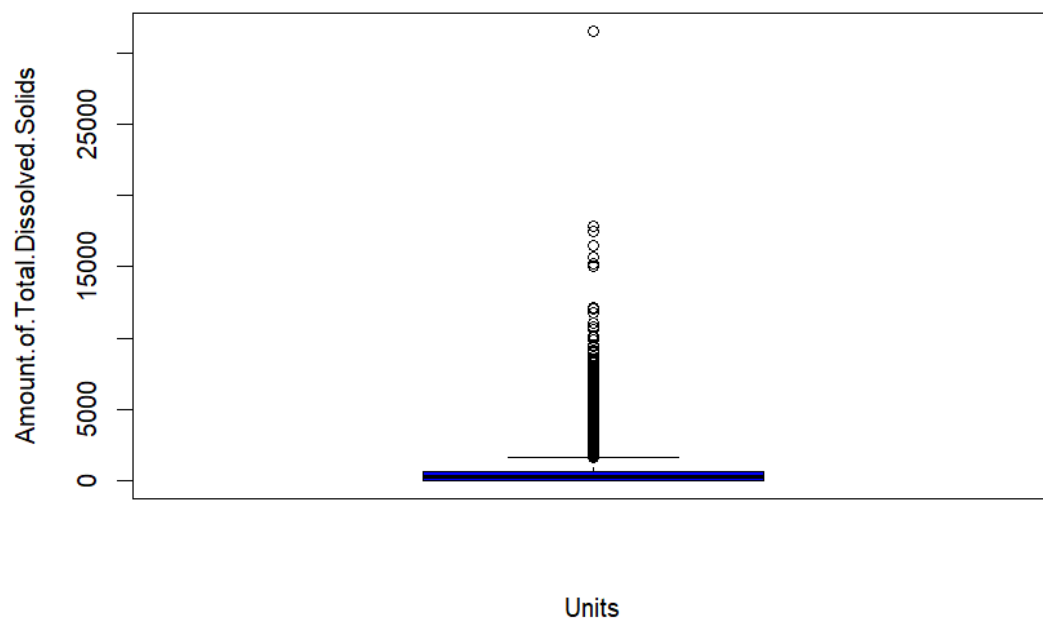
Amount.of.Silicon.dioxide

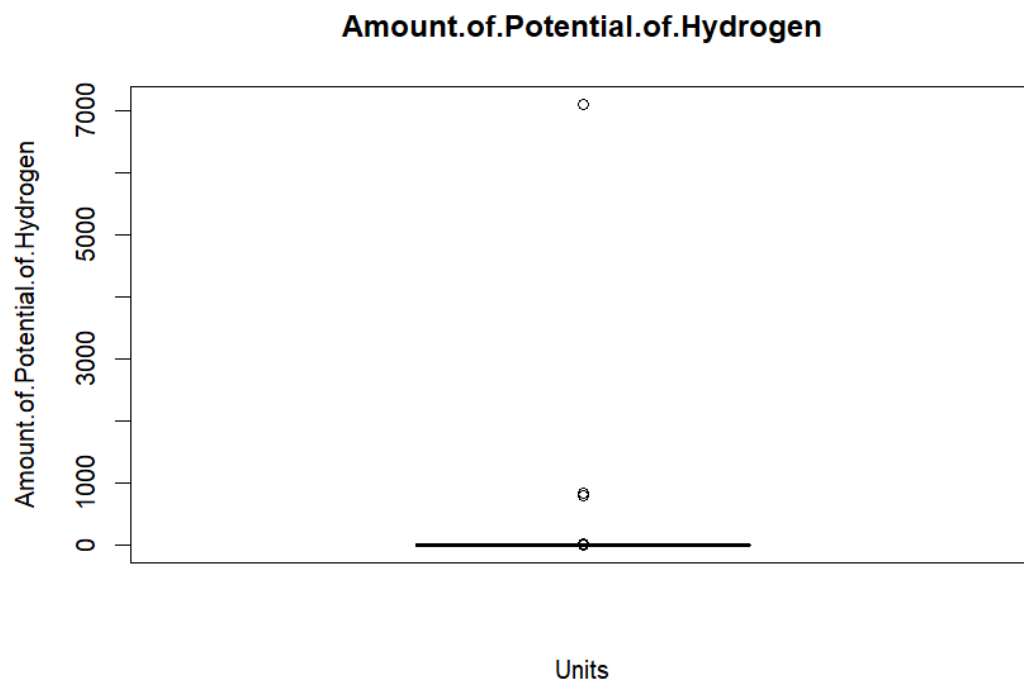


Amount.of.Hardness.Total

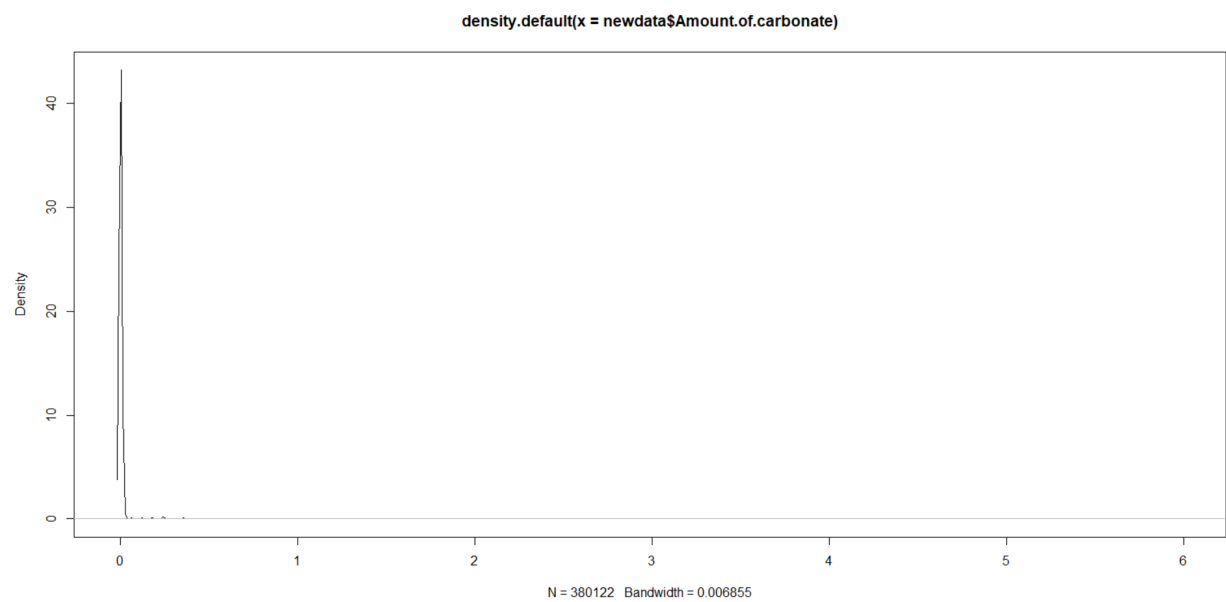
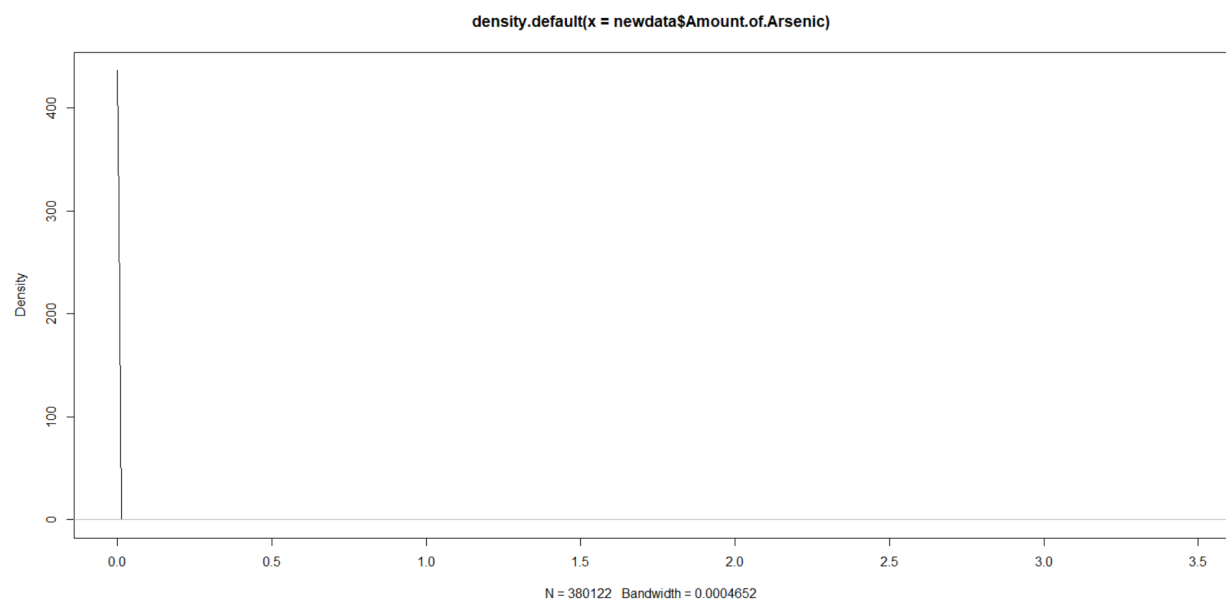


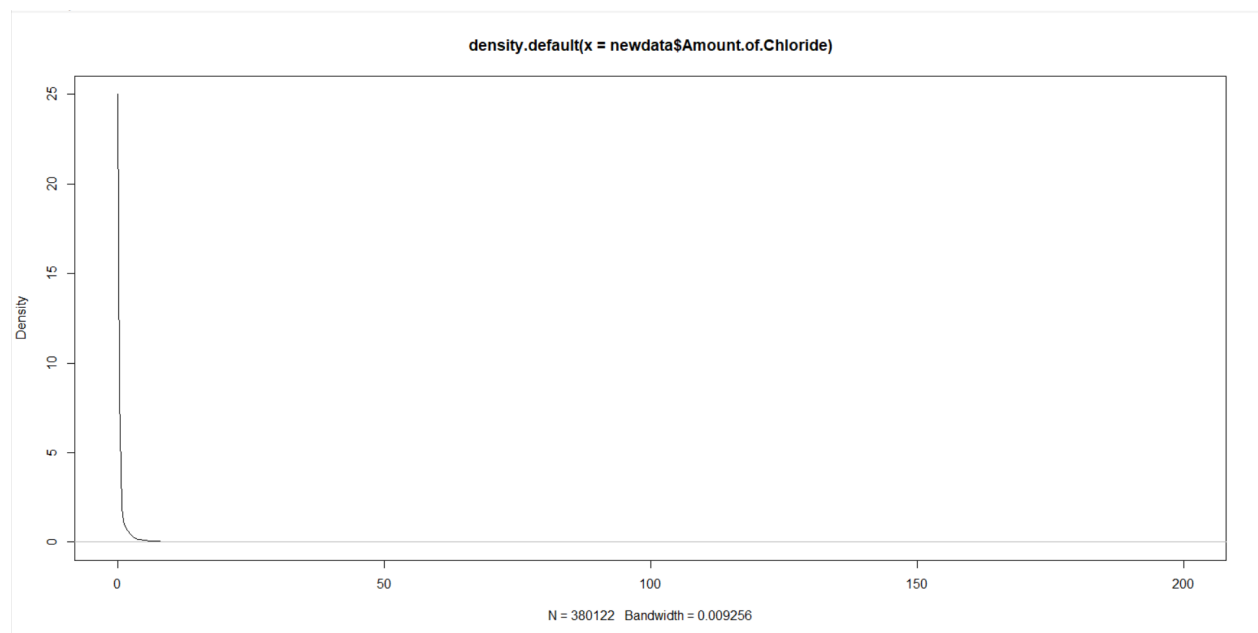
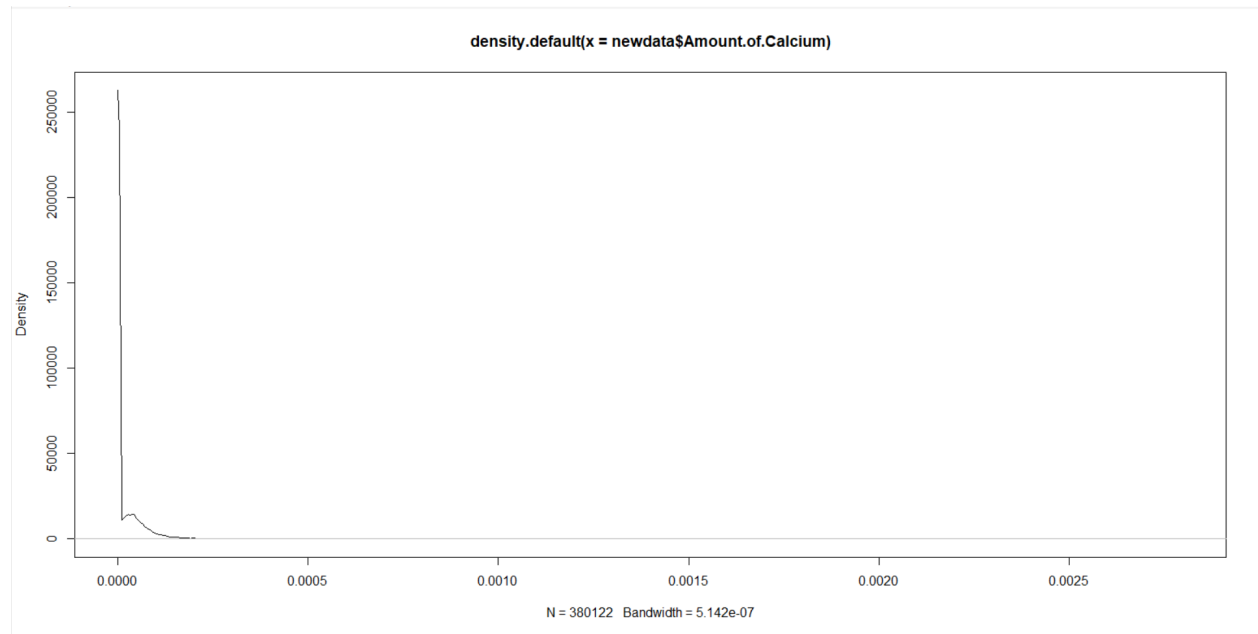
Amount.of.Total.Dissolved.Solids

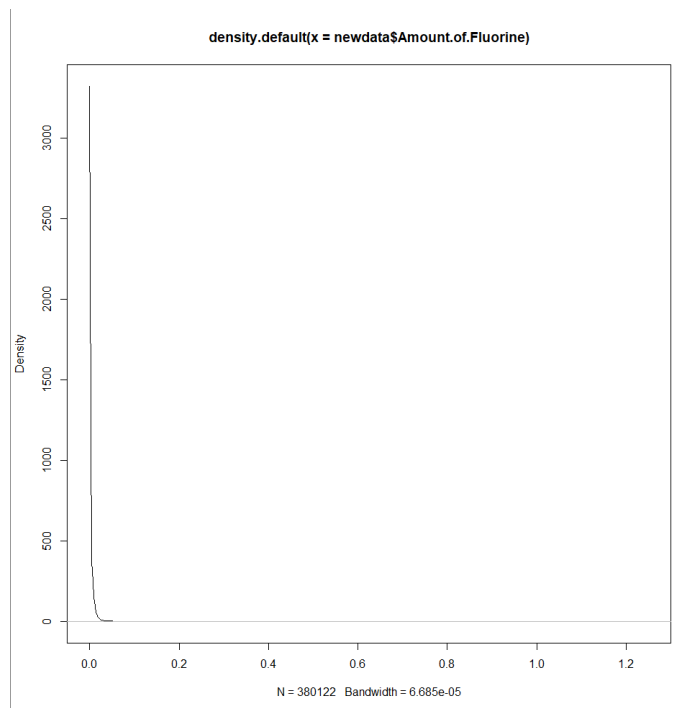
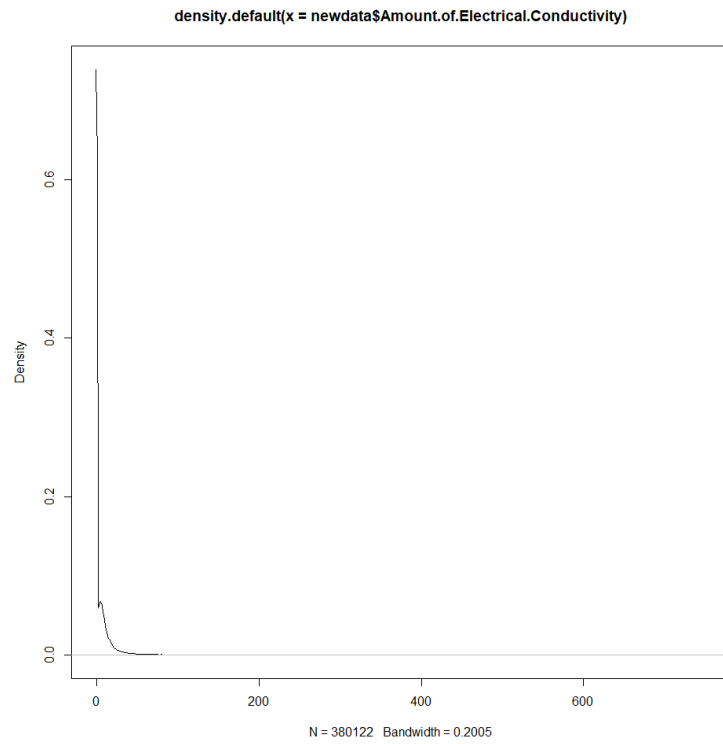


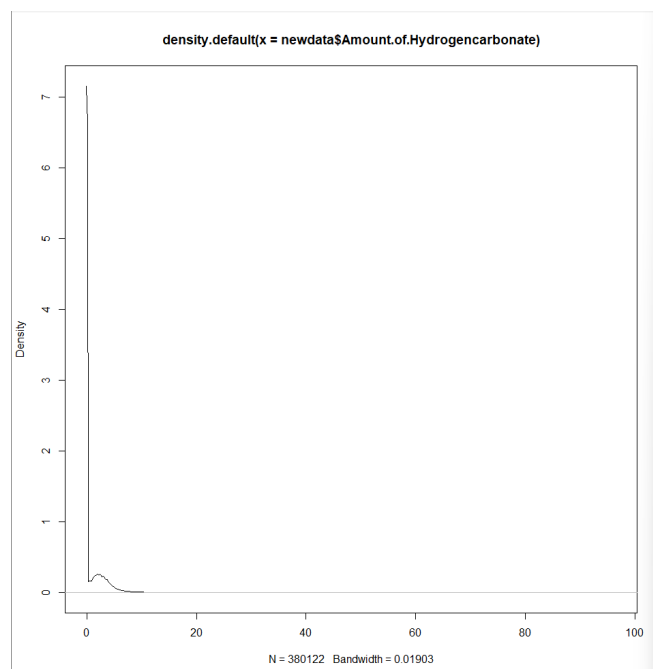
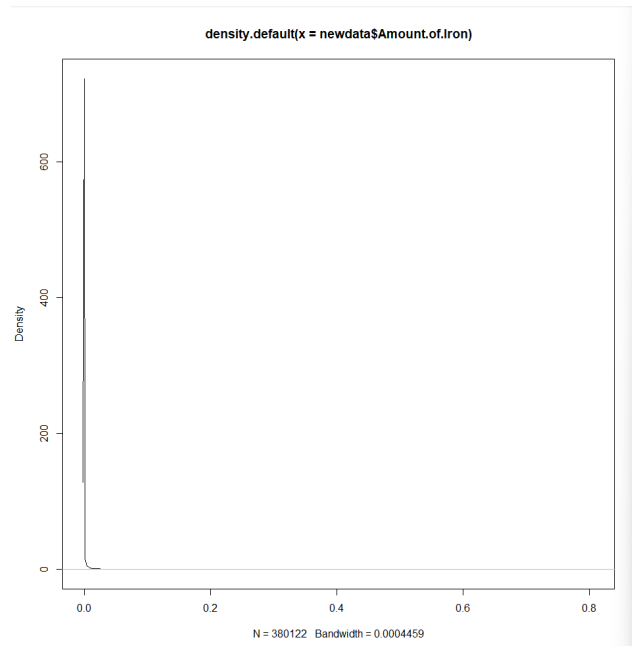


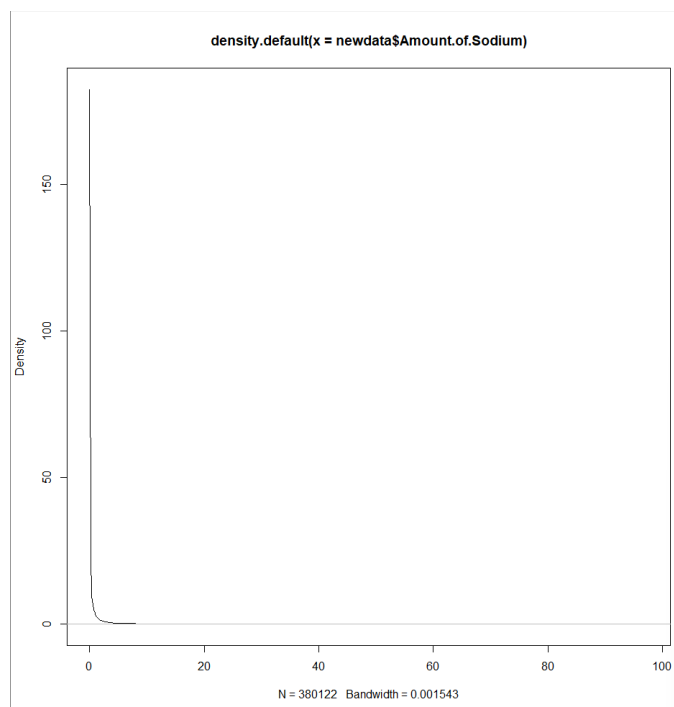
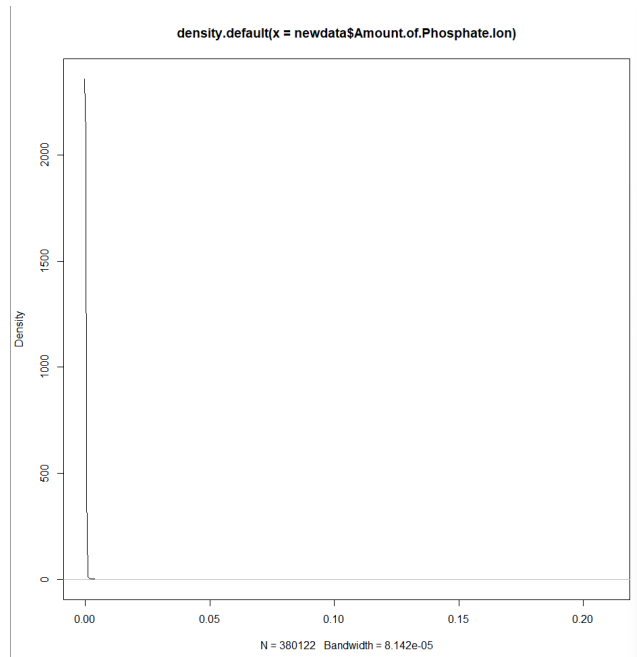
Shape of the Distribution:

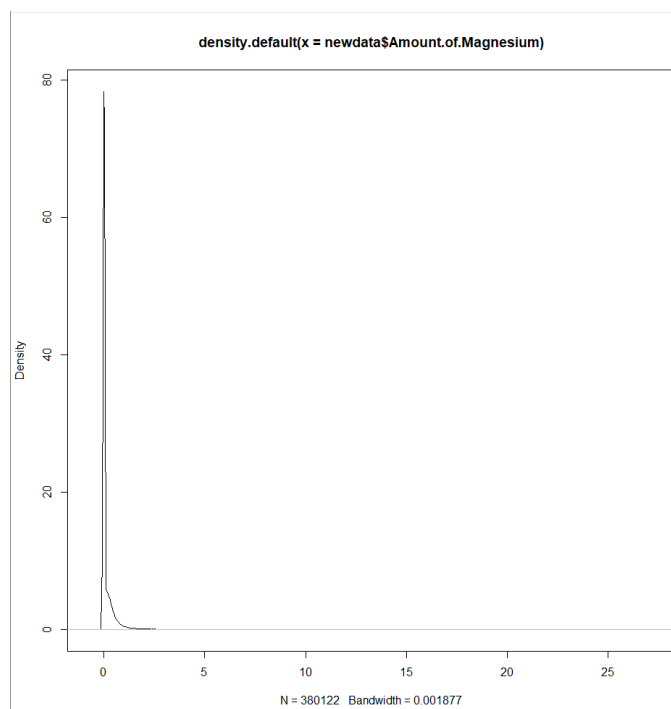
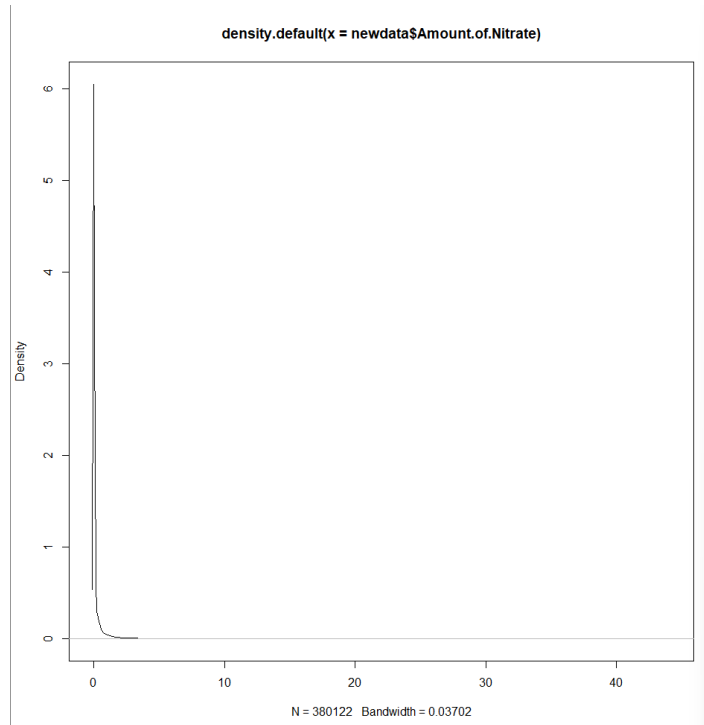


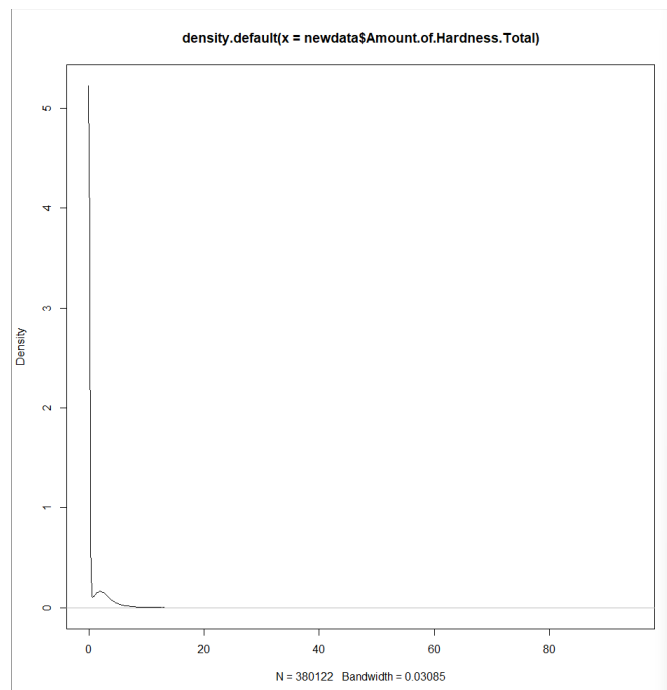
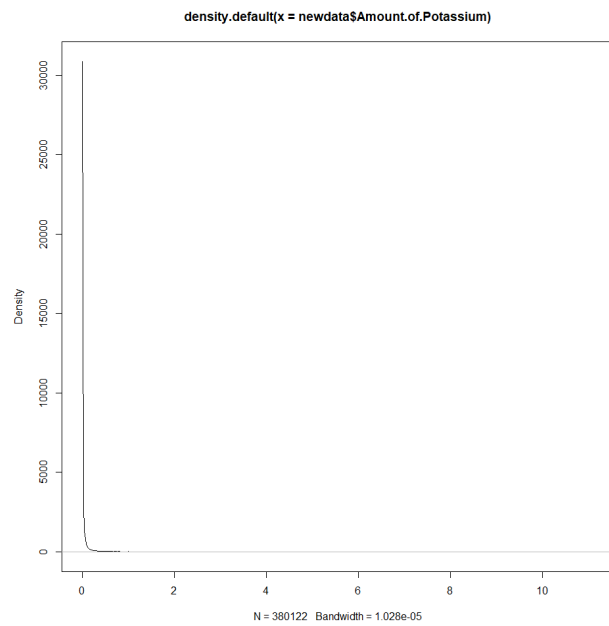


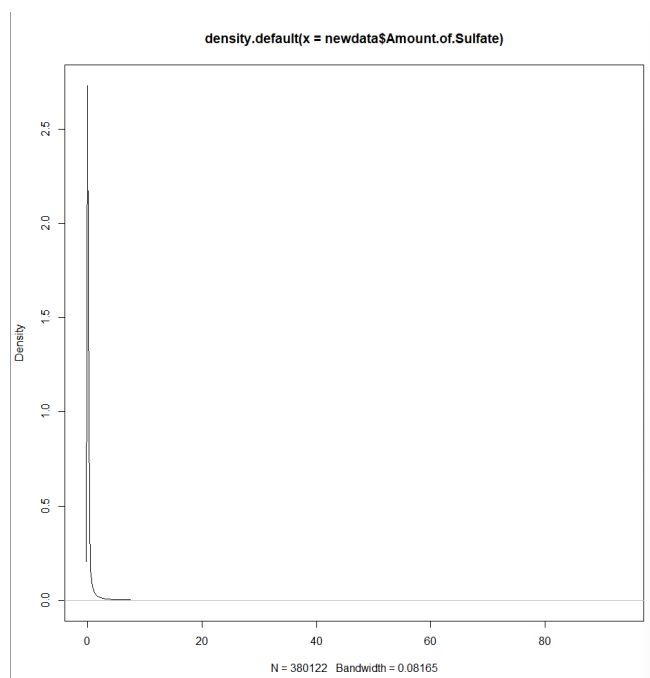
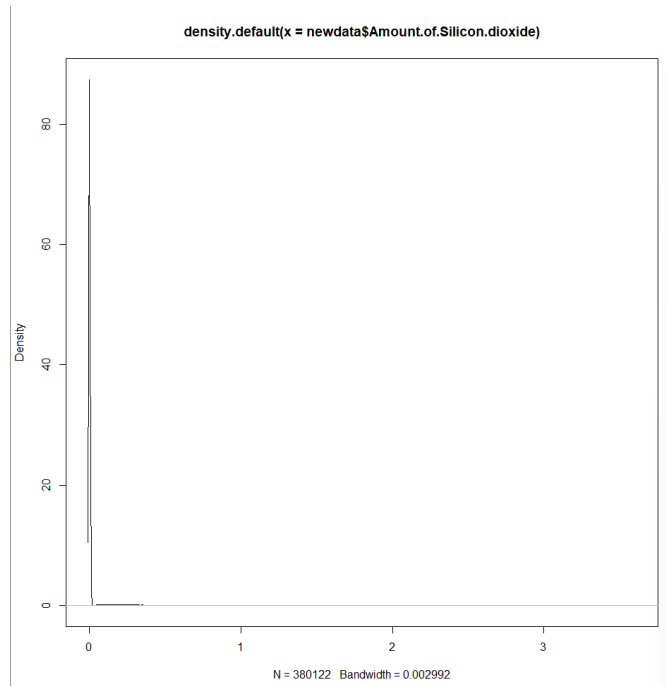


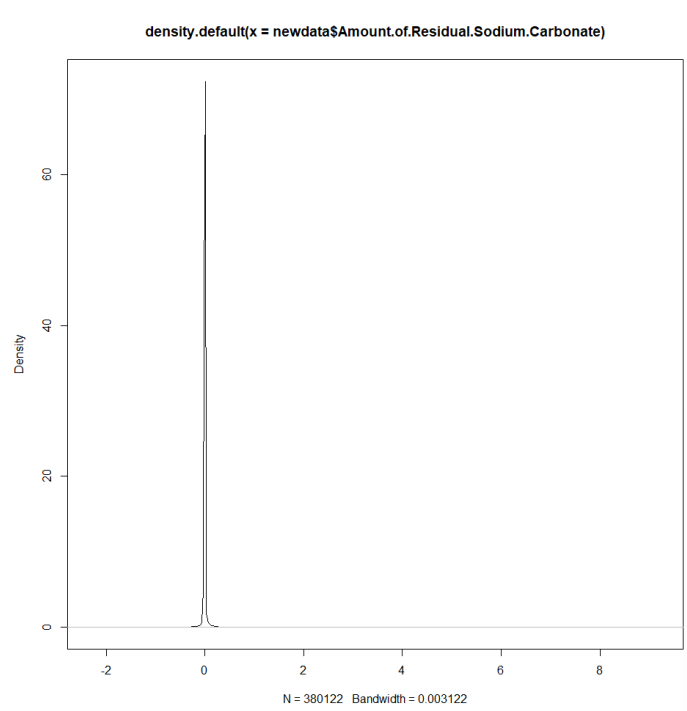
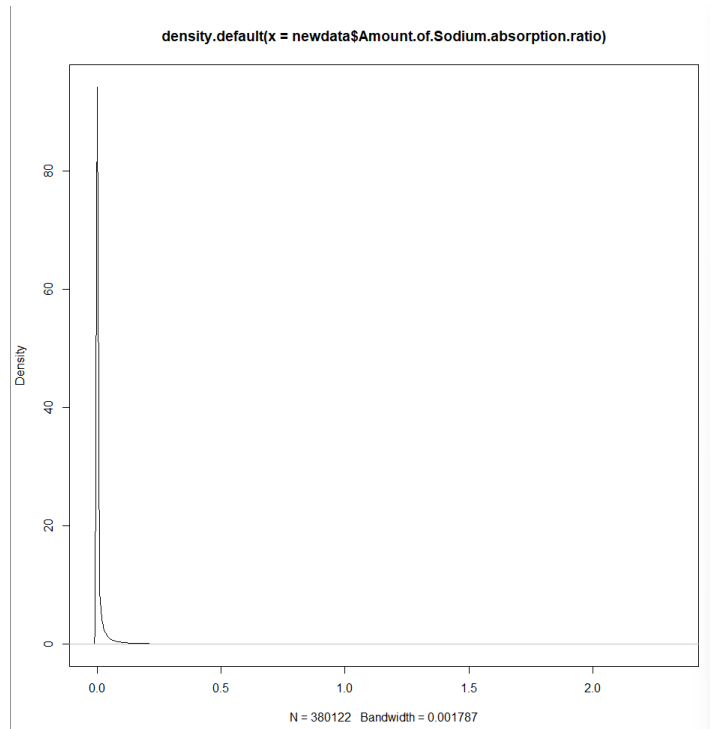


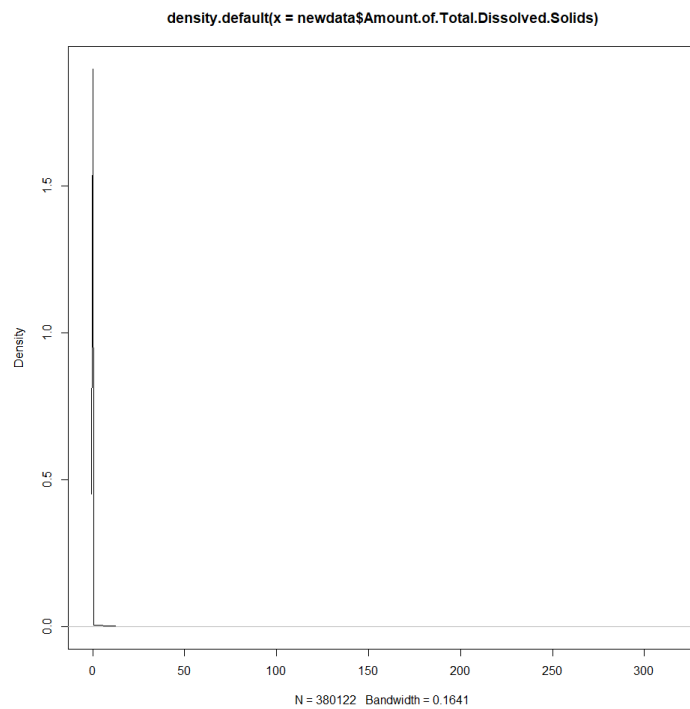
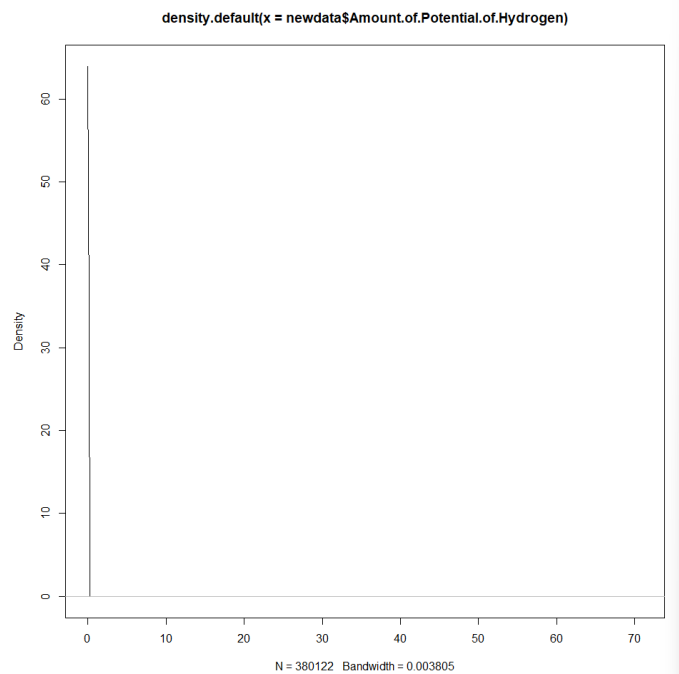


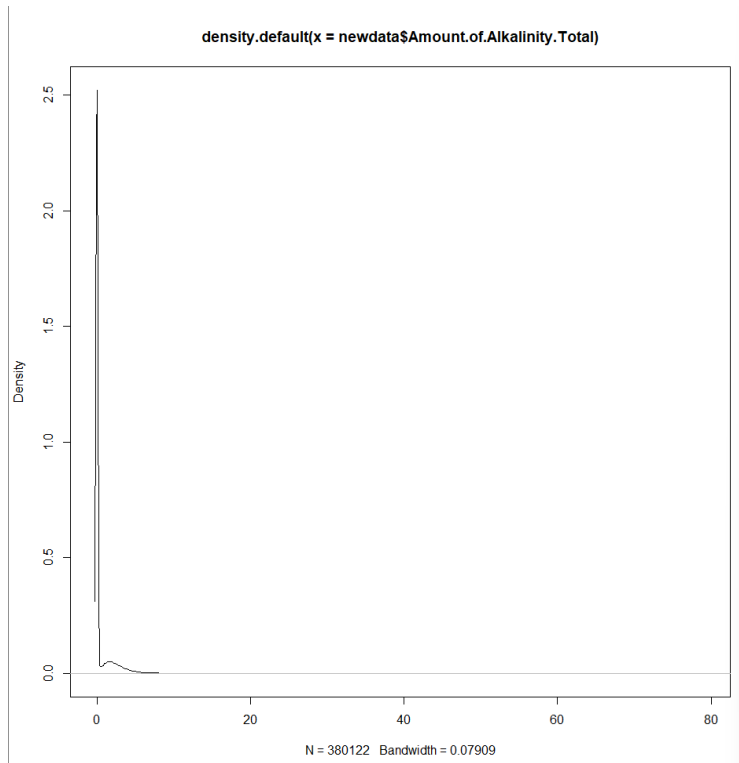












Skewness:

Amount of Arsenic = 400.6072

Amount of carbonate = 14.32542

Amount of Calcium = 9.420497

Amount of Chloride = 19.476

Amount of Electrical.Conductivity = 9.712217

Amount of Fluorine = 38.56824

Amount of Iron = 37.74283

Amount of Hydrogencarbonate = 3.451901

Amount of Potassium = 16.9747

Amount of Magnesium = 14.80591

Amount of Nitrate = 15.89685

Amount of Sodium = 14.5218

Amount of Phosphate.Ion = 67.55734

Amount of Residual.Sodium.Carbonate = 62.36075

Amount of Sodium.absorption.ratio = 10.10406

Amount of Sulfate = 16.72734

Amount of Silicon.dioxide = 14.8222

Amount of Hardness.Total = 7.539842

Amount of Alkalinity.Total = 4.798605

Amount of Total.Dissolved.Solids = 26.37784

Amount of Potential of Hydrogen = 519.0455

Yes, there are following no of outliers for all variables:

Arsenic = 110

carbonate = 11603

Calcium = 77534

Chloride = 70391

Electrical Conductivity = 48993

Fluorine = 68962

Iron = 25387

Hydrogencarbonate = 86240

Potassium = 94814

Magnesium = 84088

Nitrate = 94033

Sodium = 89833

Phosphate Ion = 6834

Residual Sodium Carbonate = 39769

Sodium absorption ratio = 81290

Sulfate = 90315

Silicon dioxide = 9565

Hardness Total = 76034

Alkalinity Total = 62239

Total Dissolved Solids = 11063

Potential of Hydrogen = 3

Mean:

Amount.of.Arsenic: 4.479559e-05

Amount.of.carbonate: 0.01227982
Amount.of.Calcium: 1.68819e-05
Amount.of.Chloride: 0.5772182
Amount.of.Electrical.Conductivity: 3.924686
Amount.of.Fluorine: 0.002037034
Amount.of.Iron: 0.0006247889
Amount.of.Hydrogencarbonate: 0.7494007
Amount.of.Potassium: 0.03841907
Amount.of.Magnesium: 0.1160776
Amount.of.Nitrate: 0.1247351
Amount.of.Sodium: 0.3664106
Amount.of.Phosphate.Ion: 4.937487e-05
Amount.of.Residual.Sodium.Carbonate: 0.0005206389
Amount.of.Sodium.absorption.ratio:0.006810049
Amount.of.Sulfate: 0.2456147
Amount.of.Silicon.dioxide: 0.005729157
Amount.of.Hardness.Total: 0.8963687
Amount.of.Alkalinity.Total: 0.3999567
Amount.of.Total.Dissolved.Solids: 0.2461375

Median:

Amount.of.Arsenic: 0
Amount.of.carbonate: 0
Amount.of.Calcium: 0
Amount.of.Chloride:0
Amount.of.Electrical.Conductivity:0
Amount.of.Fluorine:0
Amount.of.Iron:0
Amount.of.Hydrogencarbonate:0
Amount.of.Potassium:0
Amount.of.Magnesium:0
Amount.of.Nitrate:0
Amount.of.Sodium:0
Amount.of.Phosphate.Ion:0
Amount.of.Residual.Sodium.Carbonate:0
Amount.of.Sodium.absorption.ratio:0
Amount.of.Sulfate:0
Amount.of.Silicon.dioxide:0

Amount.of.Hardness.Total:0
Amount.of.Alkalinity.Total:0
Amount.of.Total.Dissolved.Solids:0

Standard deviation:

Amount.of.Arsenic: 0.006750972
Amount.of.carbonate: 0.0994766
Amount.of.Calcium:4.620545e-05
Amount.of.Chloride:2.602605
Amount.of.Electrical.Conductivity:10.78155
Amount.of.Fluorine: 0.007179138
Amount.of.Iron:0.006471257
Amount.of.Hydrogen Carbonate: 1.613891
Amount.of.Potassium:0.2512055
Amount.of.Magnesium: 0.3845186
Amount.of.Nitrate: 0.5372712
Amount.of.Sodium: 1.560392
Amount.of.Phosphate.Ion: 0.001181585
Amount.of.Residual.Sodium.Carbonate:0.04531287
Amount.of.Sodium.absorption.ratio: 0.02593723
Amount.of.Sulfate:1.185014
Amount.of.Silicon.dioxide:0.04342095
Amount.of.Hardness.Total: 2.368146
Amount.of.Alkalinity.Total: 1.147759
Amount.of.Total.Dissolved.Solids: 2.381736

6)

The lowest and maximum residuals are -197 and 19804.2, respectively. Half of the residuals are smaller than -125.8 and half are bigger than -125.8, which is the median residual. -168.7 is the first quartile (or the 25th percentile), and -10 is the third quartile (or the 75th percentile).

The anticipated value of the environmental quality indicator when the SDP is zero is represented by the model's intercept, which is $1.981e+02$. The environmental quality indicator is anticipated to fall by $-5.314e-06$ units for every unit increase in SDP.

Residual Standard error is 449.4 which tells us that the predicted values differ by 449.4 on an average.

Residuals:

Min	1Q	Median	3Q	Max
-197.0	-168.7	-125.8	-10.0	19804.2

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	1.981e+02	2.069e+00	95.774	<2e-16 ***
sdp	-5.314e-06	3.903e-06	-1.361	0.173

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

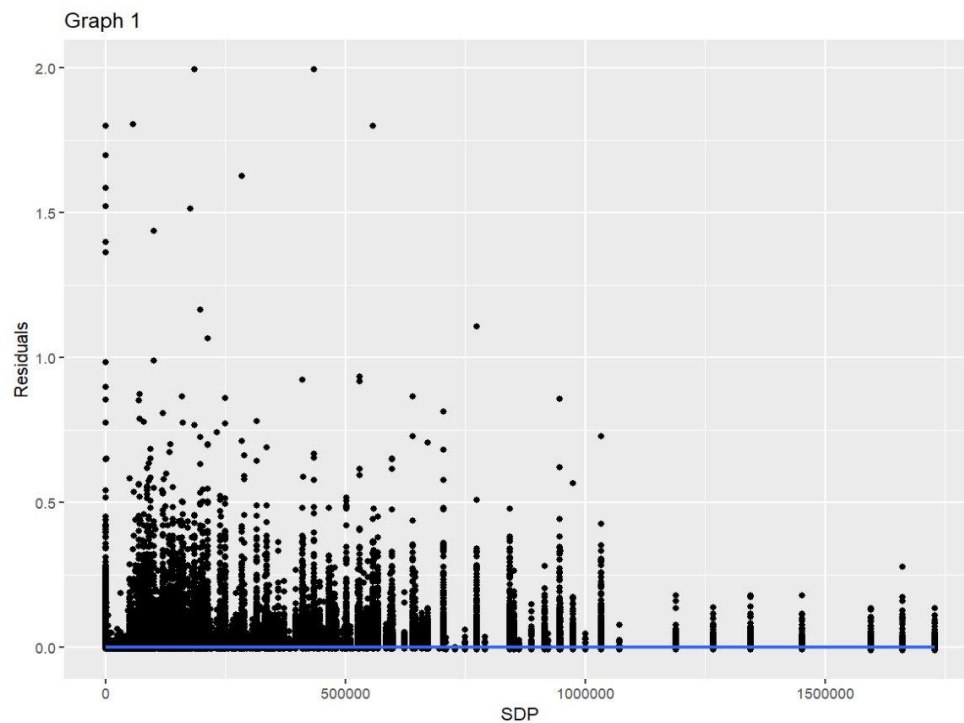
Residual standard error: 449.4 on 95334 degrees of freedom

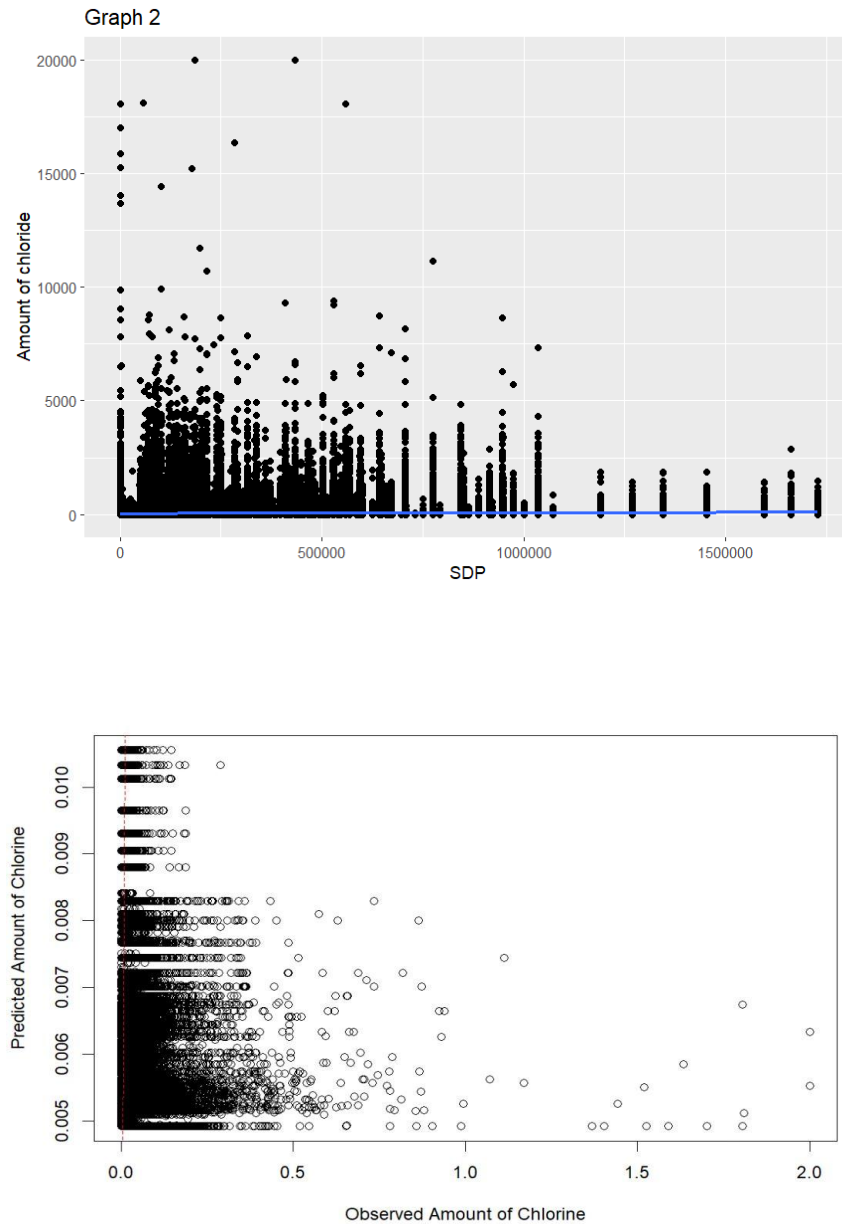
(284786 observations deleted due to missingness)

Multiple R-squared: 1.944e-05, Adjusted R-squared: 8.95e-06

F-statistic: 1.853 on 1 and 95334 DF, p-value: 0.1734

7)





The graph show how the SDP and model residuals—the discrepancies between the environmental quality indicator's observed and anticipated values—relate to one another. The SDP values are represented on the x-axis, and the model residuals are represented on the y-axis. Each point on the plot reflects the discrepancy between the environmental quality indicator's actual and expected values for a specific SDP value.

We would anticipate no discernible pattern or trend in the residuals as the SDP values rise if the model is a good fit for the data. However, if there is a pattern or trend in the residuals, it may be a sign that not all of the significant elements that affect the environmental quality

indicator are being well captured by the model. Moreover, it can imply the presence of heteroscedasticity (i.e., the variance of the residuals is not constant across all levels of SDP).

First plot shows residuals on y-axis and sdp on x axis. This shows the residual values corresponding to the various sdp values.

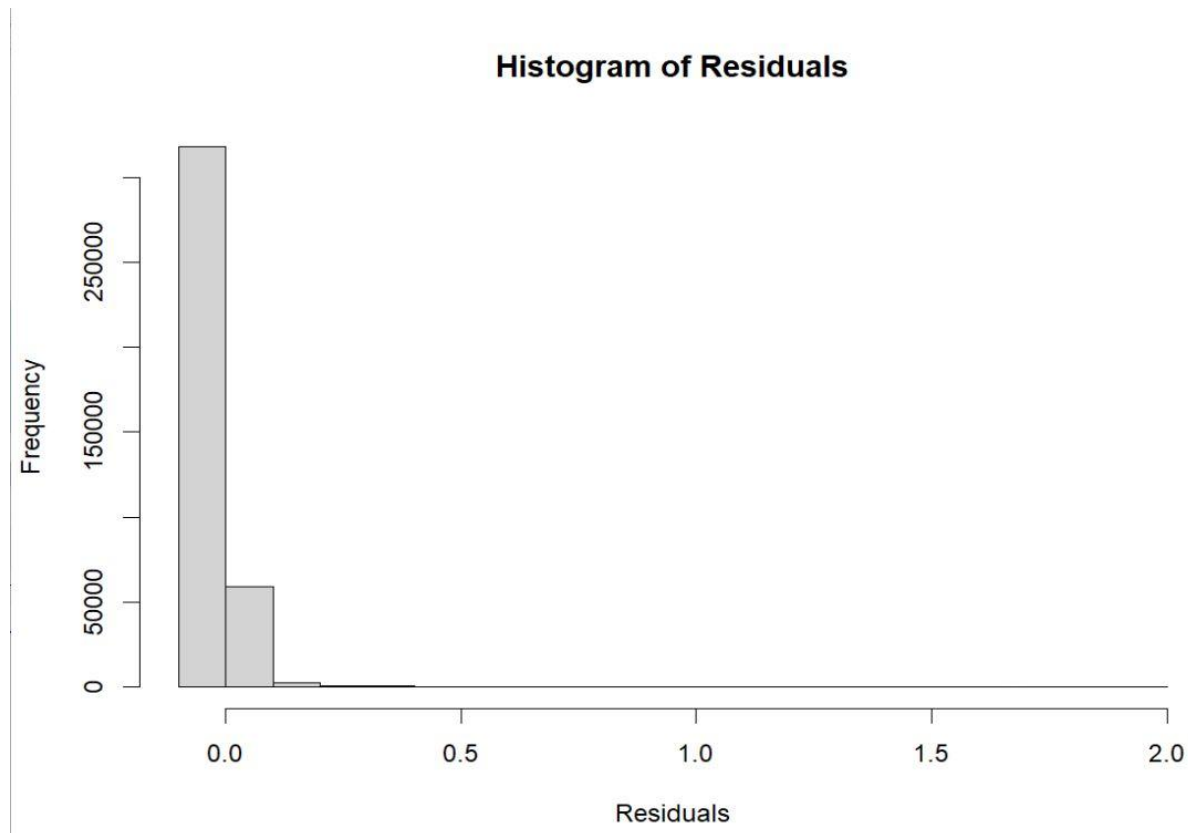
Second plot shows the EQI(Amount of chloride on y-axis) and SDP on x axis and this plots the residuals corresponding to both.

The third plot shows the predicted and the observed values.

The link between all three of these models indicates that the projected model is highly similar to the real model, which is a measure of its efficacy. We may infer that all three models are the quantifiers of the same approximation and that any one of them will yield results that are roughly equivalent because they all point to the same analysis.

Together, they enable us to assess the model's overall performance as well as the areas where it may be strengthened.

8)



The summation of residuals comes out to be **2.294001e-06** which is very close to zero thus depicting that our model has captured most of the variations.

9)

The model includes five predictor variables: a constant (intercept), a variable called "newdata\$sdp", the square of "newdata\$sdp" (represented as "I(sdp^2)"), the cube of "newdata\$sdp" (represented as "I(sdp^3)"), and a variable called "Gini".

The intercept is 3.080e+02, indicating that when all predictor variables are 0, the predicted level is 3.080e+02.

The coefficient estimates indicate how much the response variable is expected to change for a one-unit increase in each predictor variable, holding all other predictors constant. For example, for a one-unit increase in "newdata\$sdp", we expect the response variable to increase by 2.484e-04.

The p-values associated with each coefficient indicate the level of significance of each predictor variable. A low p-value (less than 0.05) indicates that the predictor variable is statistically significant and likely to have a non-zero effect on the response variable.

The multiple R-squared values (0.009021) indicate that the model explains only a small portion of the variability in the response variable.

The adjusted R-squared value takes into account the number of predictor variables in the model and is slightly lower than the multiple R-squared values.

The residual standard error is a measure of the variability of the response variable that the model does not explain. In this case, the residual standard error is 448.3.

The F-statistic and its associated p-value test whether the model as a whole is statistically significant. In this case, the F-statistic is 214.9, and the p-value is less than 2.2e-16, indicating that the model is statistically significant.

Residuals:

Min	1Q	Median	3Q	Max
-278.3	-163.6	-112.6	-7.9	19788.4

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)	
(Intercept)	3.080e+02	8.210e+00	37.516	< 2e-16	***
newdata\$sdp	2.484e-04	2.786e-05	8.915	< 2e-16	***
I(sdp^2)	-2.296e-10	4.577e-11	-5.015	5.3e-07	***
I(sdp^3)	3.580e-17	1.922e-17	1.863	0.0625	.
Gini	-5.391e+02	2.497e+01	-21.590	< 2e-16	***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 448.3 on 94444 degrees of freedom
(285673 observations deleted due to missingness)

Multiple R-squared: 0.009021, Adjusted R-squared: 0.008979

F-statistic: 214.9 on 4 and 94444 DF, p-value: < 2.2e-16