# Rain Prediction Using Machine Learning

**Tools**

# Table of Contents
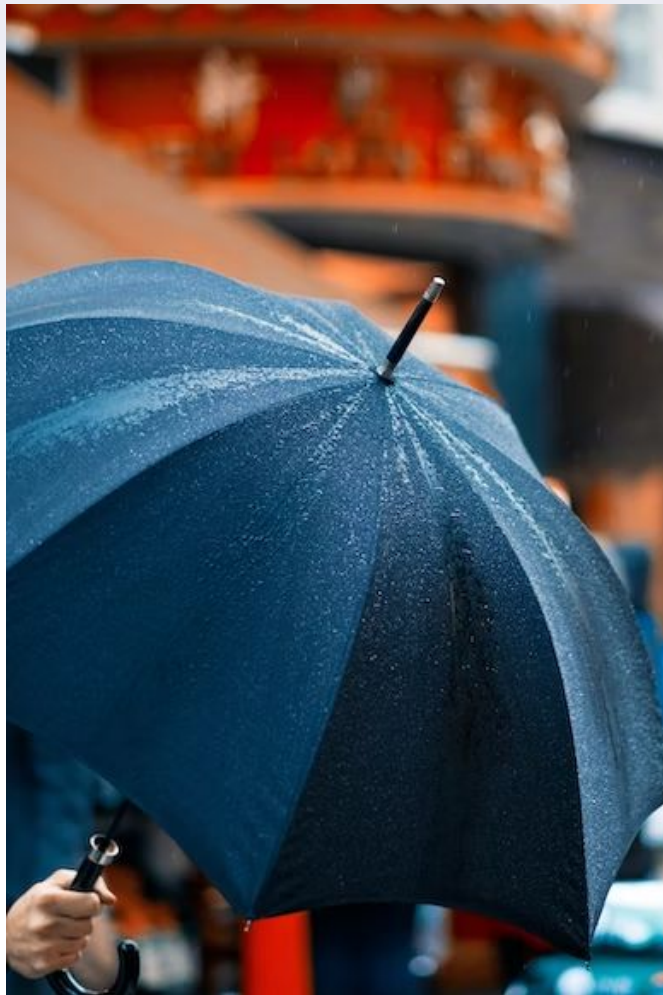
# 1 Business Background and Objectives

# Introduction and Problems

Weather has a significant impact on many life aspects, one of which is agricultural industry and because of that, being able to predict it helps farmers in their day-to-day decisions such as how to plan efficiently, minimize costs and maximize yields.

A major agricultural company needs to have an accurate **rain prediction algorithm** that will improve their decision-making on typical farming activities such as planting and irrigating.

Using historical rain information from Australia regions in 10 years as research data, it is necessary to **predict** weather(**rain**) in the **next day**.
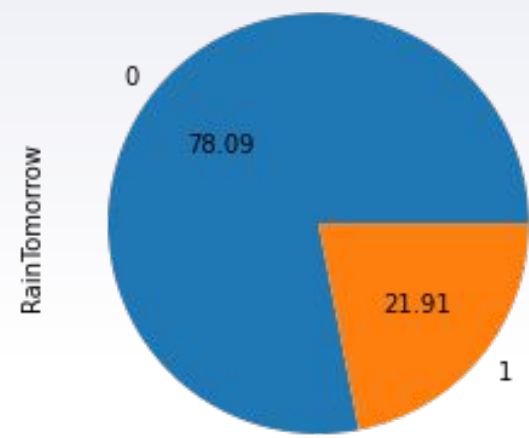
# Objectives

- What factors and conditions in current day that will cause rain in the next day?
- What machine learning algorithms are suitable for predicting rain in the next day?
- Where is the location that has the highest frequency of rain?
- What is the impact of the predictive model for business problems that operating in the agricultural sector?

**2**

# Data Preparation

# Dataset Overview



Target distribution is **imbalance**, I will use **ROC** curve and **AUC** score as model metric evaluation.

ROC-AUC is performance measurement for the classification problems at various threshold settings.

Data contains **145460 rows** with **22 features** and **1 binary column** with **RainTomorrow** as **target**.

All features are the weather and climate elements that occur on that day in certain location that are used to predict rain in the next day.

Features consists **16 numerical features** and **6 categorical features** (Date will be transform to datetime format).
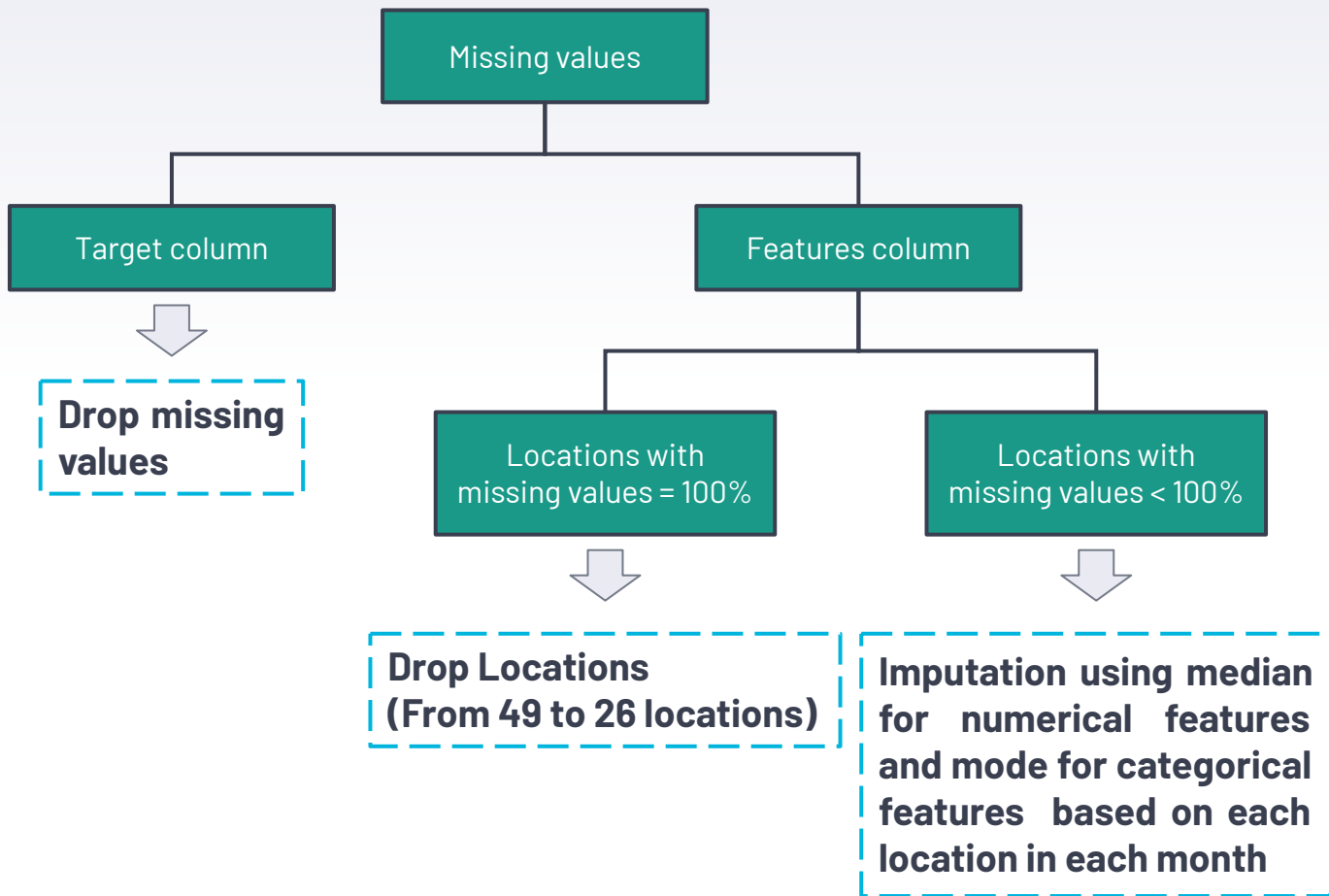
| Numerical Features | | Categorical Features |
|---|---|---|
| <ul><li>MinTemp</li><li>MaxTemp</li><li>Rainfall</li><li>Evaporation</li><li>Sunshine</li><li>WindGustSpeed</li><li>WindSpeed9am</li><li>WindSpeed3pm</li></ul> | <ul><li>Humidity9am</li><li>Humidity3pm</li><li>Pressure9am</li><li>Pressure3pm</li><li>Cloud9am</li><li>Cloud3pm</li><li>Temp9am</li><li>Temp3pm</li></ul> | <ul><li>Date</li><li>Location</li><li>WindGustDir</li><li>WindDir9am</li><li>WindDir3pm</li><li>RainToday</li></ul> |

Dataset obtained from **kaggle**.

# Missing Values Handling

## Missing Values Ratio

```
Location          0.000000
Month             0.000000
Year              0.000000
Day               0.000000
MinTemp           1.020899
MaxTemp           0.866905
Rainfall          2.241853
Evaporation      43.166506
Sunshine         48.009762
WindGustDir       7.098859
WindGustSpeed     7.055548
WindDir9am        7.263853
WindDir3pm        2.906641
WindSpeed9am      1.214767
WindSpeed3pm      2.105046
Humidity9am       1.824557
Humidity3pm       3.098446
Pressure9am      10.356799
Pressure3pm      10.331363
Cloud9am         38.421559
Cloud3pm         40.807095
Temp9am           1.214767
Temp3pm           2.481094
RainToday         2.241853
RainTomorrow      2.245978
dtype: float64
```
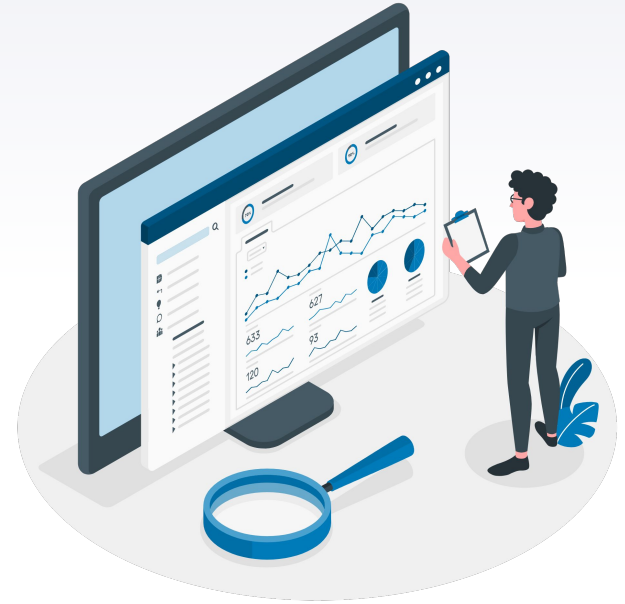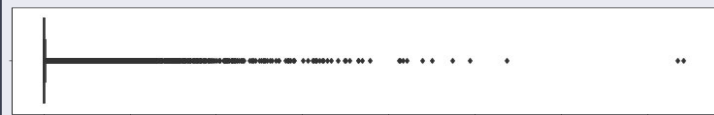
Missing values

Target column

Features column

Drop missing values

Locations with missing values = 100%

Locations with missing values < 100%

Drop Locations (From 49 to 26 locations)

Imputation using median for numerical features and mode for categorical features based on each location in each month

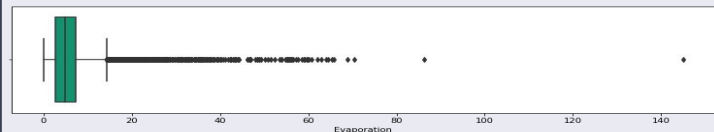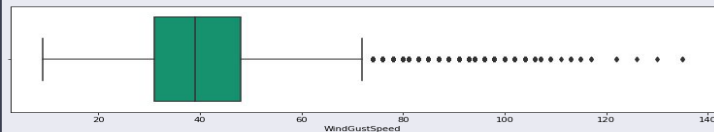**3** **Exploratory Data Analysis (EDA)**

# Outliers Check



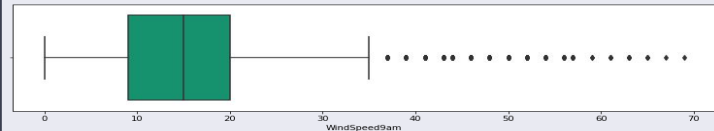Rainfall column outliers

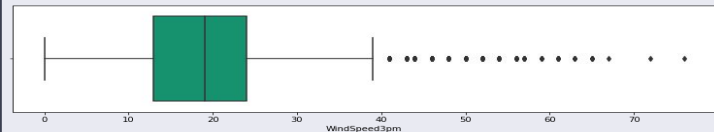Values > 2.4

Evaporation column outliers

Values > 21.2

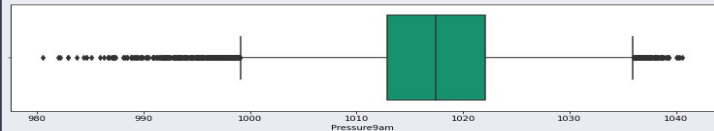WindGustSpeed column outliers

Values > 99.0

WindSpeed9am column outliers
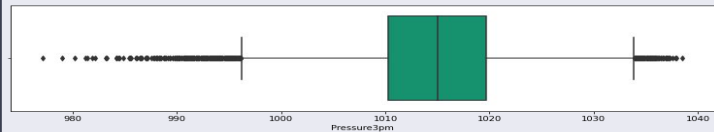
Values > 53.0

WindSpeed3pm column outliers

Values > 57.0

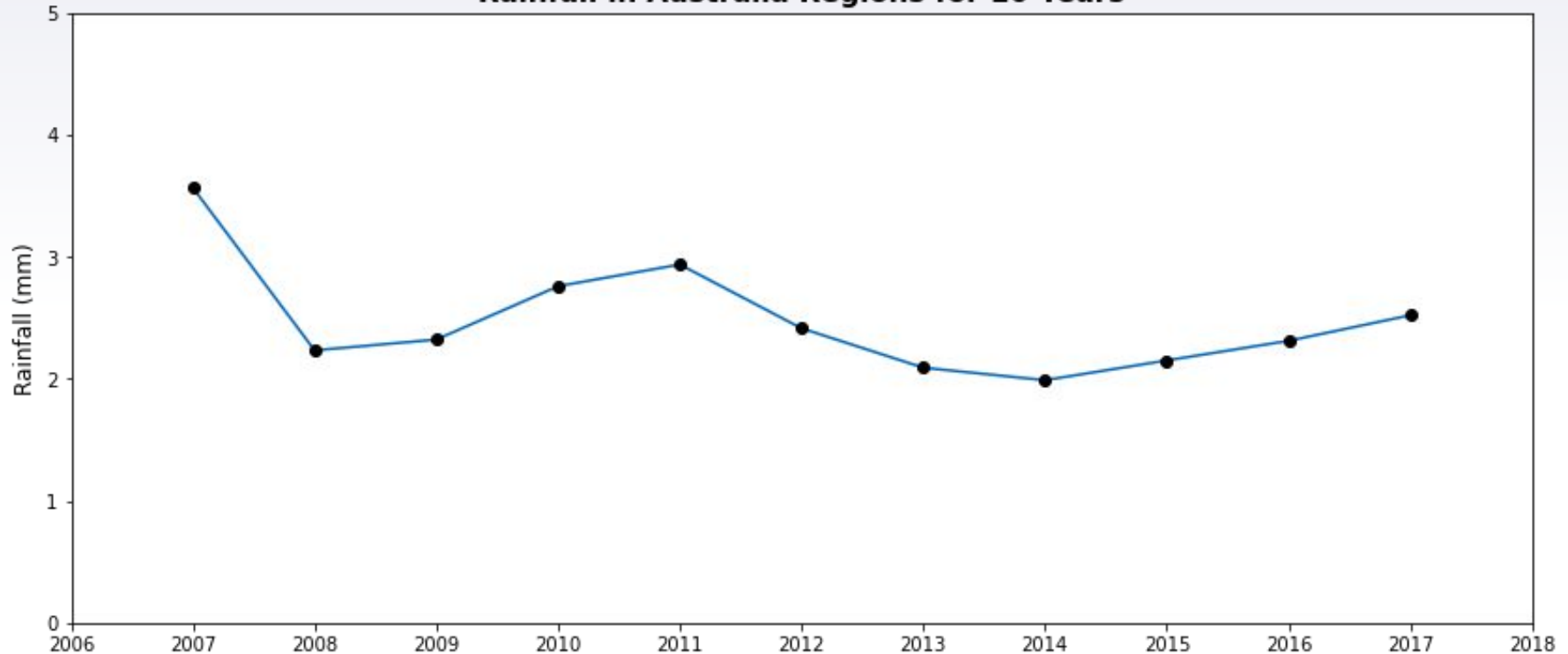Pressure9am column outliers

Values < 985.3

Pressure3pm column outliers

Values < 982.1

**Seven features** have extreme **outliers** and need to removed based **IQR** (Interquartile Range) **upper** and **lower** limit.
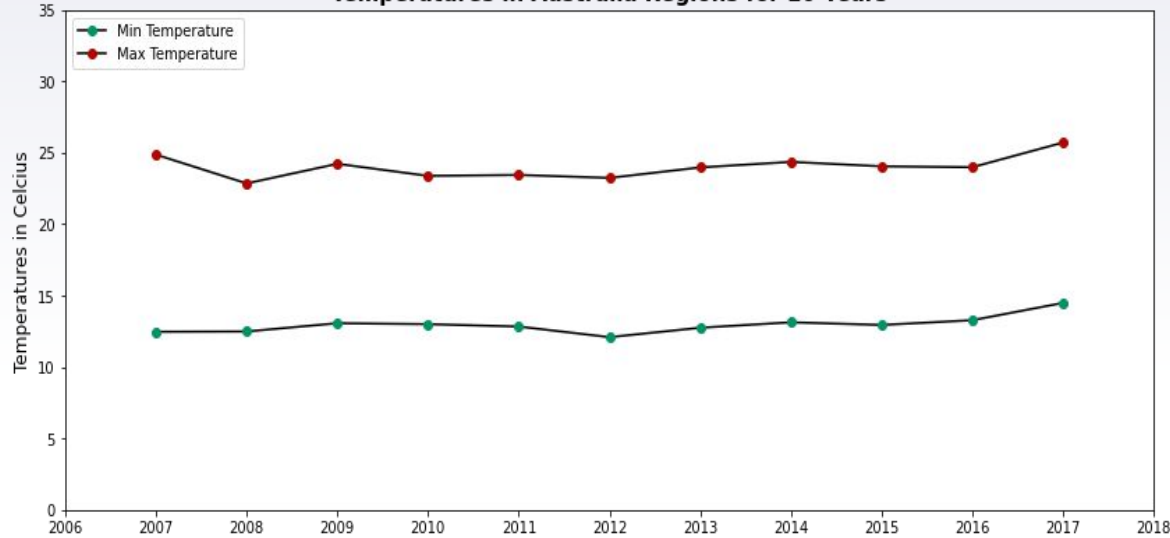
# Rainfall in Australia Regions



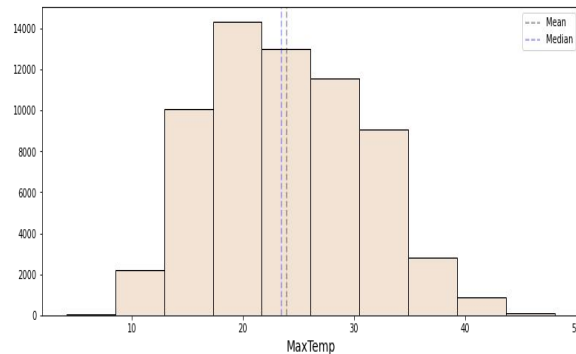**Rainfall in Australia Regions for 10 Years**

**Average rainfall** in Australia regions in 10 years relatively **decreased** over years. It seems the decreased rainfall caused by global warming issue.

# Temperatures in Australia Regions



**Average minimum** and **maximum temperatures** for 10 years **slightly increase**, due to global warming issue. For temperatures **data distribution**, relatively **normal distributed**. It can be concluded that Australia regions have a **relatively stable minimum** and **maximum temperatures**.

# Fraction of Sky by Cloud (in Oktas)



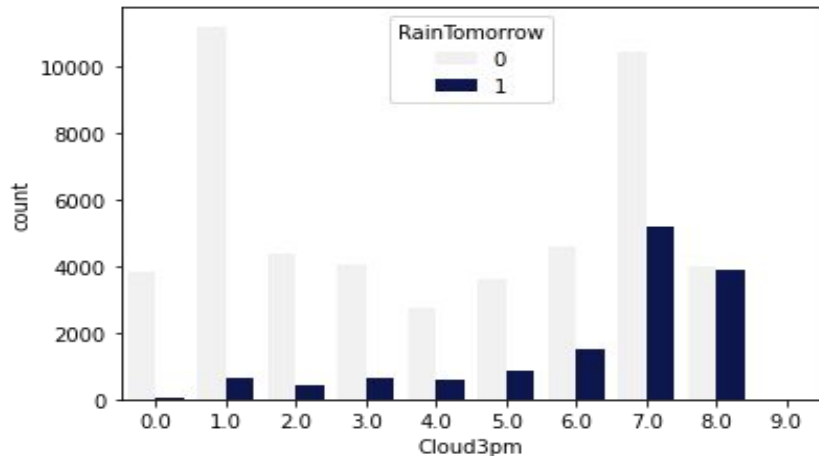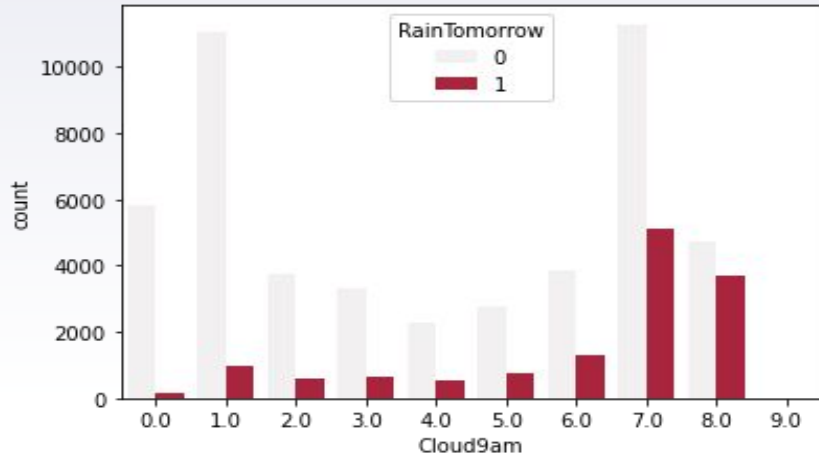**Sky** in Australia **mostly cloudy** (**Fraction** of **sky** by **cloud** is **7 oktas**). **Rain** in the **next day** tend to **happen** when **cloud** in **9am** and **3pm** is **7** or **8 oktas**.

| Oktas | Definition | Category |
|-------|------------|----------|
| 0 | Sky clear | Fine |
| 1 | 1/8 of sky covered or less, but not zero | Fine |
| 2 | 2/8 of sky covered | Fine |
| 3 | 3/8 of sky covered | Partly Cloudy |
| 4 | 4/8 of sky covered | Partly Cloudy |
| 5 | 5/8 of sky covered | Partly Cloudy |
| 6 | 6/8 of sky covered | Cloudy |
| 7 | 7/8 of sky covered or more, but not 8/8 | Cloudy |
| 8 | 8/8 of sky completely covered, no breaks | Overcast |

(Reference: Worldweather)

# Rain Frequency Based on Locations

## Top 10 Location With the Most Experiences Rain

| Location | Value |
|---|---|
| Portland | 887 |
| Cairns | 765 |
| NorfolkIsland | 741 |
| MountGambier | 739 |
| CoffsHarbour | 700 |
| Sydney | 692 |
| Darwin | 686 |
| SydneyAirport | 624 |
| Hobart | 610 |
| Watsonia | 577 |

**Portland** is the **location** in Australia that **experiences** the **most rain**. It means Portland is the location with the enough water source.

# Rain Condition in Portland
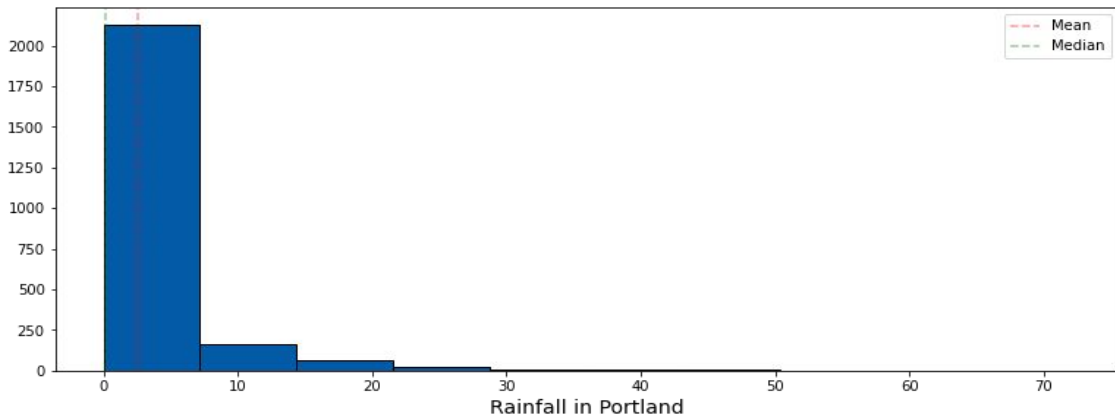


**Rain Weather Condition in Portland**

## Rain Period in Portland

**Rain in Portland** is most common to **happen** from **period May** to **August**. We can conclude that these period is the rainy season in Portland.
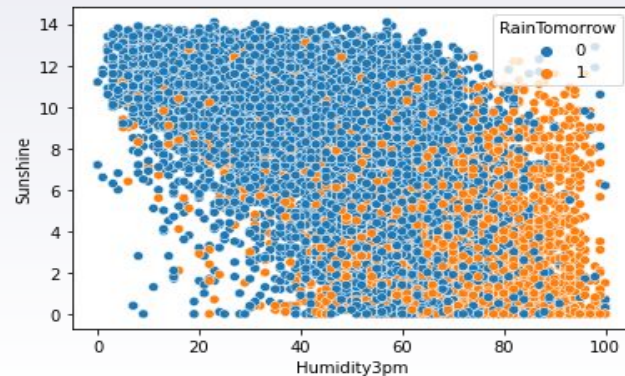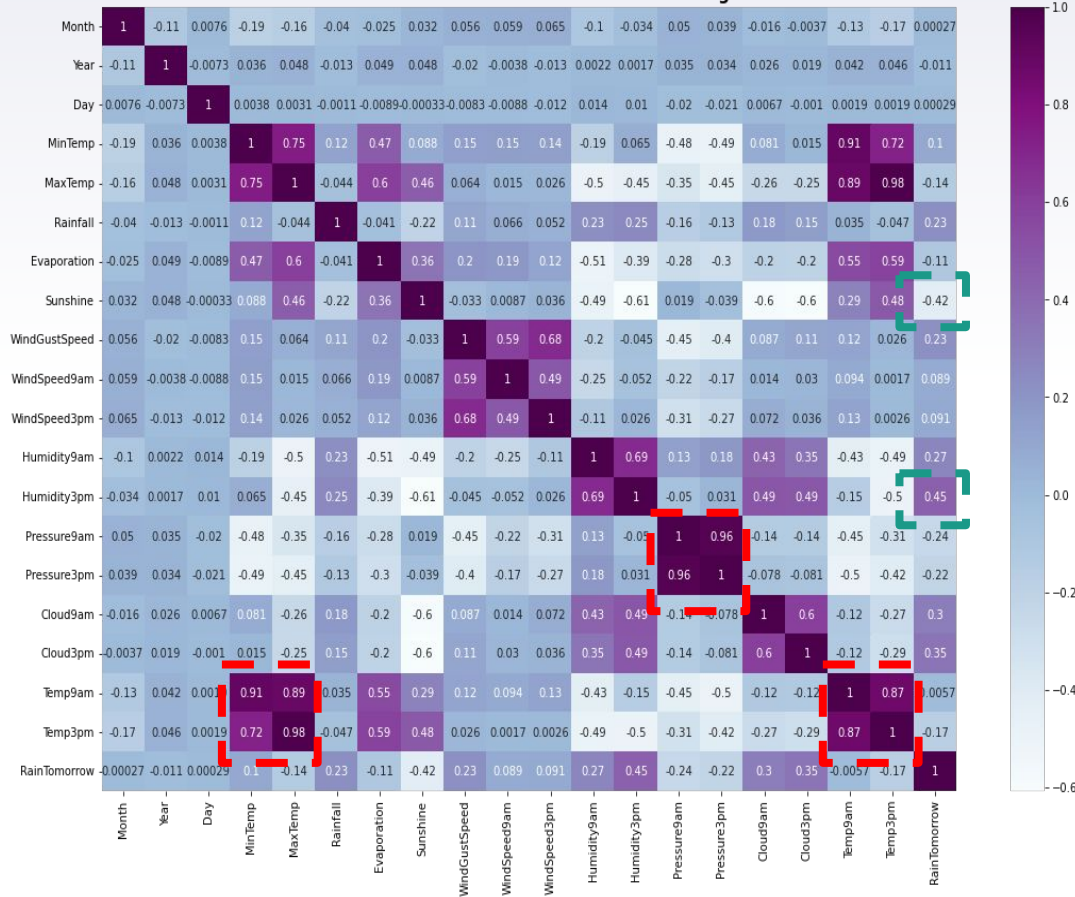
## Rainfall in Portland

**Rainfall in Portland** majority in **range** around **0 - 8 mm** with **average** of **rainfall** is **2.49 mm**.

# Multivariate Analysis


Correlation Between Numerical Value and Target



- **Humidity3pm** (**Positive**) and **Sunshine** (**Negative**) features have the **highest correlation** with **target**.
- **Pressure9am** and **Pressure3pm** then **MinTemp**, **MaxTemp**, **Temp9am**, and **Temp3pm** features have **high correlate each other**. **Pressure9am** and **Temp3pm** have **higher correlation with target** and will be kept for modelling.
- When **Humidity3pm ratio** is **high** and number of **hours** of **bright Sunshine** is **low**, **rain** in the **next day** tend to **happen**.
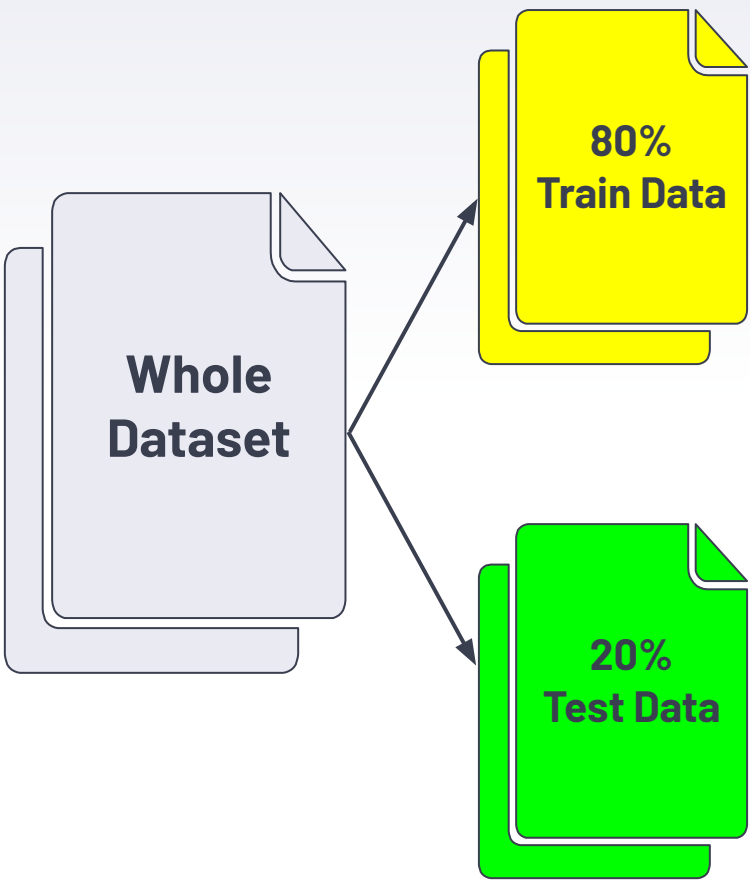
**4**

# Feature
# Engineering

# Dataset Split



**Whole Dataset**

**80% Train Data**

**20% Test Data**

# Feature Engineering

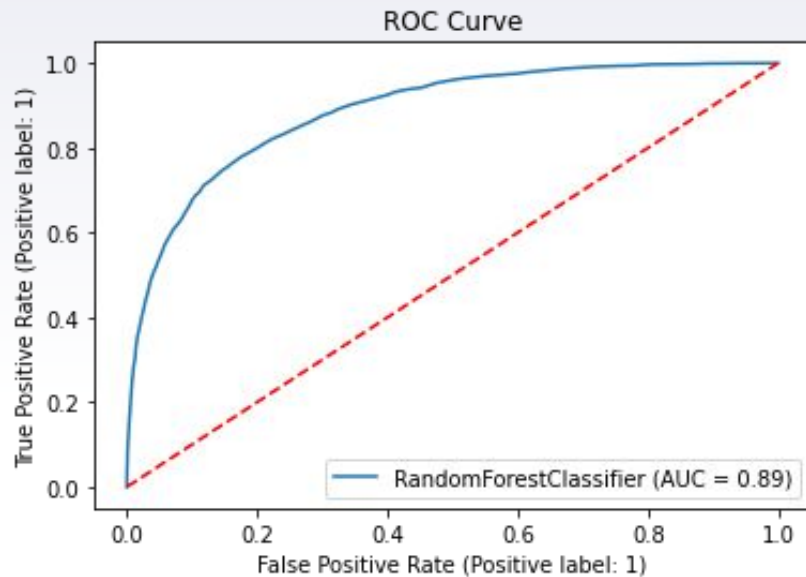| Aspects | Action | Treatment |
|---|---|---|
| Categorical features | Categorical feature with 2 distinct values -> Binary Encoding<br>Categorical feature with more than 2 distinct values -> One Hot Encoding | Train Data<br>Test Data |
| Features Drop | Multicollinearity -> Pressure3pm, MinTemp, MaxTemp, and Temp9am<br>Don't contribute for modelling -> Year | Train Data<br>Test Data |
| Imbalance Data | Resampling use Undersampling | Train Data |
| Scaling | Normalization using Min-Max Scaler | Train Data<br>Test Data |

# 5 Modelling and Evaluation

# Modelling

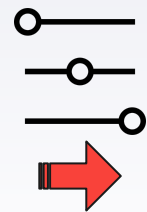| Model | AUC Score |
|---|---|
| K-Nearest Neighbors | 79% |
| Logistic Regression | 88% |
| Decision Tree | 72% |
| **Random Forest** | **89%** |
| XG-Boost | 88% |



ROC Curve

I **trained** data using **5 classification algorithms**. Since the **data** is **imbalance**, I use **ROC** curve and **AUC** score as **metric evaluation**. **Random Forest** have the **best AUC score** compared than other algorithms.

# Hyperparameter Tuning

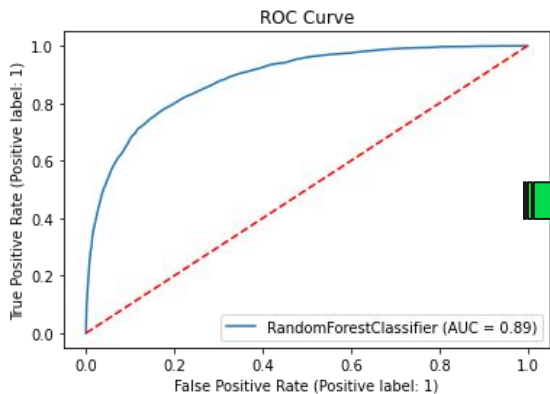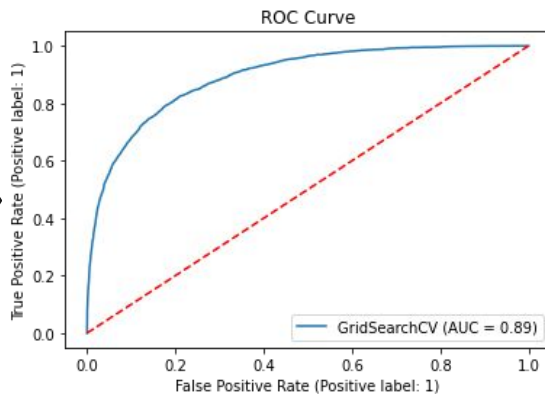| Parameter | Value |
|---|---|
| n_estimators (Number of trees) | 600, 700, 800, 900, 1000 |
| max_depth (Longest path between the root node and the leaf node) | 60, 70, 80, 90, 100 |
| criterion (Function to measure the quality of a split) | 'gini', 'entropy' |

**Tuning**

**GridSearch CV**

| Parameter | Best Value |
|---|---|
| n_estimators | 800 |
| max_depth | 60 |
| criterion | 'entropy' |

**Before Tuned**

**After Tuned**

ROC Curve

RandomForestClassifier (AUC = 0.89)

False Positive Rate (Positive label: 1)

True Positive Rate (Positive label: 1)

ROC Curve

GridSearchCV (AUC = 0.89)

False Positive Rate (Positive label: 1)

True Positive Rate (Positive label: 1)
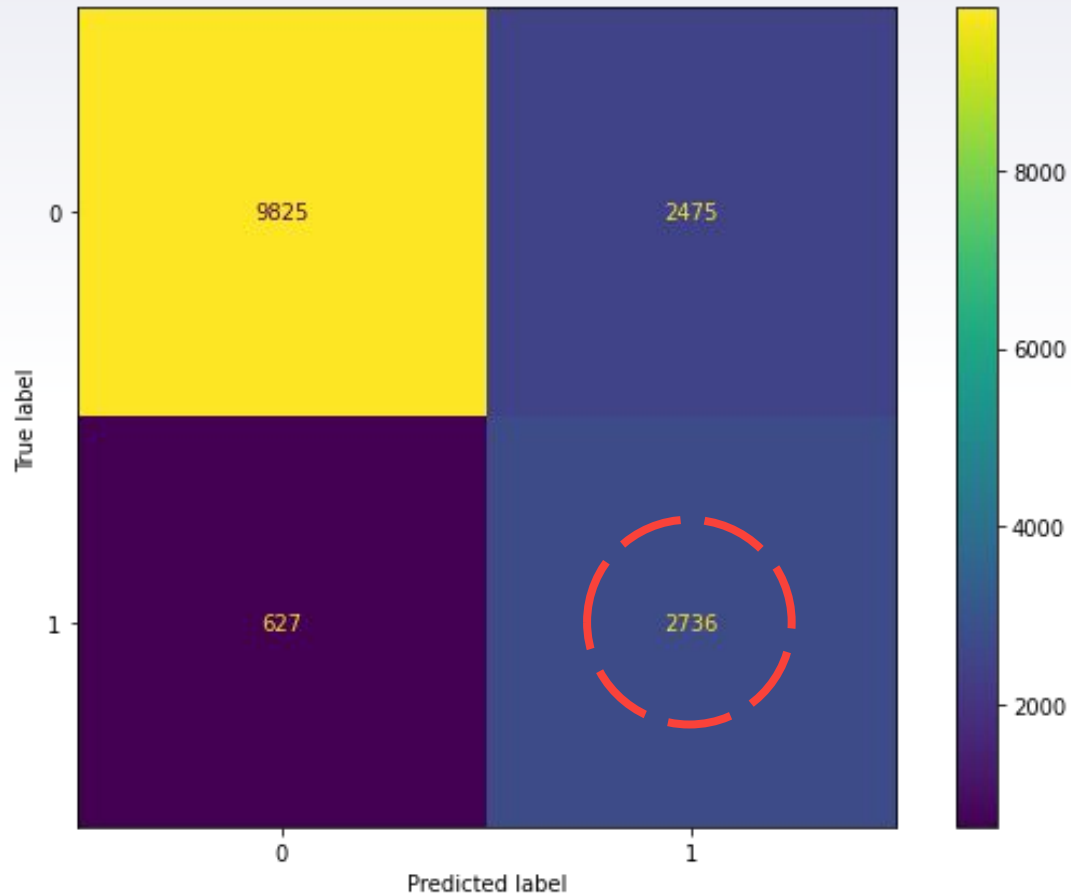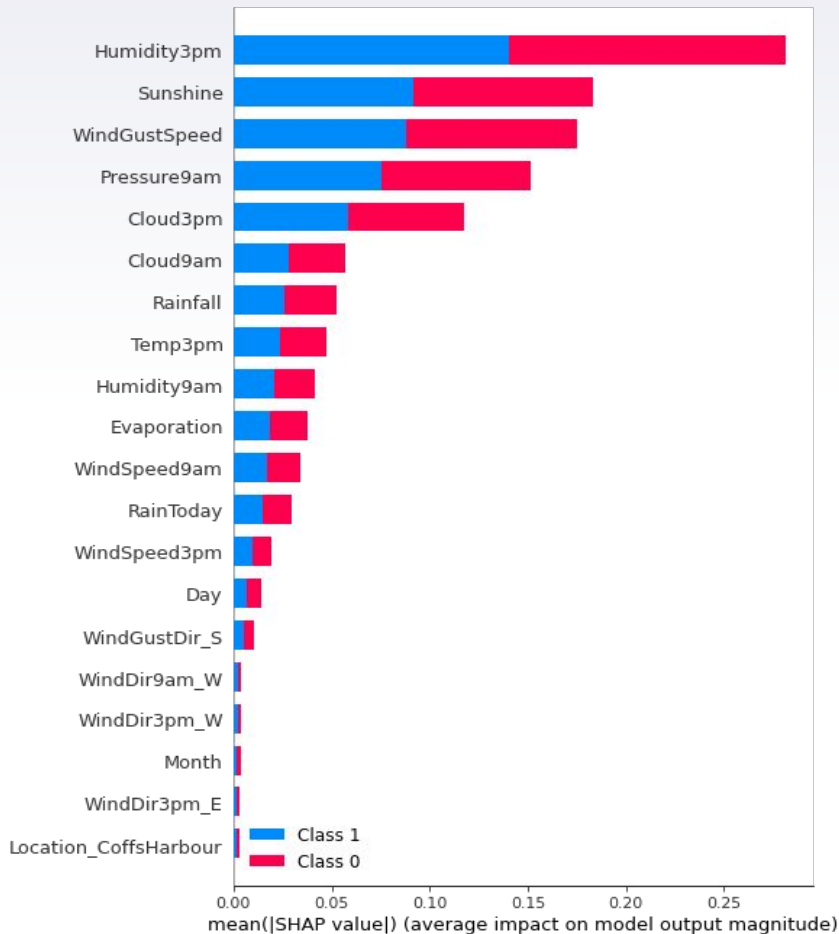
Overall, **ROC** curve and **AUC** score **before** and **after tuned** is **relatively same**. It means **before tuned**, the **model** already **have optimum performance**. But still, **hyperparameter tuning** is **essentials** part for **controlling model behavior**. It's also very important to **avoid** model **overfitting**.

# Confusion Matrix



From confusion matrix we can analyze model have **True Positive Rate** (**TPR**) by **81%**. It means, from 3363 days when our model predict will be rain, 2736 days it really rained.

# Feature Importances



Feature values in pink cause to increase the prediction. Size of the bar shows the magnitude of the feature's effect. Feature values in blue cause to decrease the prediction. Based on SHAP value, Humidity3pm and Sunshine feature have the highest effect on the prediction. It's same with heatmap correlation.

# Potential Impact in Business

We can do simulation our predictive model for business in the agricultural sector in Australia. For example, water supply costs for irrigation 1000 hectares rice field in dry season.

## Before

| | |
|---|---|
| Water needs | = 1917 Liter/month |
| Price of water | = 0.29 USD/m$^3$ |
| Cost for water | = 555,930 USD |

## After

| | |
|---|---|
| Water needs | = 1312 Liter/month |
| Price of water | = 0.29 USD/m$^3$ |
| Cost for water | = 380,460 USD |

### SAVING

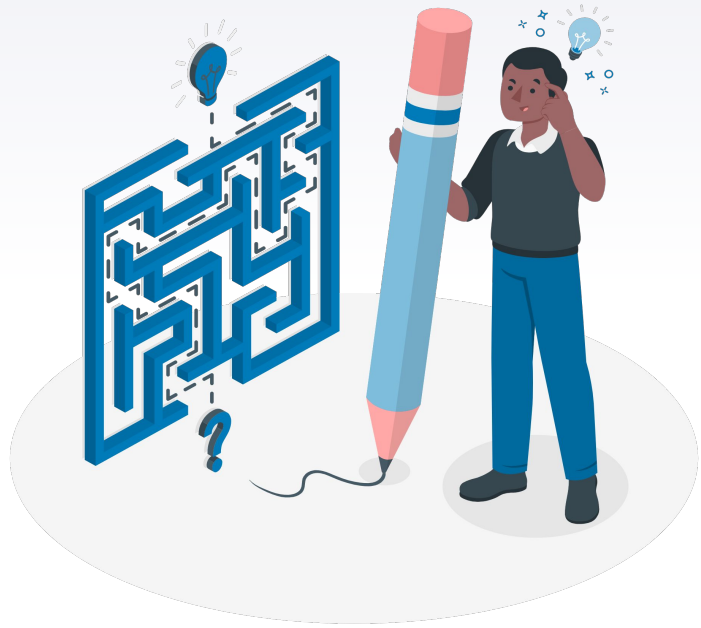## $ 175,470 or 32%

on water supply for irrigation per month

# 6 Conclusion and Recommendation

# Conclusion

- Based on Heatmap and SHAP values for feature importances, Humidity at 3pm and Sunshine has the big impact to cause rain in the next day.
- Random Forest classifier is the best model algorithm for predicting rain the next day because have the highest AUC score than other classifier algorithm.
- The location with the most rain frequency is Portland with rain season tend to happen in May until August.
- Based on simulation, model performance can help save company cost for water supply by 32%.

# Recommendation

If we want to start running the company that operating in agricultural sector, I think Australia regions is become one of the good choice since the temperature is relatively stable. Many agricultural products that can grown well based on those temperatures range. Also, I think Portland is the best location due to high rain frequency in a year, because it's very helpful for doing farming activities, such as planting and irrigation.

The best time for harvesting and make as much water stock is May until August. So, when the dry season comes, we won't too worried about lack of water and still can do activities like irrigation and farm will not easily to drought.

As a model evaluation, model have good performance and as performance trial, it would be better if the model was applied in Portland as the location with the highest rainfall frequency.

THANK YOU