

## Data Collection and Preprocessing Phase

Date	07 JULY 2024
Team ID	739850
Project Title	Air Quality Index Analyzer using machine learning
Maximum Marks	6 Marks

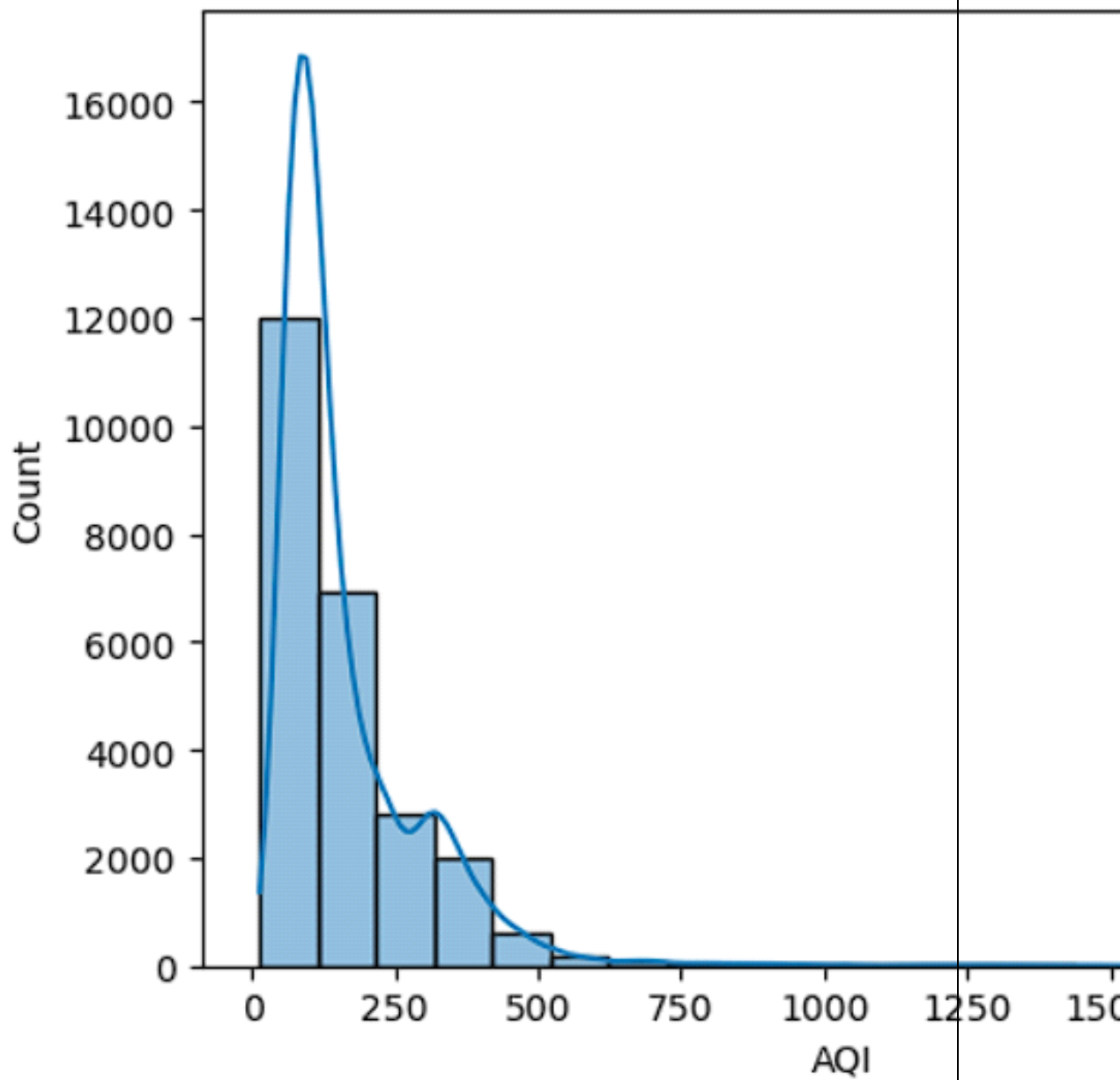
# Data Exploration and Preprocessing Report

Identifies data sources, assesses quality issues like missing values and duplicates, and implements resolution plans to ensure accurate and reliable analysis

[illegible]

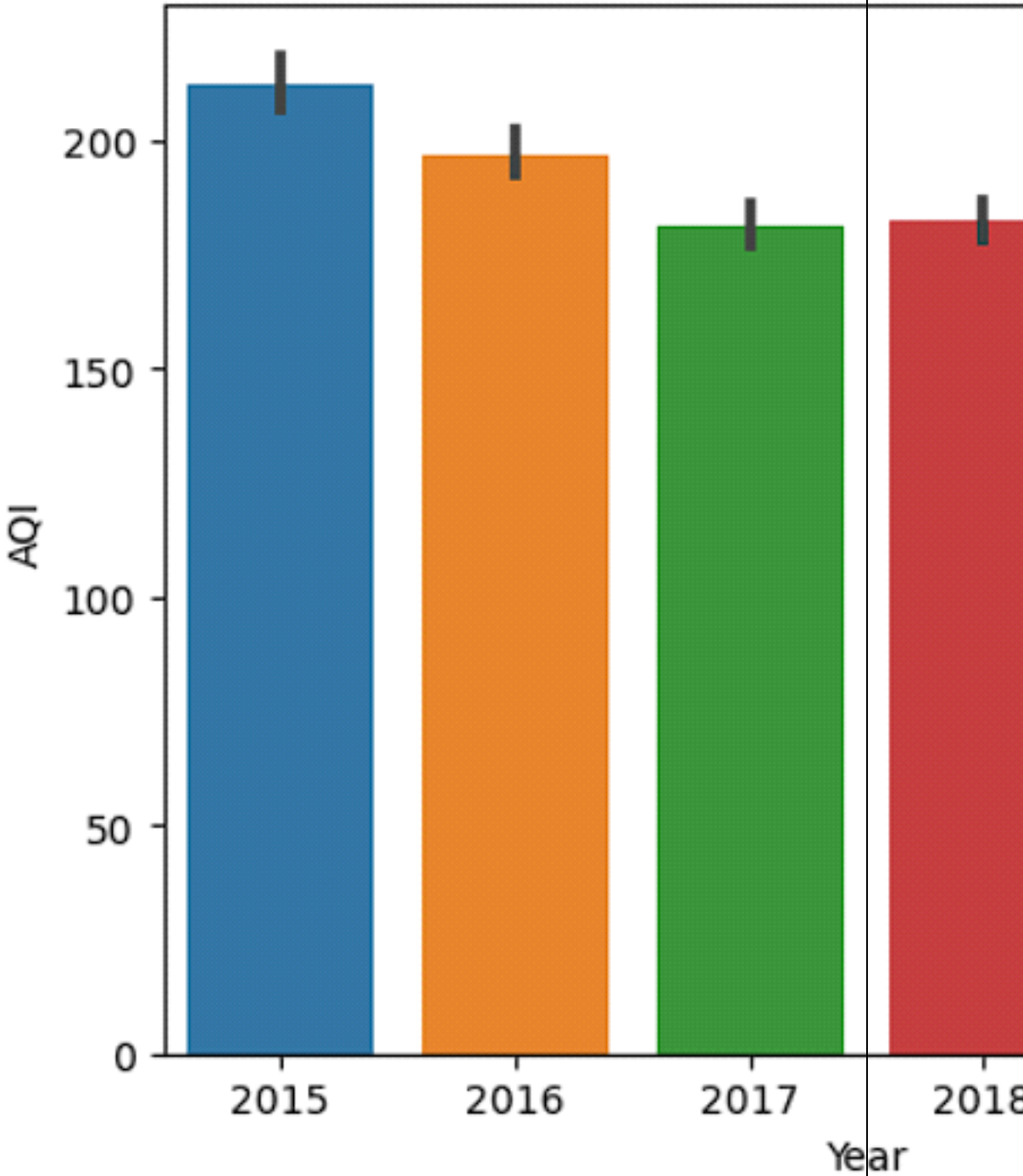
Uni  
vari  
ate  
Ana  
lysi  
s

Exploration of individual variables (mean, median, mode, etc.).

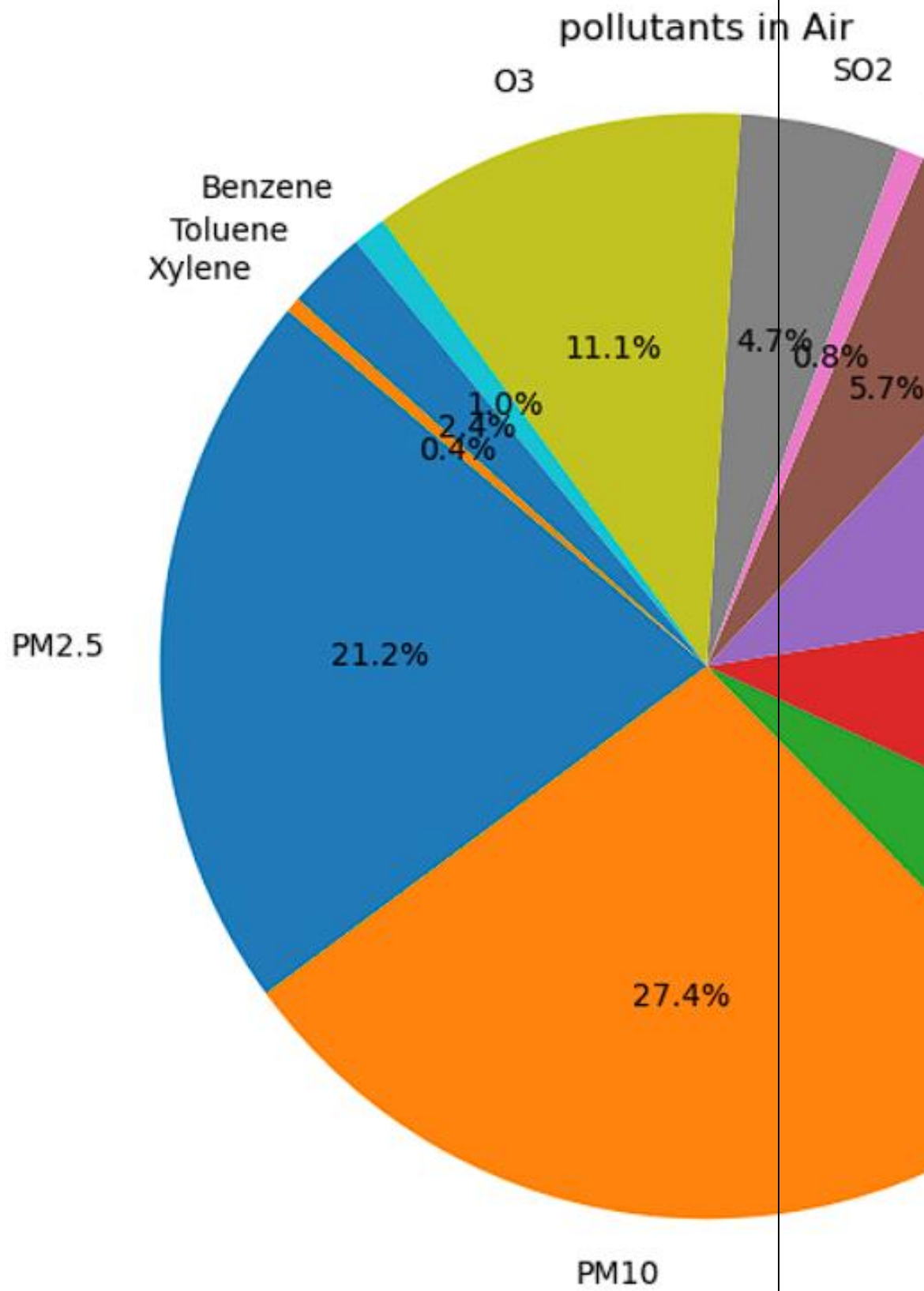


Biv  
aria  
te  
Ana  
lysi  
s

Relationships between two variables (correlation, scatter plots).

	 <table border="1"><thead><tr><th>Year</th><th>AQI</th></tr></thead><tbody><tr><td>2015</td><td>215</td></tr><tr><td>2016</td><td>198</td></tr><tr><td>2017</td><td>182</td></tr><tr><td>2018</td><td>184</td></tr></tbody></table>	Year	AQI	2015	215	2016	198	2017	182	2018	184
Year	AQI										
2015	215										
2016	198										
2017	182										
2018	184										
Mul tiva riat e Ana lysi	Patterns and relationships involving multiple variables.										

S



Outliers and Anomalies	<p>Identification and treatment of outliers.</p> <pre> 26]: import pandas as pd import numpy as np import matplotlib.pyplot as plt  def handle_outliers(df):     # Plot boxplots before handling outliers     plt.figure(figsize=(15, 10))     df.boxplot(rot=90)     plt.title('Boxplot Before Handling Outliers')     plt.show()      for column in df.columns:         if pd.api.types.is_numeric_dtype(df[column]):             Q1 = df[column].quantile(0.25)             Q3 = df[column].quantile(0.75)             IQR = Q3 - Q1             lower_bound = Q1 - 1.5 * IQR             upper_bound = Q3 + 1.5 * IQR              # Cap the outliers             df[column] = np.where(df[column] &lt; lower_bound,                                   lower_bound,                                   np.where(df[column] &gt; upper_bound,   upper_bound,   df[column]))      # Plot boxplots after handling outliers     plt.figure(figsize=(15, 10))     df.boxplot(rot=90)     plt.title('Boxplot After Handling Outliers')     plt.show()      return df </pre>
Data Preprocessing Code Screenshots	
Loading Data	<p>Code to load the dataset into the preferred environment (e.g., Python, R).</p>

```
In [48]: X=data_city.drop('AQI',axis=1)
y=data_city['AQI']
```

```
In [49]: X
```

```
Out[49]:
```

	City	PM2.5	PM10	NO	NO2	NOx
0	0	34.515	154.750	0.92	18.22	17.15
1	0	25.830	226.235	0.97	15.69	16.46
2	0	36.205	72.125	17.40	19.30	29.70
3	0	25.830	226.235	1.70	18.48	17.97
4	0	54.440	72.125	22.10	21.42	37.76
...	...	...	...	...	...	...
29526	25	15.020	50.940	7.68	25.06	19.54
29527	25	24.380	74.090	3.42	26.06	16.53
29528	25	22.910	65.730	3.45	29.53	18.33
29529	25	16.640	49.970	4.05	29.26	18.80
29530	25	15.000	66.000	0.40	26.85	14.05

29531 rows × 15 columns



Smart  
Internz



Smart  
Internz

Handling Missing Data	<p>Code for identifying and handling missing values.</p> <h2>Handling Null Values</h2> <pre>2]: data_city.isna().sum()</pre> <pre>2]: City                0       Date                0       PM2.5             4598       PM10             11140       NO                3582       NO2              3585       NOx              4185       NH3             10328       CO               2059       SO2              3854       O3              4022       Benzene          5623       Toluene          8041       Xylene          18109       AQI              4681       AQI_Bucket       4681       dtype: int64</pre>
Data	Code for transforming variables (scaling, normalization).

Transformation																																																																																																																																																													
Feature Engineering	<div>Code for creating new features or modifying existing ones.</div> <div><div>In [48]:</div><div>x=data_city.drop('AQI',axis=1) y=data_city['AQI']</div></div> <div><div>In [49]:</div><div>X</div></div> <div><div>Out[49]:</div><table><tr><th></th><th>City</th><th>PM2.5</th><th>PM10</th><th>NO</th><th>NO2</th><th>NOx</th><th>NH3</th><th>CO</th><th>SO2</th><th>O3</th><th>Benzene</th><th>To</th></tr><tr><td>0</td><td>0</td><td>34.515</td><td>154.750</td><td>0.92</td><td>18.22</td><td>17.15</td><td>8.975</td><td>0.92</td><td>27.640</td><td>85.635</td><td></td><td>0.00</td></tr><tr><td>1</td><td>0</td><td>25.830</td><td>226.235</td><td>0.97</td><td>15.69</td><td>16.46</td><td>9.095</td><td>0.97</td><td>24.550</td><td>34.060</td><td></td><td>3.68</td></tr><tr><td>2</td><td>0</td><td>36.205</td><td>72.125</td><td>17.40</td><td>19.30</td><td>29.70</td><td>6.880</td><td>2.86</td><td>29.070</td><td>30.700</td><td></td><td>6.80</td></tr><tr><td>3</td><td>0</td><td>25.830</td><td>226.235</td><td>1.70</td><td>18.48</td><td>17.97</td><td>9.085</td><td>1.70</td><td>18.590</td><td>36.080</td><td></td><td>4.43</td></tr><tr><td>4</td><td>0</td><td>54.440</td><td>72.125</td><td>22.10</td><td>21.42</td><td>37.76</td><td>7.915</td><td>2.86</td><td>29.545</td><td>39.310</td><td></td><td>7.01</td></tr><tr><td>...</td><td>...</td><td>...</td><td>...</td><td>...</td><td>...</td><td>...</td><td>...</td><td>...</td><td>...</td><td>...</td><td></td><td>...</td></tr><tr><td>29526</td><td>25</td><td>15.020</td><td>50.940</td><td>7.68</td><td>25.06</td><td>19.54</td><td>12.470</td><td>0.47</td><td>8.550</td><td>23.300</td><td></td><td>2.24</td></tr><tr><td>29527</td><td>25</td><td>24.380</td><td>74.090</td><td>3.42</td><td>26.06</td><td>16.53</td><td>11.990</td><td>0.52</td><td>12.720</td><td>30.140</td><td></td><td>0.74</td></tr><tr><td>29528</td><td>25</td><td>22.910</td><td>65.730</td><td>3.45</td><td>29.53</td><td>18.33</td><td>10.710</td><td>0.48</td><td>8.420</td><td>30.960</td><td></td><td>0.01</td></tr><tr><td>29529</td><td>25</td><td>16.640</td><td>49.970</td><td>4.05</td><td>29.26</td><td>18.80</td><td>10.030</td><td>0.52</td><td>9.840</td><td>28.300</td><td></td><td>0.00</td></tr><tr><td>29530</td><td>25</td><td>15.000</td><td>66.000</td><td>0.40</td><td>26.85</td><td>14.05</td><td>5.200</td><td>0.59</td><td>2.100</td><td>17.050</td><td></td><td>0.00</td></tr></table><div>29531 rows × 15 columns</div></div>		City	PM2.5	PM10	NO	NO2	NOx	NH3	CO	SO2	O3	Benzene	To	0	0	34.515	154.750	0.92	18.22	17.15	8.975	0.92	27.640	85.635		0.00	1	0	25.830	226.235	0.97	15.69	16.46	9.095	0.97	24.550	34.060		3.68	2	0	36.205	72.125	17.40	19.30	29.70	6.880	2.86	29.070	30.700		6.80	3	0	25.830	226.235	1.70	18.48	17.97	9.085	1.70	18.590	36.080		4.43	4	0	54.440	72.125	22.10	21.42	37.76	7.915	2.86	29.545	39.310		7.01	...	...	...	...	...	...	...	...	...	...	...		...	29526	25	15.020	50.940	7.68	25.06	19.54	12.470	0.47	8.550	23.300		2.24	29527	25	24.380	74.090	3.42	26.06	16.53	11.990	0.52	12.720	30.140		0.74	29528	25	22.910	65.730	3.45	29.53	18.33	10.710	0.48	8.420	30.960		0.01	29529	25	16.640	49.970	4.05	29.26	18.80	10.030	0.52	9.840	28.300		0.00	29530	25	15.000	66.000	0.40	26.85	14.05	5.200	0.59	2.100	17.050		0.00
	City	PM2.5	PM10	NO	NO2	NOx	NH3	CO	SO2	O3	Benzene	To																																																																																																																																																	
0	0	34.515	154.750	0.92	18.22	17.15	8.975	0.92	27.640	85.635		0.00																																																																																																																																																	
1	0	25.830	226.235	0.97	15.69	16.46	9.095	0.97	24.550	34.060		3.68																																																																																																																																																	
2	0	36.205	72.125	17.40	19.30	29.70	6.880	2.86	29.070	30.700		6.80																																																																																																																																																	
3	0	25.830	226.235	1.70	18.48	17.97	9.085	1.70	18.590	36.080		4.43																																																																																																																																																	
4	0	54.440	72.125	22.10	21.42	37.76	7.915	2.86	29.545	39.310		7.01																																																																																																																																																	
...	...	...	...	...	...	...	...	...	...	...		...																																																																																																																																																	
29526	25	15.020	50.940	7.68	25.06	19.54	12.470	0.47	8.550	23.300		2.24																																																																																																																																																	
29527	25	24.380	74.090	3.42	26.06	16.53	11.990	0.52	12.720	30.140		0.74																																																																																																																																																	
29528	25	22.910	65.730	3.45	29.53	18.33	10.710	0.48	8.420	30.960		0.01																																																																																																																																																	
29529	25	16.640	49.970	4.05	29.26	18.80	10.030	0.52	9.840	28.300		0.00																																																																																																																																																	
29530	25	15.000	66.000	0.40	26.85	14.05	5.200	0.59	2.100	17.050		0.00																																																																																																																																																	