

Rainfall Weather forecasting Project

In this Project, I will be implementing a predictive model on Rain Dataset to predict whether or not it will rain tomorrow in Australia. The Dataset contains about 10 years of daily weather observations of different locations in Australia. By the end of this article, you will be able to build a predictive model.

Design a predictive model with the use of machine learning algorithms to forecast whether or not it will rain tomorrow in Australia.

Importing Libraries

The first step in any Data Analysis step is importing necessary libraries.

```
In [5]: import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
```

Loading the Dataset

Dataset can be loaded using a method read_csv

```
In [2]: df=pd.read_csv("https://raw.githubusercontent.com/dscrentlist/dataset3/main/weatherAUS.csv")
df.head()
```

```
Out[2]:
```

	Date	Location	MinTemp	MaxTemp	Rainfall	Evaporation	Sunshine	WindGustDir	WindGustSpeed	WindDirbm	...	Humiditybm	Humidity3pm	Pressurebm	Pressure3pm	Cloudbm	Cloud3pm	Tempbm	Temp3pm	RainToday	RainTomorrow		
0	2008-12-01	Albury	13.4	22.9	0.0	0.0	NAN	NAN	W	44.0	W	...	71.0	22.0	1007.7	1007.1	8.0	NAN	16.9	21.8	No	No	
1	2008-12-02	Albury	7.4	25.1	0.0	0.0	NAN	NAN	WNW	44.0	NW	...	44.0	25.0	1010.6	1007.8	NAN	NAN	17.2	24.3	No	No	
2	2008-12-03	Albury	12.9	25.7	0.0	0.0	NAN	NAN	WSW	46.0	W	...	38.0	30.0	1007.6	1008.7	NAN	NAN	2.0	21.0	23.2	No	No
3	2008-12-04	Albury	8.2	28.0	0.0	0.0	NAN	NAN	NE	24.0	SE	...	45.0	16.0	1017.8	1012.8	NAN	NAN	18.1	26.5	No	No	
4	2008-12-05	Albury	17.5	32.3	1.0	0.0	NAN	NAN	W	41.0	ENE	...	82.0	33.0	1010.8	1006.0	7.0	8.0	17.8	29.7	No	No	

5 rows × 23 columns

Checking the Dimensions of Dataset

The shape property is used to find the dimenrions of the dataset

```
In [15]: df.shape
Out[15]: (8425, 23)
```

Data Preprocessing

Real-world data is often messy, incomplete, unstructured, inconsistent, redundant, sprinkled with wacky values. So, without deploying any Data Preprocessing techniques, it is almost impossible to gain insights from raw data.

What exactly is Data Preprocessing?

Data preprocessing is a process of converting raw data to a suitable format to extract insights. It is the first and foremost step in the Data Science life cycle. Data Preprocessing makes sure that data is clean, organize and read-to-feed to the Machine Learning model.

In detailed summary of a Dataset

```
In [17]: df.info()
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 8425 entries, 0 to 8424
Data columns (total 23 columns):
 #   Column      Non-Null Count  Dtype  
---  --
 0   Date        8425 non-null   object 
 1   Location    8425 non-null   object 
 2   MinTemp     8390 non-null   float64
 3   MaxTemp     8365 non-null   float64
 4   Rainfall    8185 non-null   float64
 5   Evaporation 4913 non-null   float64
 6   Sunshine    4423 non-null   float64
 7   WindGustDir 7434 non-null   object 
 8   WindDirbm   7596 non-null   object 
 9   WindDir3pm  8117 non-null   object 
10   WindSpeed9am 8349 non-null   float64
11   WindSpeed3pm 8318 non-null   float64
12   Humidity9am 8366 non-null   float64
13   Humidity3pm  8323 non-null   float64
14   Pressure9am 7113 non-null   float64
15   Pressure3pm 7113 non-null   float64
16   Cloud9am    6864 non-null   float64
17   Cloud3pm    6979 non-null   float64
18   Temp9am     8369 non-null   float64
19   Temp3pm     8329 non-null   float64
20   RainToday   8329 non-null   object 
21   RainTomorrow 8185 non-null   object 
dtypes: float64(16), object(7)
memory usage: 1.1+ MB

Dataset has two data types, float64, object

Except for the Date, Location columns, every column has missing values. Let's generate descriptive statistics for the dataset using the function describe() in pandas.

Descriptive Statistics: It is used to summarize and describe the features of data in a meaningful way to extract insights. It uses two types of statistic to describe or summarize data:
```

```
In [18]: df.describe()
```

```
Out[18]:
```

	MinTemp	MaxTemp	Rainfall	Evaporation	Sunshine	WindGustDir	WindSpeed9am	WindSpeed3pm	Humidity9am	Humidity3pm	Pressure9am	Pressure3pm	Cloud9am	Cloud3pm	Temp9am	Temp3pm	RainToday	RainTomorrow
count	8390.000000	8365.000000	8185.000000	4913.000000	4423.000000	7434.000000	8349.000000	8318.000000	8366.000000	8323.000000	7113.000000	7113.000000	6864.000000	6979.000000	8369.000000	8329.000000	8425.000000	8325.000000
mean	13.193205	23.699760	2.869913	5.389395	7.632205	40.174469	13.847646	18.533652	67.822496	51.249790	1017.640233	1023.6	4.566222	4.503183	17.762015	22.442934		
std	5.403696	6.136408	10.459279	5.044444	3.861223	14.166721	10.174579	9.766866	18.832283	18.423774	6.828669	6.787658	2.877658	2.731659	5.827035	5.980200		
min	-2.000000	19.300000	0.000000	0.000000	0.000000	7.000000	0.000000	0.000000	0.000000	10.000000	989.800000	982.500000	0.000000	0.000000	1.900000	7.300000		
25%	9.200000	18.300000	0.000000	0.000000	4.750000	30.000000	11.000000	11.000000	56.000000	39.000000	1013.000000	1012.000000	1.000000	2.000000	13.800000	18.000000		
50%	13.300000	23.300000	0.000000	4.600000	6.700000	39.000000	13.000000	19.000000	68.000000	51.000000	1017.700000	1015.300000	5.000000	5.000000	17.800000	21.900000		
75%	17.400000	28.000000	1.000000	7.000000	10.700000	50.000000	24.000000	24.000000	80.000000	63.000000	1022.300000	1019.800000	7.000000	7.000000	21.900000	26.400000		
max	28.300000	45.900000	371.000000	145.000000	13.900000	107.000000	63.000000	83.000000	100.000000	99.000000	1039.000000	1036.000000	8.000000	8.000000	39.400000	44.100000		

```
In [20]: print(df.describe(include=object))

count      Date      Location  WindGustDir  WindDirbm  WindDir3pm  RainToday  \
unique      8425      8425      7434      7596      8117      8185
top      2011-02-01  Melbourne      N      N      SE      No
freq         5      1622      713      906      813      6195

count      RainTomorrow
unique      2
top      No
freq      6195
```

```
In [61]: df.tail(5)
```

```
Out[6]:
```

	Date	Location	MinTemp	MaxTemp	Rainfall	Evaporation	Sunshine	WindGustDir	WindGustSpeed	WindDirbm	...	Humidity9am	Humidity3pm	Pressure9am	Pressure3pm	Cloud9am	Cloud3pm	Temp9am	Temp3pm	RainToday	RainTomorrow
8419	2017-06-30	Ulluru	3.5	21.8	0.0	NAN	NAN	E	31.0	ESE	...	99.0	27.0	1024.7	1021.2	NAN	NAN	9.4	20.9	No	No
8420	2017-06-21	Ulluru	2.8	23.4	0.0	NAN	NAN	E	31.0	SE	...	51.0	24.0	1024.6	1020.3	NAN	NAN	10.1	22.4	No	No
8421	2017-06-22	Ulluru	3.6	25.3	0.0	NAN	NAN	NW	22.0	SE	...	96.0	21.0	1019.1	1019.3	NAN	NAN	10.9	24.5	No	No
8422	2017-06-23	Ulluru	5.4	26.9	0.0	NAN	NAN	N	37.0	SE	...	53.0	24.0	1021.0	1016.8	NAN	NAN	12.5	26.1	No	No
8423	2017-06-24	Ulluru	7.8	27.0	0.0	NAN	NAN	SE	28.0	SSE	...	51.0	24.0	1019.4	1016.5	3.0	2.0	15.1	26.0	No	No

6 rows × 23 columns

```
In [18]: df.isnull().sum()
```

```
Out[18]:
```

	Date	Location	MinTemp	MaxTemp	Rainfall	Evaporation	Sunshine	WindGustDir	WindSpeed9am	WindSpeed3pm	Humidity9am	Humidity3pm	Pressure9am	Pressure3pm	Cloud9am	Cloud3pm	Temp9am	Temp3pm	RainToday	RainTomorrow
count	8425	8425	7434	7596	8117	8185														
unique	3804	72	16	16	16	2														
top	2011-02-01	Melbourne	N	N	SE	No														
freq	5	1622	713	906	813	6195														

```
In [11]: df.isnull().sum().sum()
```

```
Out[11]: 19472
```

```
In [14]: df.columns
```

```
Out[14]: Index(['Date', 'Location', 'MinTemp', 'MaxTemp', 'Rainfall', 'Evaporation',
        'Sunshine', 'WindGustDir', 'WindGustSpeed', 'WindDirbm', 'WindDir3pm', '...', 'Pressure3pm', 'Cloud9am', 'Cloud3pm', 'Temp9am', 'Temp3pm', 'RainToday', 'RainTomorrow',
        'year', 'month', 'day'],
        dtype='object')
```

Finding Categorical and Numerical Features in a Data set

Categorical Features in Dataset:

```
In [21]: categorical_features = [column_name for column_name in df.columns if df[column_name].dtype == 'O']
print("Number of Categorical Features: {}".format(len(categorical_features)))
print("Categorical Features: {}".format(categorical_features))

Number of Categorical Features: 7
Categorical Features: ['Date', 'Location', 'WindGustDir', 'WindDirbm', 'WindDir3pm', 'RainToday', 'RainTomorrow']

Numerical Features in Dataset
```

```
In [22]: numerical_features = [column_name for column_name in df.columns if df[column_name].dtype != 'O']
print("Number of Numerical Features: {}".format(len(numerical_features)))
print("Numerical Features: {}".format(numerical_features))

Number of Numerical Features: 16
Numerical Features: ['MinTemp', 'MaxTemp', 'Rainfall', 'Evaporation', 'Sunshine', 'WindGustSpeed', 'WindSpeed9am', 'WindSpeed3pm', 'Humidity9am', 'Humidity3pm', 'Pressure9am', 'Pressure3pm', 'Cloud9am', 'Cloud3pm', 'Temp9am', 'Temp3pm']
```

Many machine learning algorithms like Linear Regression, Logistic Regression, k-nearest neighbors, etc. can handle only numerical data, so encoding categorical data to numeric becomes a necessary step. But before jumping into encoding, check the cardinality of each categorical feature.

This high cardinality feature poses many serious problems like it will increase the number of dimensions of data when that feature is encoded. This is not good for the model.

```
In [23]: for each_feature in categorical_features:
unique_values = len(df[each_feature].unique())
print("Cardinality(no. of unique values) of {} are: {}".format(each_feature, unique_values))

Cardinality(no. of unique values) of Date are: 3804
Cardinality(no. of unique values) of Location are: 12
Cardinality(no. of unique values) of WindGustDir are: 17
Cardinality(no. of unique values) of WindDirbm are: 27
Cardinality(no. of unique values) of WindDir3pm are: 37
Cardinality(no. of unique values) of RainToday are: 3
Cardinality(no. of unique values) of RainTomorrow are: 3

Date column has high cardinality which poses several problems to the model in terms of efficiency and also dimensions of data increase when encoded to numerical data.
```

```
In [2]: df.head()
```

```
Out[2]:
```

	Location	MinTemp	MaxTemp	Rainfall	Evaporation	Sunshine	WindGustDir	WindGustSpeed	WindDirbm	WindDir3pm	...	Pressure3pm	Cloud9am	Cloud3pm	Temp9am	Temp3pm	RainToday	RainTomorrow	year	month	day		
0	Albury	13.4	22.9	0.0	0.0	NAN	NAN	W	44.0	W	WNW	...	1007.1	8.0	NAN	16.9	21.8	No	No	2008	12	1	
1	Albury	7.4	25.1	0.0	0.0	NAN	NAN	WNW	44.0	NW	WSW	...	1007.8	NAN	NAN	17.2	24.3	No	No	2008	12	2	
2	Albury	12.9	25.7	0.0	0.0	NAN	NAN	WSW	46.0	W	WSW	...	1008.7	NAN	NAN	2.0	21.0	23.2	No	No	2008	12	3
3	Albury	9.2	28.0	0.0	0.0	NAN	NAN	NE	24.0	SE	E	...	1012.8	NAN	NAN	18.1	26.5	No	No	2008	12	4	
4	Albury	17.5	32.3	1.0	0.0	NAN	NAN	W	41.0	ENE	NW	...	1006.0	7.0	8.0	17.8	29.7	No	No	2008	12	5	

5 rows × 25 columns

Handling Missing Values

Machine learning algorithms can't handle missing values and cause problems. So they need to be addressed in the first place. There are many techniques to identify and impute missing values.

If a dataset contains missing values and loaded using pandas, then missing values get replaced with NaN(Not a Number) values. These NaN values can be identified using methods like isna() or isnull() and they can be imputed using fillna(). This process is known as Missing Data Imputation.

Handling Missing values in Categorical Features

```
In [39]: categorical_features = [column_name for column_name in df.columns if df[column_name].dtype == 'O']
df[categorical_features].isnull().sum()
```

```
Out[39]:
```

	Location	MinTemp	MaxTemp	Rainfall	Evaporation	Sunshine	WindGustDir	WindDirbm	WindDir3pm	...	Pressure3pm	Cloud9am	Cloud3pm	Temp9am	Temp3pm	RainToday	RainTomorrow
count	8																
unique	991																
top	829																
freq	248																
dtype	int64																

Imputing the missing values in categorical features using the most frequent value which is mode

```
In [42]: categorical_features_with_null = [feature for feature in categorical_features if df[feature].isnull().sum()]
for each_feature in categorical_features_with_null:
mode_val = df[each_feature].mode()[0]
df[each_feature].fillna(mode_val,inplace=True)

Handling Missing values in Numerical features
```

```
In [44]: numerical_features = [column_name for column_name in df.columns if df[column_name].dtype != 'O']
df[numerical_features].isnull().sum()
```

```
Out[44]:
```

	MinTemp	MaxTemp	Rainfall	Evaporation	Sunshine	WindGustDir	WindSpeed9am	WindSpeed3pm	Humidity9am	Humidity3pm	Pressure9am	Pressure3pm	Cloud9am	Cloud3pm	Temp9am	Temp3pm	RainToday	RainTomorrow	year	month	day
count	75																				
unique	69																				
top	249																				
freq	3512																				
dtype	float64																				

Missing values in Numerical Features can be imputed using Mean and Median. Mean is sensitive to outliers and median is immune to outliers. If you want to impute the missing values with mean values, then outliers in numerical features need to be addressed properly.

Outliers detection and treatment

What is an outlier?

An Outlier is an observation that lies an abnormal distance from other values in a given sample. They can be detected using visualization(like boxplots, scatter plots), Z-score, statistical and probabilistic algorithms

Outlier Treatment to remove outliers from Numerical Features

```
In [ ]:

In [ ]:

In [ ]: features_with_outliers = ['MinTemp', 'MaxTemp', 'Rainfall', 'Evaporation', 'WindGustSpeed', 'WindSpeed9am', 'WindSpeed3pm', 'Humidity9am', 'Pressure9am', 'Pressure3pm', 'Temp9am', 'Temp3pm']
for feature in features_with_outliers:
q1 = df[feature].quantile(0.25)
q3 = df[feature].quantile(0.75)
IQR = q3 - q1
lower_limit = q1 - (IQR*1.5)
upper_limit = q3 + (IQR*1.5)
df.loc[df[feature]<lower_limit, feature] = lower_limit
df.loc[df[feature]>upper_limit, feature] = upper_limit

Now, numerical features are free from outliers. Let's Impute missing values in numerical features using mean.
```

```
In [45]: numerical_features_with_null = [feature for feature in numerical_features if df[feature].isnull().sum()]
for feature in numerical_features_with_null:
mean_value = df[feature].mean()
df[feature].fillna(mean_value,inplace=True)

It's time to do some analysis on each feature to understand about data and get some insights.
```

Exploratory Data Analysis

Exploratory Data Analysis(EDA) is a technique used to analyze, visualize, investigate, interpret, discover and summarize data. It helps Data Scientists to extract trends, patterns, and relationships in data.

Exploring target variable

```
In [46]: df['RainTomorrow'].value_counts().plot(kind='bar')
```

```
Out[46]:
```



Sunshine vs Rainfall

```
In [47]: sns.lineplot(data=df,x='Sunshine',y='Rainfall',color='green')
```

```
Out[47]: <AxesSubplot: xlabel='Sunshine', ylabel='Rainfall'>
```



Sunshine vs Evaporation

```
In [48]: sns.lineplot(data=df,x='Sunshine',y='Evaporation',color='blue')
```

```
Out[48]: <AxesSubplot: xlabel='Sunshine', ylabel='Evaporation'>
```



In the above line plot, the Sunshine feature is proportional to the Evaporation feature

Correlation

Correlation is a statistic that helps to measure the strength of the relationship between two features. It is used in bivariate analysis. Correlation can be calculated with method corr() in pandas.

```
In [66]: df.corr()
```

```
Out[66]:
```

	Location	MinTemp	MaxTemp	Rainfall	Evaporation	Sunshine	WindGustDir	WindGustSpeed	WindDirbm	WindDir3pm	...	Pressure3pm	Cloud9am	Cloud3pm	Temp9am	Temp3pm	RainToday	RainTomorrow	year	month	day
Location	1.000000	0.000000	0.142541	-0.044802	0.031072	0.099510	0.188644	0.142151	-0.009839	-0.112645	...	-0.049051	-0.121723	-0.134822	0.179210	0.137959	-0.044982	-0.044719	0.009071	0.017288	0.020599
MinTemp	0.156760	1.000000	0.717522	0.097086	0.254965	0.056194	0.182399	0.230056	-0.065052	0.044324	...	-0.422991	0.091604	0.082614	0.867119	0.685					