# Data Exploration and Preprocessing Report
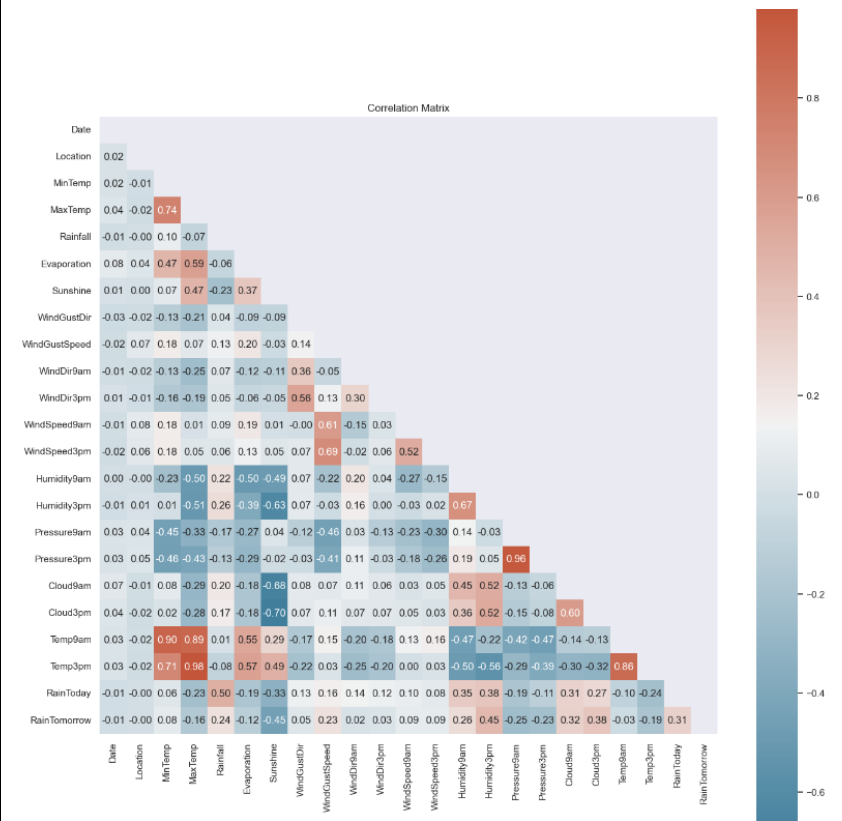
| Date | 15 April 2024 |
| --- | --- |
| Team ID | Team-738164 |
| Project Title | Rainfall Prediction Using Machine Learning |
| Maximum Marks | 6 Marks |

**Data Exploration Screenshots:**

| Section | Description |
| --- | --- |
| Data Overview | <u>Dimensions:</u><br>145460 rows x 23 columns<br><u>Descriptive Statistics:</u><br> |
| Univariate Analysis |  |
| Bivariate Analysis | |

## Rainy Days by Location



## Seasonality of Rainfall

| | |
|---|---|
| Multivariate Analysis |  Correlation Matrix |
| Outliers and Anomalies | 1. Multiple columns have clear outliers (e.g., the max Rainfall value is 371.0 despite the 75th percentile being 0.8) <br> 2. Not seeing any values that are immediate cause for concern (such as a negative value for minimum Rainfall) |

**Data Preprocessing Code Screenshots:**

| | |
|---|---|
| Loading Data | ```# Loading the dataset
df = pd.read_csv('weatherAUS.csv')

df.head()``` <br><br>  |

Table (df.head() output):

| | Date | Location | MinTemp | MaxTemp | Rainfall | Evaporation | Sunshine | WindGustDir | WindGustSpeed | WindDir9am | WindDir3p |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 2008-12-01 | Albury | 13.4 | 22.9 | 0.6 | NaN | NaN | W | 44.0 | W | WN |
| 1 | 2008-12-02 | Albury | 7.4 | 25.1 | 0.0 | NaN | NaN | WNW | 44.0 | NNW | WS |
| 2 | 2008-12-03 | Albury | 12.9 | 25.7 | 0.0 | NaN | NaN | WSW | 46.0 | W | WS |
| 3 | 2008-12-04 | Albury | 9.2 | 28.0 | 0.0 | NaN | NaN | NE | 24.0 | SE | |
| 4 | 2008-12-05 | Albury | 17.5 | 32.3 | 1.0 | NaN | NaN | W | 41.0 | ENE | N |

| | |
|---|---|
| Handling Missing Data | ```python\ndf_imputed = df.dropna(axis=0, subset=['RainTomorrow'])\n\n\ncont_feats = [col for col in df_imputed.columns if df_imputed[col].dtype != object]\ncont_feats.remove('RainTomorrow')\ncont_feats.remove('RainToday')\n\n\nimputer = IterativeImputer(random_state=42)\ndf_imputed_cont = imputer.fit_transform(df_imputed[cont_feats])\ndf_imputed_cont = pd.DataFrame(df_imputed_cont, columns=cont_feats)\n\n\ncat_feats = [col for col in df_imputed.columns if col not in cont_feats]\ncat_feats.remove('RainTomorrow')\n\n# Also removing Date and Location since no values are missing\ncat_feats.remove('Date')\ncat_feats.remove('Location')\n\n\nimport numpy as np\n\ndf_imputed_cat = df_imputed[cat_feats]\n\nfor col in df_imputed_cat.columns:\n    # Find missing values in the current column\n    missing_values = df_imputed_cat[col].isnull()\n\n    # Calculate probabilities based on non-missing values\n    probabilities = df_imputed_cat[col][~missing_values].value_counts(normalize=True)\n\n    # Replace missing values with random choice based on probabilities\n    df_imputed_cat.loc[missing_values, col] = np.random.choice(probabilities.index,\n                                                              size=np.sum(missing_values),\n                                                              p=probabilities.values)\n\n\ndf_date_loc = df_imputed[['Date', 'Location']]\ndf_target = df_imputed.RainTomorrow\n\n\ndf_imputed_final = pd.concat(objs=[df_date_loc.reset_index(drop=True), df_imputed_cont.reset_index(drop=True),\n                                   df_imputed_cat.reset_index(drop=True), df_target.reset_index(drop=True)], axis=1)\n``` |
| Data Transformation | ```python\ndf_month = df_imputed_final.copy()\ndf_month.insert(1, 'Month', df_month.Date.apply(lambda x: int(str(x)[5:7])))\ndf_month.drop(columns='Date', inplace=True)\n\n\nfrom sklearn.preprocessing import LabelEncoder\nle=LabelEncoder()\n\ndf_month['Month']=le.fit_transform(df_month['Month'])\n\ndf_month['Location']=le.fit_transform(df_month['Location'])\n\ndf_month['WindGustDir']=le.fit_transform(df_month['WindGustDir'])\n\ndf_month['WindDir9am']=le.fit_transform(df_month['WindDir9am'])\n\ndf_month['WindDir3pm']=le.fit_transform(df_month['WindDir3pm'])\n\ndf_month['RainToday']=le.fit_transform(df_month['RainToday'])\n\ndf_month['RainTomorrow']=le.fit_transform(df_month['RainTomorrow'])\n``` |
| Feature Engineering | Attached the codes in final submission. |
| Save Processed Data | ```python\n# Saving the preprocessed data\ndf_final = df_month.copy()\n``` |