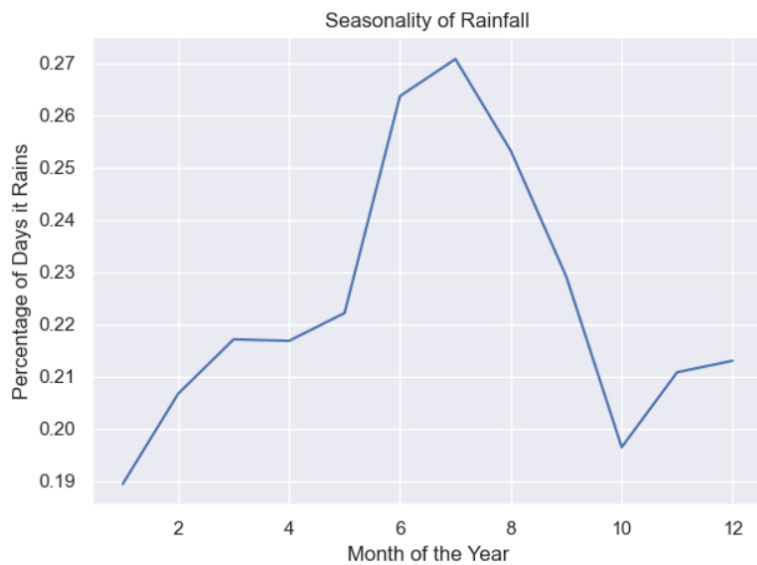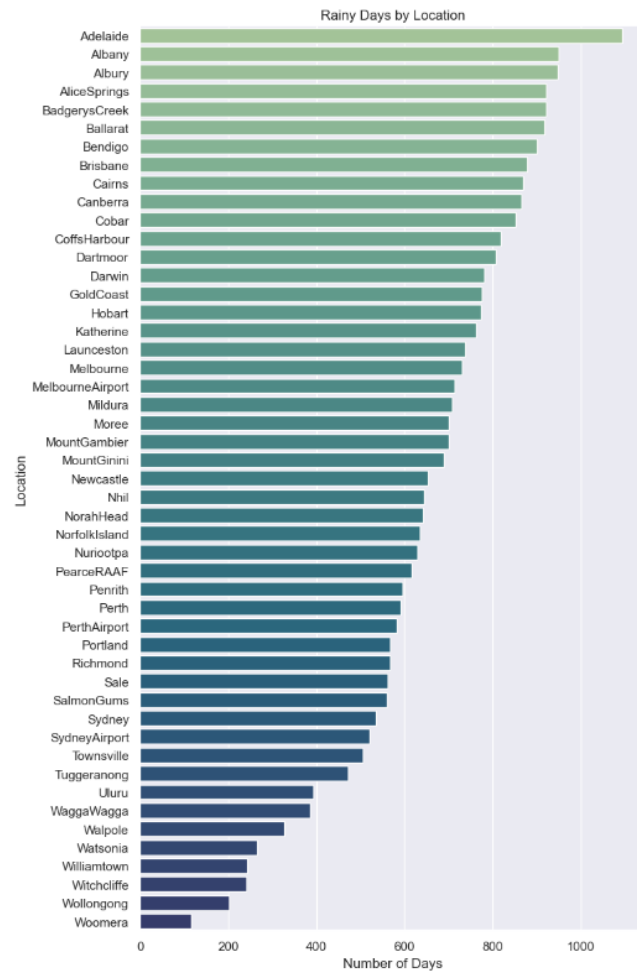# Data Exploration and Preprocessing Report

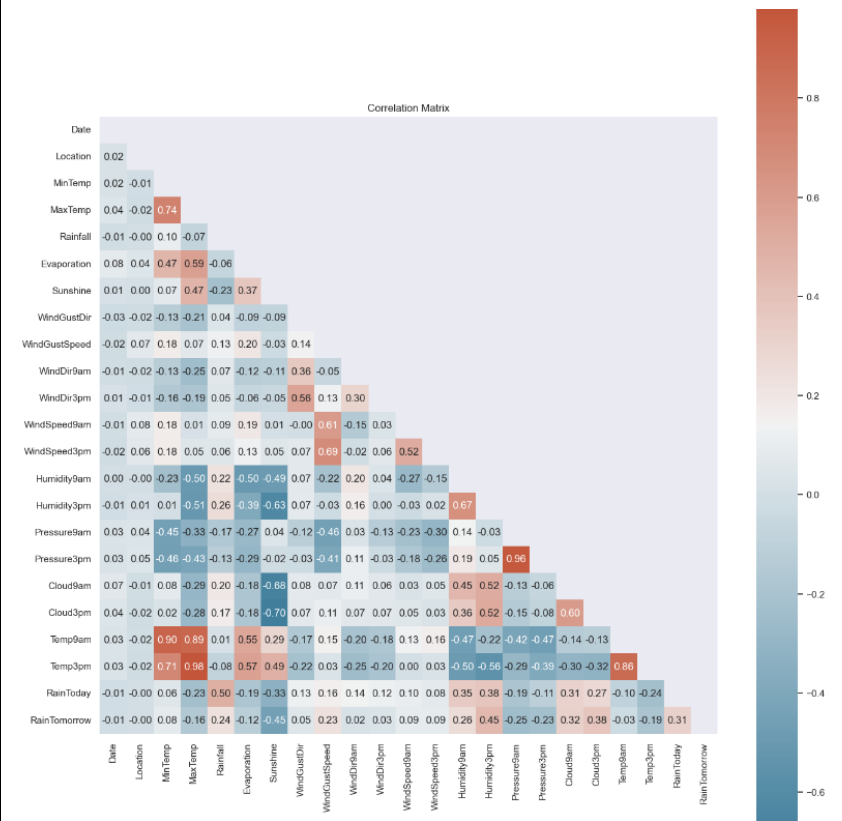| Date | 15 April 2024 |
|---|---|
| Team ID | Team-738164 |
| Project Title | Rainfall Prediction Using Machine Learning |
| Maximum Marks | 6 Marks |

**Data Exploration Screenshots:**

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

Bivariate Analysis

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

df.head()
```

| | 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
df_imputed = df.dropna(axis=0, subset=['RainTomorrow'])


cont_feats = [col for col in df_imputed.columns if df_imputed[col].dtype != object]
cont_feats.remove('RainTomorrow')
cont_feats.remove('RainToday')


imputer = IterativeImputer(random_state=42)
df_imputed_cont = imputer.fit_transform(df_imputed[cont_feats])
df_imputed_cont = pd.DataFrame(df_imputed_cont, columns=cont_feats)


cat_feats = [col for col in df_imputed.columns if col not in cont_feats]
cat_feats.remove('RainTomorrow')

# Also removing Date and Location since no values are missing
cat_feats.remove('Date')
cat_feats.remove('Location')


import numpy as np

df_imputed_cat = df_imputed[cat_feats]

for col in df_imputed_cat.columns:
    # Find missing values in the current column
    missing_values = df_imputed_cat[col].isnull()

    # Calculate probabilities based on non-missing values
    probabilities = df_imputed_cat[col][~missing_values].value_counts(normalize=True)

    # Replace missing values with random choice based on probabilities
    df_imputed_cat.loc[missing_values, col] = np.random.choice(probabilities.index,
                                                               size=np.sum(missing_values),
                                                               p=probabilities.values)


df_date_loc = df_imputed[['Date', 'Location']]
df_target = df_imputed.RainTomorrow


df_imputed_final = pd.concat(objs=[df_date_loc.reset_index(drop=True), df_imputed_cont.reset_index(drop=True),
                                   df_imputed_cat.reset_index(drop=True), df_target.reset_index(drop=True)], axis=1)
``` |
|---|---|
| Data Transformation | ```python
df_month = df_imputed_final.copy()
df_month.insert(1, 'Month', df_month.Date.apply(lambda x: int(str(x)[5:7])))
df_month.drop(columns='Date', inplace=True)


from sklearn.preprocessing import LabelEncoder
le=LabelEncoder()

df_month['Month']=le.fit_transform(df_month['Month'])

df_month['Location']=le.fit_transform(df_month['Location'])

df_month['WindGustDir']=le.fit_transform(df_month['WindGustDir'])

df_month['WindDir9am']=le.fit_transform(df_month['WindDir9am'])

df_month['WindDir3pm']=le.fit_transform(df_month['WindDir3pm'])

df_month['RainToday']=le.fit_transform(df_month['RainToday'])

df_month['RainTomorrow']=le.fit_transform(df_month['RainTomorrow'])
``` |
| Feature Engineering | Attached the codes in final submission. |
| Save Processed Data | ```python
# Saving the preprocessed data
df_final = df_month.copy()
``` |