

Data Collection and Preprocessing Phase

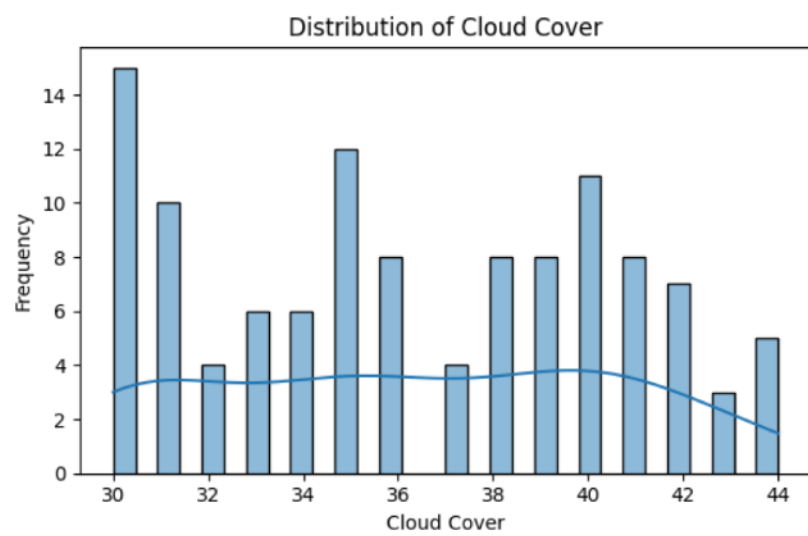
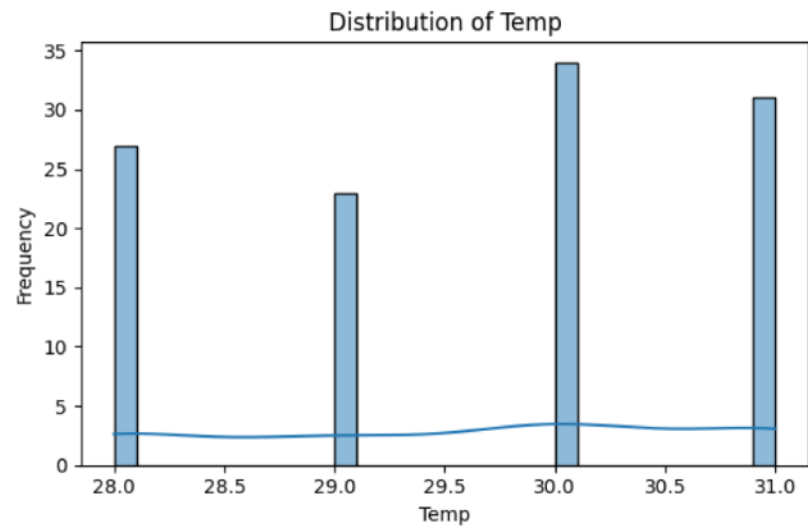
Date	13 June 2025
Team ID	SWTID1749618778
Project Title	Rising Waters: A Machine Learning Approach To Flood Prediction
Maximum Marks	6 Marks

Data Exploration and Preprocessing Template

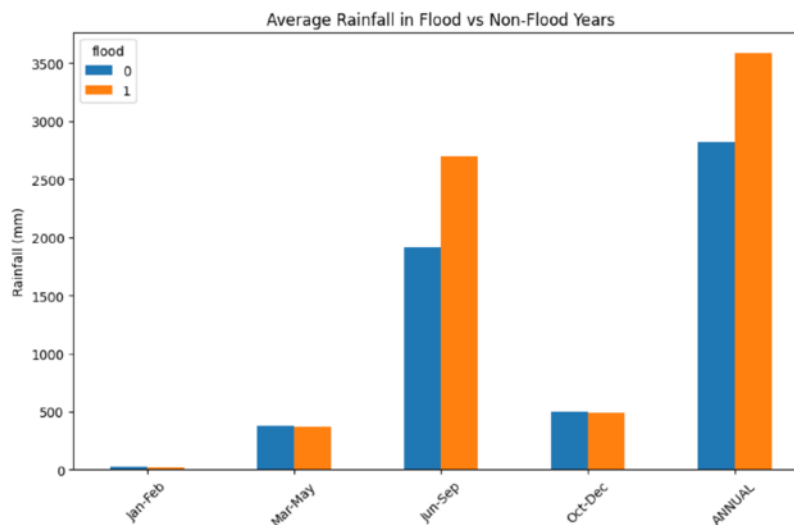
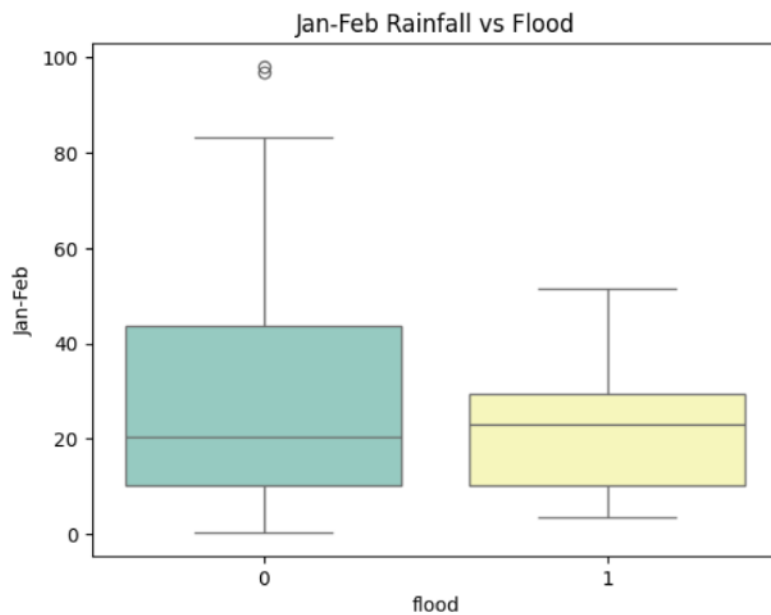
Identifies data sources, assesses quality issues like missing values and duplicates, and implements resolution plans to ensure accurate and reliable analysis.

Section	Description
Data Overview	<u>Dimension:</u> 115 rows x 11 columns
	<u>Descriptive statistics:</u>

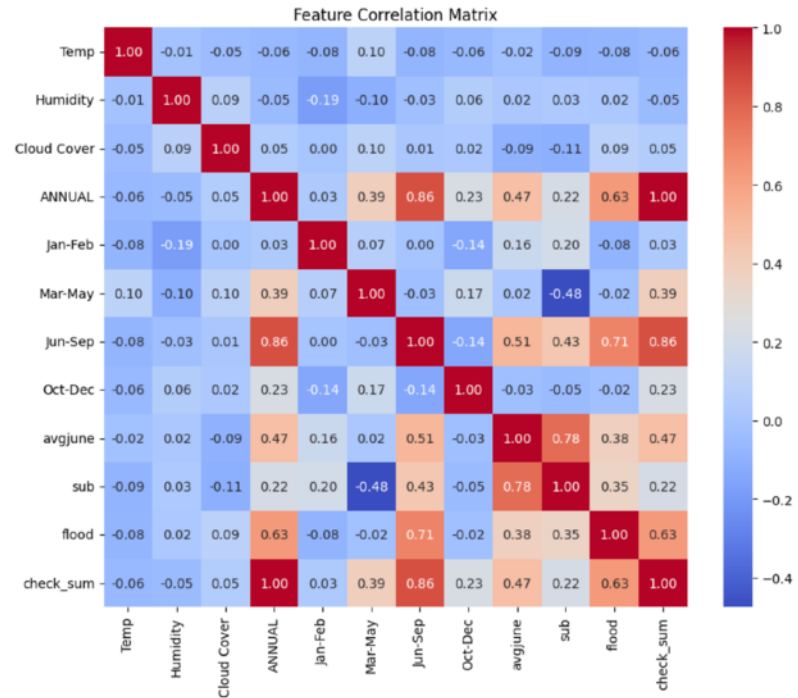
Univariate Analysis



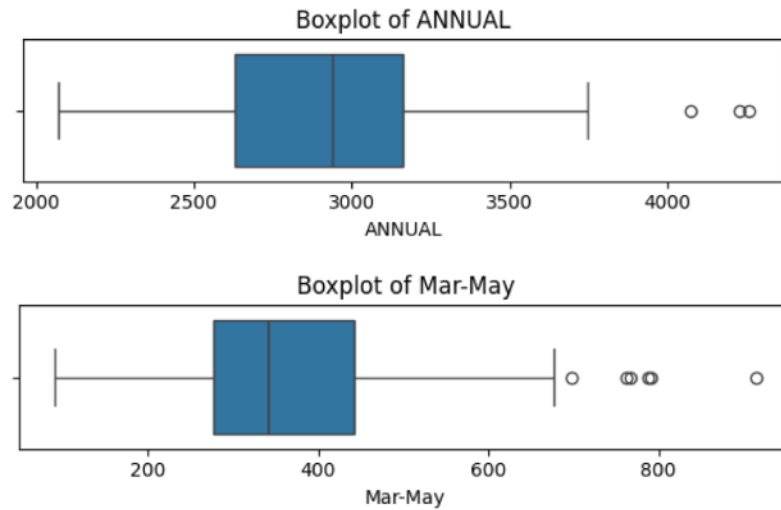
Bivariate Analysis



Multivariate Analysis



Outliers and Anomalies



Data Preprocessing Code Screenshots

Loading Data	<pre>df = pd.read_excel('../data/flood_dataset_raw.xlsx')</pre> <pre>df.head()</pre> <table><thead><tr><th></th><th>Temp</th><th>Humidity</th><th>Cloud Cover</th><th>ANNUAL</th><th>Jan-Feb</th><th>Mar-May</th><th>Jun-Sep</th><th>Oct-Dec</th><th>avgjune</th><th>sub</th><th>flood</th></tr></thead><tbody><tr><td>0</td><td>29</td><td>70</td><td>30</td><td>3248.6</td><td>73.4</td><td>386.2</td><td>2122.8</td><td>666.1</td><td>274.866667</td><td>649.9</td><td>0</td></tr><tr><td>1</td><td>28</td><td>75</td><td>40</td><td>3326.6</td><td>9.3</td><td>275.7</td><td>2403.4</td><td>638.2</td><td>130.300000</td><td>256.4</td><td>1</td></tr><tr><td>2</td><td>28</td><td>75</td><td>42</td><td>3271.2</td><td>21.7</td><td>336.3</td><td>2343.0</td><td>570.1</td><td>186.200000</td><td>308.9</td><td>0</td></tr><tr><td>3</td><td>29</td><td>71</td><td>44</td><td>3129.7</td><td>26.7</td><td>339.4</td><td>2398.2</td><td>365.3</td><td>366.066667</td><td>862.5</td><td>0</td></tr><tr><td>4</td><td>31</td><td>74</td><td>40</td><td>2741.6</td><td>23.4</td><td>378.5</td><td>1881.5</td><td>458.1</td><td>283.400000</td><td>586.9</td><td>0</td></tr></tbody></table>		Temp	Humidity	Cloud Cover	ANNUAL	Jan-Feb	Mar-May	Jun-Sep	Oct-Dec	avgjune	sub	flood	0	29	70	30	3248.6	73.4	386.2	2122.8	666.1	274.866667	649.9	0	1	28	75	40	3326.6	9.3	275.7	2403.4	638.2	130.300000	256.4	1	2	28	75	42	3271.2	21.7	336.3	2343.0	570.1	186.200000	308.9	0	3	29	71	44	3129.7	26.7	339.4	2398.2	365.3	366.066667	862.5	0	4	31	74	40	2741.6	23.4	378.5	1881.5	458.1	283.400000	586.9	0
	Temp	Humidity	Cloud Cover	ANNUAL	Jan-Feb	Mar-May	Jun-Sep	Oct-Dec	avgjune	sub	flood																																																														
0	29	70	30	3248.6	73.4	386.2	2122.8	666.1	274.866667	649.9	0																																																														
1	28	75	40	3326.6	9.3	275.7	2403.4	638.2	130.300000	256.4	1																																																														
2	28	75	42	3271.2	21.7	336.3	2343.0	570.1	186.200000	308.9	0																																																														
3	29	71	44	3129.7	26.7	339.4	2398.2	365.3	366.066667	862.5	0																																																														
4	31	74	40	2741.6	23.4	378.5	1881.5	458.1	283.400000	586.9	0																																																														
Handling Missing Data	<pre>df = pd.read_excel('../data/flood_dataset_raw.xlsx')</pre> <pre>df.isnull().any()</pre> <pre>Temp False Humidity False Cloud Cover False ANNUAL False Jan-Feb False Mar-May False Jun-Sep False Oct-Dec False avgjune False sub False flood False dtype: bool</pre>																																																																								
Data Transformation	<pre>X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)</pre> <pre>scaler = StandardScaler()</pre> <pre>X_train_scaled = scaler.fit_transform(X_train)</pre> <pre>X_test_scaled = scaler.transform(X_test)</pre>																																																																								
Feature Engineering	Attached the codes in final submission.																																																																								
Save Processed Data	-																																																																								