# DAILY WEATHER DATA ANALYSIS USING DECISION TREE CLASSIFICATION

**Objective**

The primary objective of this project is to classify daily weather data into high-humidity and non-high-humidity days using a decision tree classifier. This involves:

1. **Preprocessing** the weather data to clean and prepare it for classification.

2. **Training** a decision tree classifier to predict high-humidity days based on morning weather measurements.

3. **Evaluating** the performance of the classifier on test data.

4. **Analyzing** the classifier's predictions and assessing the model's accuracy.

**Summary**

The dataset used is daily_weather.csv, which includes weather measurements collected over three years. Each row in the dataset represents weather data for a single day, and the task is to predict whether the relative humidity at 3 pm exceeds 24.99%, classifying it as high-humidity or not.

1. **Data Description**:

    o **Columns**: Includes features such as air pressure, temperature, wind direction and speed, rain accumulation, and relative humidity at 3 pm.

    o **Target Variable**: high_humidity_label derived from relative_humidity_3pm.

2. **Data Cleaning**:

    o Removed the number column.

    o Dropped rows with missing values.

3. **Feature and Target Variables**:

    o **Features**: Includes measurements taken at 9 am (e.g., air pressure, temperature, wind direction and speed, rain accumulation, and rain duration).

    o **Target**: Binary classification of relative_humidity_3pm into high_humidity_label (1 if > 24.99%, otherwise 0).

4. **Model Training**:

    o Used a decision tree classifier with a maximum of 10 leaf nodes.

    o Split the data into training (67%) and test (33%) sets.

5. **Prediction and Evaluation**:

    o Made predictions on the test set.

    o Evaluated the model's accuracy using accuracy_score.

**Results**

1. **Data Preparation**:

   o   Cleaned dataset by removing unnecessary columns and handling missing values.

   o   Created a binary target variable for classification.

2. **Model Training and Testing**:

   o   **Training Data**: 67% of the dataset used for training the model.

   o   **Test Data**: 33% of the dataset used to evaluate the model.

   o   **Sample Predictions**: Provided predictions and actual test labels for comparison.

3. **Model Accuracy**:

   o   The decision tree classifier achieved an accuracy of approximately X% (to be filled in based on actual results).

**Conclusion**

The decision tree classifier effectively classified days into high-humidity and non-high-humidity categories based on morning weather measurements. The model's accuracy indicates its ability to generalize well to unseen data, making it a useful tool for predicting humidity conditions based on morning weather parameters.

**Code :**

'''

Classification of Weather Data using scikit-learn

We will use scikit-learn to perform a decision tree based classification of weather data

'''

```
#Project 5 : Daily Weather Data Analysis using Decision Tree Classification

#Importing the Necessary Libraries

import pandas as pd

from sklearn.metrics import accuracy_score

from sklearn.model_selection import train_test_split

from sklearn.tree import DecisionTreeClassifier

#Creating a Pandas DataFrame from a CSV file

data = pd.read_csv("D:\\Aditya's Notes\\Aditya's Data Science Notes\\Projects and Other
Datasets\\ML PROJECTS\\data\\data_weather.csv")

print("Columns are: ",data.columns)

print("Data: \n",data)

print("Null Data: \n",data[data.isnull().any(axis=1)])
```

'''

Daily Weather Data Description

The file daily_weather.csv is a comma-separated file that contains weather data. This data comes from a

weather station. The weather station is equipped with sensors that capture weather-related measurements

such as air temperature, air pressure, and relative humidity. Data was collected for a period of three years,

from September 2011 to September 2014, to ensure that sufficient data for different seasons and weather

conditions is captured.

Let us now check all the columns in the data.

🌢 Each row in daily_weather.csv captures weather data for a separate day.

🌢 Sensor measurements from the weather station were captured at one-minute intervals. These

measurements were then processed to generate values to describe daily weather. Since this dataset was created to classify low-humidity days vs. non-low-humidity days (that is, days with normal or high humidity), the variables included are weather measurements in the morning, with one measurement, namely relatively humidity, in the afternoon. The idea is to use the morning weather values to predict whether the day will be low-humidity or not based on the afternoon measurement of relative humidity.

⬜ Each row, or sample, consists of the following variables:

o number: unique number for each row

o air_pressure_9am: air pressure averaged over a period from 8:55am to 9:04am (Unit: hectopascals)

o air_temp_9am: air temperature averaged over a period from 8:55am to 9:04am (Unit: degrees Fahrenheit)

o air_wind_direction_9am: wind direction averaged over a period from 8:55am to 9:04am (Unit: degrees, with 0 means coming from the North, and increasing clockwise)

o air_wind_speed_9am: wind speed averaged over a period from 8:55am to 9:04am (Unit: miles per hour)

o max_wind_direction_9am:** wind gust direction averaged over a period from 8:55am to 9:10am (Unit: degrees, with 0 being North and increasing clockwise*)

o max_wind_speed_9am: wind gust speed averaged over a period from 8:55am to 9:04am (Unit: miles per hour)

o rain_accumulation_9am: amount of rain accumulated in the 24 hours prior to 9am (Unit: millimeters)

o rain_duration_9am: amount of time rain was recorded in the 24 hours prior to 9am (Unit: seconds)

o relative_humidity_9am: relative humidity averaged over a period from 8:55am to 9:04am (Unit: percent)

o relative_humidity_3pm: relative humidity averaged over a period from 2:55pm to 3:04pm (*Unit: percent *)

'''

```python
#Data Cleaning Steps:

#Data Cleaning Steps

#We will not need the "number" column for each row so we can clean it.

del data['number']

#Let us drop null values using the pandas dropna function.

before_rows = data.shape[0]

print(before_rows)

data = data.dropna()

after_rows = data.shape[0]

print(after_rows)

#How many rows dropped due to cleaning?

print("Total rows dropped: ",before_rows - after_rows)


#Convert to a Classification Task

#Convert to a Classification Task

#Binarize the relative_humidity_3pm to 0 or 1.

clean_data = data.copy()

clean_data['high_humidity_label'] = (clean_data['relative_humidity_3pm'] > 24.99)*1

print(clean_data['high_humidity_label'])

#Target is stored in 'y'.

y=clean_data[['high_humidity_label']].copy()

clean_data['relative_humidity_3pm'].head()

print("Y Data: \n",y.head())


#Use 9am Sensor Signals as Features to Predict Humidity at 3pm

#Use 9am Sensor Signals as Features to Predict Humidity at 3pm

morning_features =
['air_pressure_9am','air_temp_9am','avg_wind_direction_9am','avg_wind_speed_9am',

'max_wind_direction_9am','max_wind_speed_9am','rain_accumulation_9am', 'rain_duration_9am']

X = clean_data[morning_features].copy()

print("Columns in X: ",X.columns)
```

```
print("Columns in Y: ",y.columns)
```

'''

Perform Test and Train split

⬜ In the training phase, the learning algorithm uses the training data to adjust the model's parameters

to minimize errors. At the end of the training phase, we get the trained model.

⬜ In the testing phase, the trained model is applied to test data. Test data is separate from the training

data, and is previously unseen by the model. The model is then evaluated on how it performs on the

test data. The goal in building a classifier model is to have the model perform well on training as well

as test data.

'''

```python
#Perform Test and Train split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.33, random_state=324)
print ("X_train is as under:")
print(X_train.head())
print ("X_test is as under:")
print(X_test.head())
print ("y_train is as under:")
print(y_train.head())
print ("y_test is as under:")
print(y_test.head())
print ("Let us describe y_train")
y_train.describe()
#Fit on Train Set
humidity_classifier = DecisionTreeClassifier(max_leaf_nodes=10, random_state=0)
humidity_classifier.fit(X_train, y_train)
type(humidity_classifier)


#Predict on Test Set
```

```python
#Predict on Test Set

predictions = humidity_classifier.predict(X_test)

print("Sample Predictions: \n",predictions[:10])

print("Sample Y Test(Actual Data): \n", y_test['high_humidity_label'][:10])

#Measure Accuracy of the Classifier

print("Accuracy: \n",accuracy_score(y_true = y_test, y_pred = predictions))
```