

# Predicting Rainfall in Australia

...

**PROJECT GROUP #3**

Hanumasri Bollepalli (801104071)

Madhuri Pawle (801083244)

Pavanitha Jampala (801131462)

# Contents

- Introduction
- About the Dataset
- Visualizations
- Data Preprocessing
- Machine Learning Model
- Evaluating the Model
- Conclusion

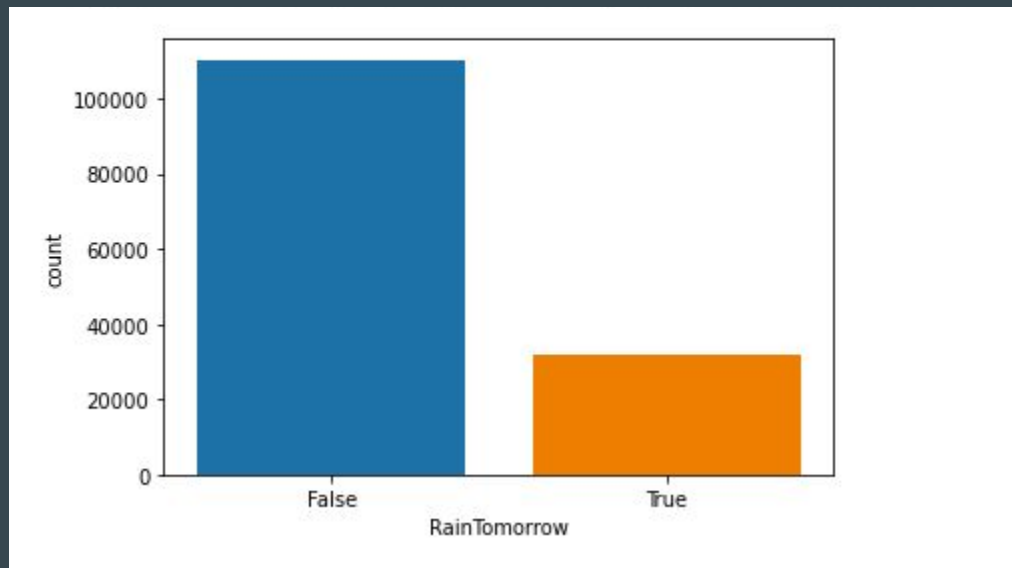
# Introduction

- The aim of the project is to predict if it will rain in Australia.
- The target audience is the general public of Australia.
- EDA (Exploratory Data Analysis) was performed to visualize various important trends.
- A logistic regression model was trained to predict rain.
- The model was evaluated using various classification metrics.

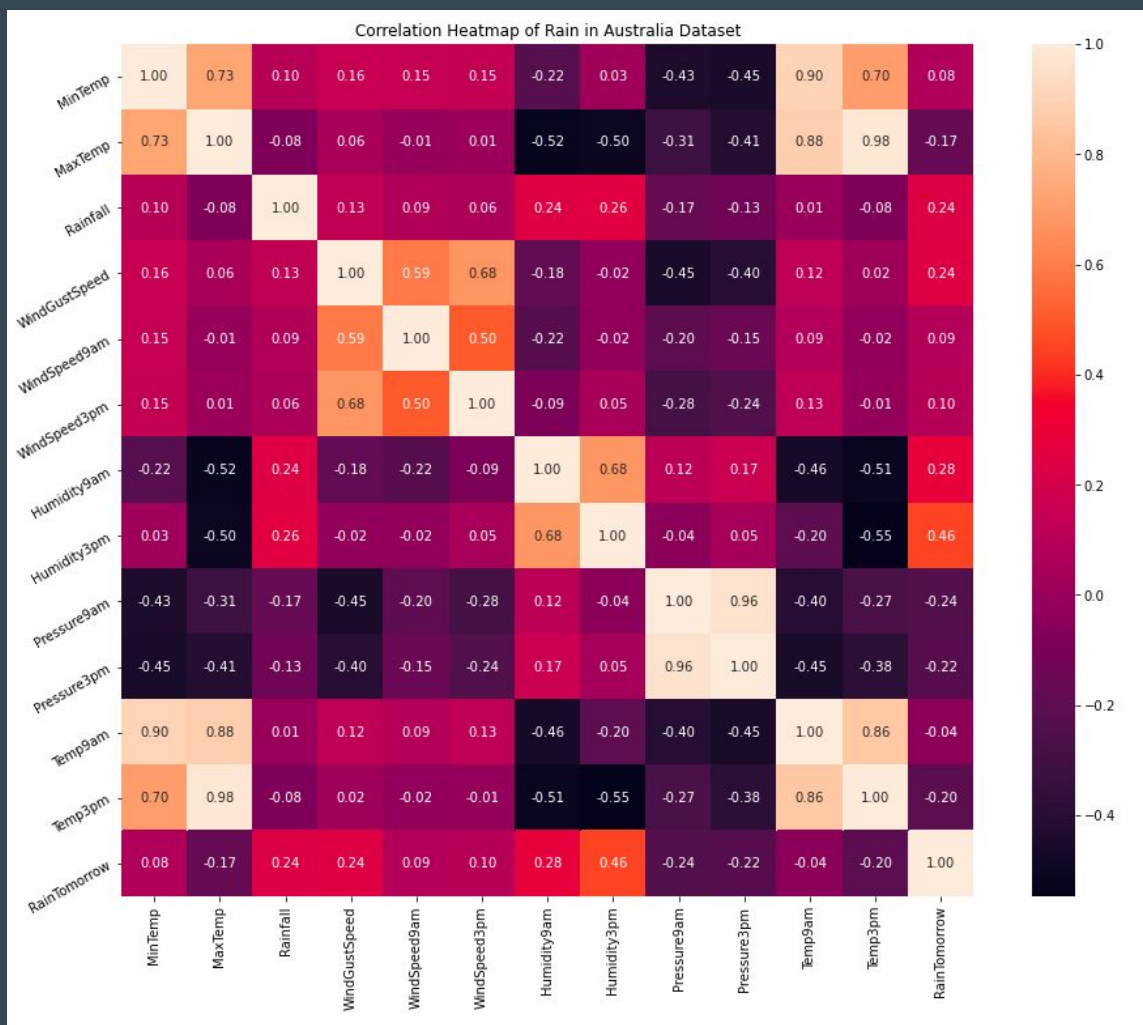
# About the Dataset

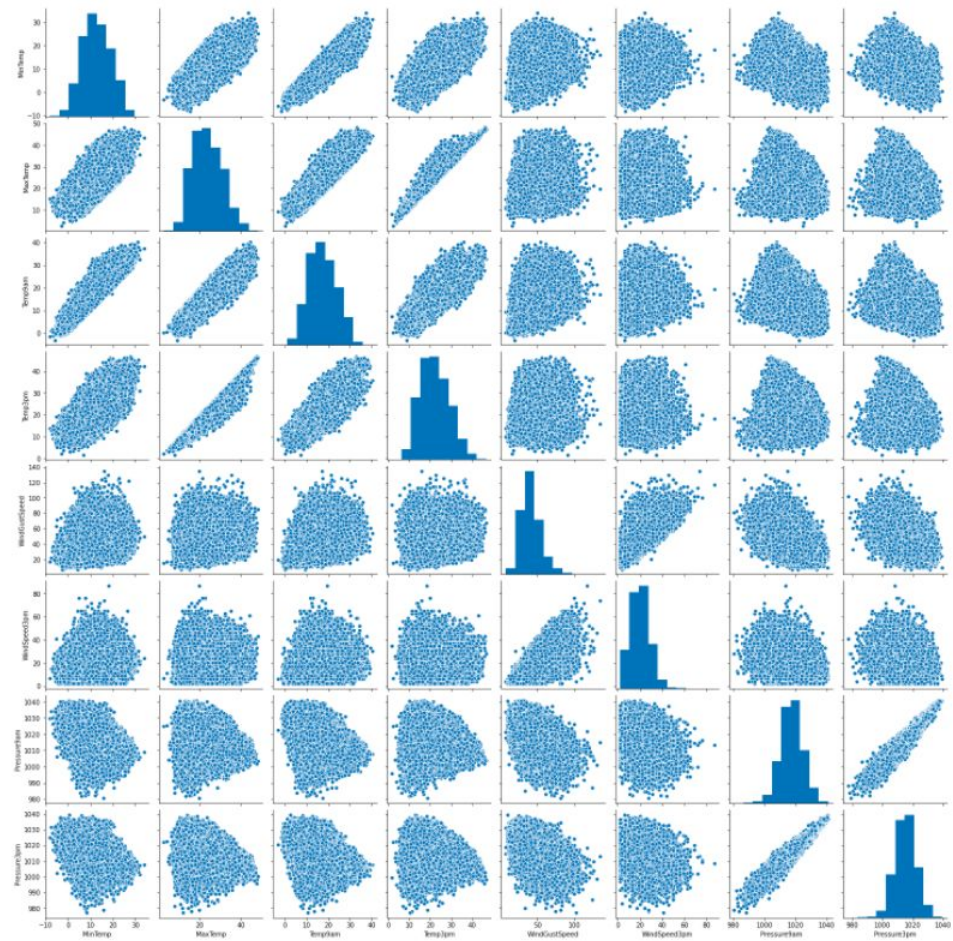
- This dataset contains about 10 years of daily weather observations from numerous Australian weather stations.
- It consists of 24 columns and 142193 rows. It has numerical, categorical and boolean data.
- It is an unbalanced dataset, as the target column RainTomorrow has 110316 "No" entries and 31877 "Yes" entries.

# Visualizations

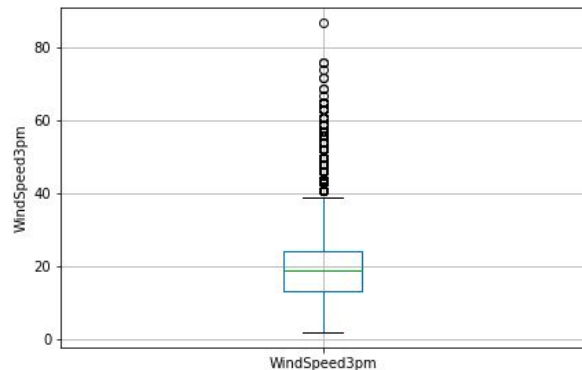
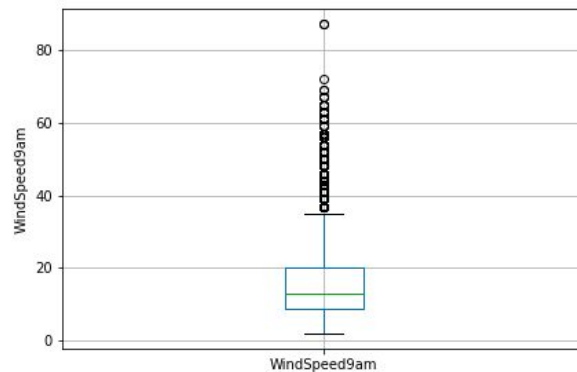
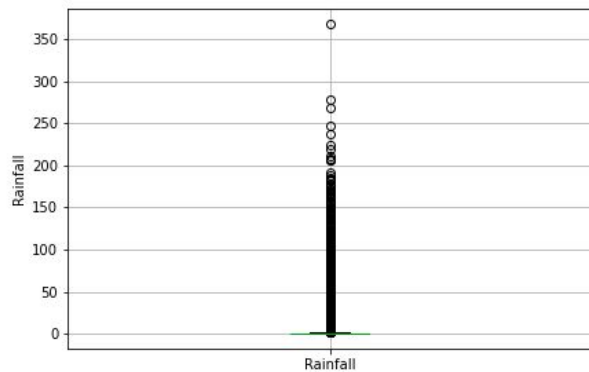


Visualizing the Target Variable





```
Text(0, 0.5, 'WindSpeed3pm')
```



Checking for Outliers



# Data Preprocessing

- The dataset is stored in BigQuery.
- The preprocessing is done using Jupyter Notebooks in AI Platform
- The datatypes of the columns are corrected.
- The columns that have less than 60% of data are dropped.
- The rows that consist of missing values in any of the columns are dropped.
- Outliers are handled using top-coding approach.
- Categorical data is encoded using One Hot Encoding as Logistic Regression cannot handle categorical data.

# Machine Learning Model

- A Logistic Regression model is trained using BigQuery ML.
- The query used to create the model is :

```
CREATE OR REPLACE MODEL
```

```
`thematic-flash-266714.australia_weather.model`
```

```
OPTIONS
```

```
( model_type="logistic_reg",  
  input_label_cols=["RainTomorrow"] ) AS
```

```
SELECT * FROM
```

```
`thematic-flash-266714.australia_weather.preprocessed`
```

## Aggregate metrics ?

Log loss ?	0.4087
ROC AUC ?	0.8432

## Score threshold

Positive class threshold ?  0.5063

Positive class true

Negative class false

Precision 0.8118

Recall 0.2501

Accuracy ? 0.8156

F1 score ? 0.3824

## Confusion matrix

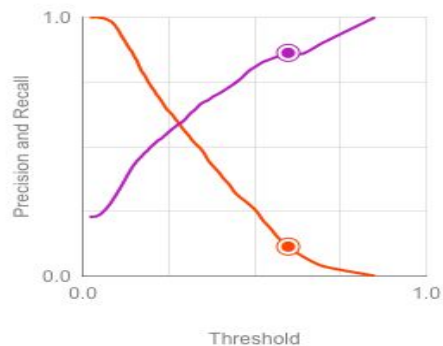
Actual labels	Predicted labels	
	true	false
true	25.01%	74.99%
false	1.71%	98.29%

Use this slider above to see which score threshold works best for your model.

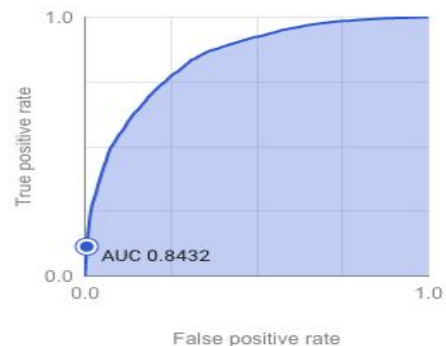
## Precision-Recall curve



## Precision and Recall vs Threshold



## ROC curve



# Evaluating the Model

- We use the following query to evaluate the model

```
SELECT * FROM ML.EVALUATE(MODEL `thematic-flash-266714.australia_weather.model`)
```

```
1 SELECT * FROM ML.EVALUATE(MODEL `thematic-flash-266714.australia_weather.model`)
```

Run

Save query

Save view

Schedule query

More

This query will process 0 B when run.

Query results

SAVE RESULTS

EXPLORE DATA

Query complete (0.0 sec elapsed, cached)

Job information

Results

JSON

Execution details

Row	precision	recall	accuracy	f1_score	log_loss	roc_auc
1	0.8041379310344827	0.2580787959274015	0.8163265306122449	0.39075067024128685	0.4087447282715583	0.8432367632367632

# Conclusion

- The accuracy of the model is 0.82 and precision is 0.80 which is good for a classification model.
- But the recall and f1-score is low, which implies that the model needs improvement.
- The performance of the model suffers due to the unbalanced dataset.

**THANK YOU!**