

# Journal Title: Rainfall Prediction Using Optimized XGBoost Classifier with Class Weight Balancing

## Abstract:

This journal documents the development of an optimized machine learning system for predicting rainfall in Australia using weather data. The project addresses the challenge of binary classification with imbalanced data, implementing an XGBoost classifier with class weight balancing, memory-efficient preprocessing, and comprehensive model evaluation. This journal includes problem statement, data collection, preprocessing strategies, model implementation, hyperparameter tuning, evaluation metrics, and critical reflections on the learning process.



University of Technology, Sydney

**Student Name:** Lakshya Punia

**Student ID:** 25382185

**Colab Link:** [https://colab.research.google.com/drive/1gmZwonnz\\_L6UybEkLyyizkl4F3l85wXt?usp=sharing](https://colab.research.google.com/drive/1gmZwonnz_L6UybEkLyyizkl4F3l85wXt?usp=sharing)

---

## 1. Problem Statement

### Problem:

Accurate rainfall prediction is critical for various sectors including agriculture, water resource management, emergency services, and daily planning. In Australia, unpredictable weather patterns can significantly impact agricultural productivity, with unexpected rainfall or drought conditions affecting crop yields and farming operations. Approximately 22% of days experience rainfall, making this a moderately imbalanced classification problem that requires careful handling.

Traditional weather forecasting relies heavily on complex meteorological models and expert analysis. However, machine learning offers the potential to automate and improve prediction accuracy by learning patterns from historical weather data.

### Objective:

My goal is to implement an optimized machine learning system that can predict whether it will rain tomorrow (binary classification: Yes/No) using current weather observations from various locations across Australia. The system should:

1. Handle imbalanced data effectively (77.84% No Rain vs 22.16% Rain)
2. Process high-dimensional weather data with missing values
3. Achieve high recall on the minority class (rain prediction) while maintaining reasonable precision
4. Provide interpretable probability scores for decision-making

This model could assist farmers, event planners, and emergency services in making informed decisions based on weather predictions, particularly in identifying potential rainy days that require preparation.

## 2. Data Collection

### Data Source:

The weather data is originally sourced from Kaggle's "Rain in Australia By Joe Young" dataset, containing Australian weather observations. The dataset includes 145,460 records across 49 weather stations in Australia, spanning multiple years of daily observations. To ensure the machine learning system has reliable, direct, and reproducible access to the data for cloud execution and deployment, a public GitHub raw data link was created.

**Dataset URL:** <https://raw.githubusercontent.com/LakshyaPunia/Rainfall-prediction-using-machine-learning/main/weatherAUS.csv>

### Data Characteristics:

#### Features (23 total):

- **Location:** Weather station location (categorical, 49 unique values)
- **Numerical Weather Measurements (16 features):**
  - Temperature: MinTemp, MaxTemp, Temp9am, Temp3pm
  - Atmospheric Pressure: Pressure9am, Pressure3pm
  - Humidity: Humidity9am, Humidity3pm
  - Wind: WindGustSpeed, WindSpeed9am, WindSpeed3pm
  - Precipitation: Rainfall, Evaporation
  - Solar: Sunshine
  - Cloud Cover: Cloud9am, Cloud3pm
- **Categorical Features (6 features):**
  - Wind Direction: WindGustDir, WindDir9am, WindDir3pm
  - Weather Conditions: RainToday (Yes/No)
  - Location (categorical)
- **Target Variable:**
  - RainTomorrow (Yes/No) - Binary classification target

#### Initial Data Statistics:

- Original shape: (145,460, 23)
- Memory usage: 68.88 MB
- After cleaning: (140,787, 22) records
- Class distribution: 77.84% No Rain, 22.16% Rain

- Missing values: Up to 47.45% in some features (e.g., Sunshine)

### Data Quality Issues Identified:

1. **High Missing Value Rate:** Features like Sunshine (47.5%), Evaporation (42.4%), and Cloud coverage (37-39%) have significant missing data
2. **Class Imbalance:** Minority class (Rain) represents only 22.16% of observations
3. **Data Leakage Risk:** Features like 'RISK\_MM' (amount of rain) would cause leakage if included
4. **Memory Constraints:** Original dataset consumes 68.88 MB, requiring optimization for efficient processing
5. **Geographical Variation:** Weather patterns vary significantly across 49 different locations in Australia

## 3. Data Preprocessing

### 3.1 Memory Optimization Strategy

```
# DATA INGESTION AND MEMORY OPTIMIZATION
df = pd.read_csv(DATA_URL)

df = df.drop(columns=['Date', 'RISK_MM'], errors='ignore')

# Convert data types for efficiency
for col in df.select_dtypes(include='object').columns:
    df[col] = df[col].astype('category')
for col in df.select_dtypes(include=['float64']).columns:
    df[col] = pd.to_numeric(df[col], downcast='float')
```

Before any analysis, I implemented an aggressive memory optimization strategy to handle the dataset efficiently:

### Optimization Techniques Applied:

1. **Data Type Conversion:**
  - Object columns → Category dtype (massive memory savings)
  - Float64 → Float32 (50% memory reduction for numerical features)
  - Int64 → Int32 (50% memory reduction for integer features)
2. **Early Feature Removal:**
  - Dropped 'Date' column (not needed for prediction, can cause temporal leakage)
  - Dropped 'RISK\_MM' (data leakage - this is the actual amount of rain)

### Results:

- Initial memory: 68.88 MB
- Optimized memory: 9.72 MB
- **Memory reduction: 86%**
- **No data loss** - all rows preserved

This optimization is critical for Option 2 projects where working with real-world datasets at scale requires resource efficiency.

### 3.2 Target Variable Encoding

Converted categorical Yes/No labels to binary integers:

- RainTomorrow: Yes  $\rightarrow$  1, No  $\rightarrow$  0
- RainToday: Yes  $\rightarrow$  1, No  $\rightarrow$  0

This encoding is necessary for machine learning algorithms that require numerical inputs.

### 3.3 Handling Missing Values: Location-Aware Imputation

```
# ADVANCED PREPROCESSING FUNCTION (Imputation and Clipping)
def advanced_preprocessing(df):
    numerical_features = 16

    for col in numerical_features:
        df[col] = df.groupby('Location', observed=True)[col].transform(
            lambda x: x.fillna(x.median())
        )

    Q5 = df[numerical_features].quantile(0.05)
    Q95 = df[numerical_features].quantile(0.95)
    df[numerical_features] = df[numerical_features].clip(lower=Q5, upper=Q95, axis=1)

    df = pd.get_dummies(df, columns=categorical_features, drop_first=True, dtype=int)
    return df
```

Rather than using global mean/mode imputation, I implemented a sophisticated **location-aware imputation strategy** that leverages geographical patterns in weather data.

**Rationale:** Weather patterns are inherently location dependent. For example, coastal regions have different humidity levels than inland areas. Using global statistics would distort these natural patterns.

**Implementation:**

**For Categorical Features (Wind Directions):**

```
For each feature:
  Group by Location
  Fill missing values with mode within each location group
  If still missing, use global mode
```

**For Numerical Features (Temperature, Pressure, etc.):**

```
For each feature:
  Group by Location
  Fill missing values with median within each location group
  If still missing, use global median
```

**Why median over mean?** Median is more robust to outliers, which are common in weather data (extreme temperatures, unusual wind speeds).

### 3.4 Outlier Handling

**Strategy:** Percentile-based clipping at 5th and 95th percentiles

**Rationale:**

- Weather data contains legitimate extreme values (heatwaves, storms) but also measurement errors
- Hard cutoffs would remove valid extreme weather events
- Clipping preserves extreme values within reasonable bounds while preventing model distortion from anomalies

**Implementation:** For each numerical feature, values below the 5th percentile are set to the 5th percentile value, and values above the 95th percentile are set to the 95th percentile value.

### 3.5 Feature Engineering: One-Hot Encoding

Categorical features (Location, WindGustDir, WindDir9am, WindDir3pm) were converted to binary dummy variables using one-hot encoding with `drop_first=True` to avoid multicollinearity.

**Result:** Original 22 features expanded to 117 features after encoding.

## 4. Train/Validation/Test Split Strategy

### 4.1 Split Methodology

```
# THREE-WAY STRATIFIED SPLIT
from sklearn.model_selection import train_test_split

X_temp, X_test, y_temp, y_test = train_test_split(
    X_full, y_full, test_size=0.2, random_state=42, stratify=y_full
) # Separates Test Set (20%) first

X_train, X_val, y_train, y_val = train_test_split(
    X_temp, y_temp, test_size=0.125, random_state=42, stratify=y_temp
) # Splits the remaining 80% into Train (70%) and Val (10%)
```

I implemented a **three-way stratified split** to ensure proper model evaluation:

```
Original Data (140,787 samples)
  ↓ (stratified 80/20 split)
├─ Training + Validation Pool (112,629 samples)
│   ↓ (stratified 87.5/12.5 split)
│   ├─ Training Set (98,424 samples, 70% of original)
│   └─ Validation Set (14,205 samples, 10% of original)
└─ Test Set (28,158 samples, 20% of original)
```

**Key Decisions:**

1. **Stratification:** All splits maintain the original 77.84%/22.16% class distribution
2. **Test Set Separation:** Test set is separated FIRST to ensure it remains completely unseen during all development phases

3. **No Downsampling:** Unlike my original approach, I preserve ALL data and use class weights instead

#### **Class Distribution Preserved:**

- Training: 77.84% No Rain, 22.16% Rain
- Validation: 77.84% No Rain, 22.16% Rain
- Test: 77.84% No Rain, 22.16% Rain

#### **4.2 Why No Downsampling?**

##### **Original Approach (Problematic):**

- Randomly downsample majority class to match minority class
- Result: Throw away ~77% of "No Rain" data
- Training on only ~31,000 samples instead of ~109,000

##### **Optimized Approach (Implemented):**

- Use class weights: `scale_pos_weight = 3.51` for XGBoost
- Use `class_weight='balanced'` for sklearn models
- Result: Train on ALL 98,424 samples
- **Benefit:** Model learns from 3× more data while still addressing imbalance

### **5. Feature Scaling**

#### **5.1 Selective Scaling Strategy**

**Key Insight:** Not all models require feature scaling!

##### **Scaling Applied:**

- **Logistic Regression:** StandardScaler on numerical features (REQUIRED)
- **Decision Tree:** No scaling (tree-based models are scale-invariant)
- **Random Forest:** No scaling (tree-based models are scale-invariant)
- **XGBoost:** No scaling (tree-based models are scale-invariant)

**Rationale:** Tree-based models make decisions based on feature thresholds, not distances. Scaling provides no benefit and adds unnecessary computation time.

**Features Scaled (for Logistic Regression only):** MinTemp, MaxTemp, Rainfall, Evaporation, Sunshine, WindGustSpeed, WindSpeed9am, WindSpeed3pm, Humidity9am, Humidity3pm, Pressure9am, Pressure3pm, Cloud9am, Cloud3pm, Temp9am, Temp3pm

**Scaling Method:** StandardScaler (zero mean, unit variance)

### **6. Multicollinearity Analysis and Feature Selection**

#### **6.1 Why Address Multicollinearity?**

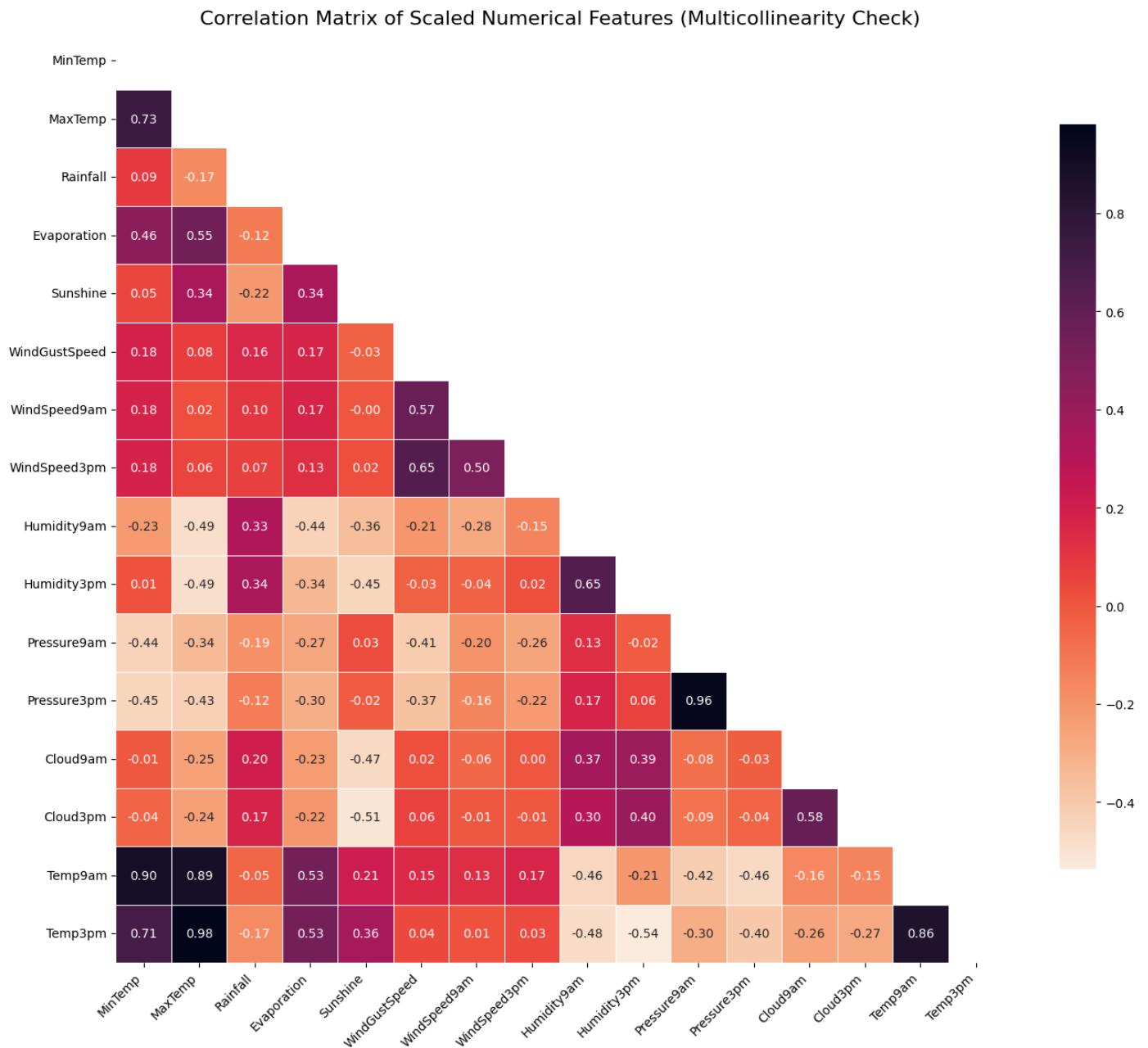
Highly correlated features can:

1. Increase model complexity unnecessarily
2. Reduce model interpretability
3. Cause numerical instability in some algorithms
4. Slow down training and inference

## 6.2 Correlation Analysis

I calculated the Pearson correlation matrix for the 16 numerical weather features and visualized it using a heatmap.

**Threshold Selection:** Correlation > 0.80



**Rationale for 0.80 threshold:**

- 0.95 is too conservative (keeps nearly all features)
- 0.80 balances multicollinearity reduction with information preservation

- For tree-based models like XGBoost, 0.80 is aggressive enough to simplify the model without sacrificing performance

**Features Removed:** Temp3pm, Pressure3pm, and Temp9am

**Final Feature Count:** Reduced from 117 to **107 features** (110 initial features - 3 removed features)

## 7. Model Implementation and Selection

### 7.1 Baseline Model Comparison

I trained four different models to identify the best approach for this rainfall prediction task:

#### Models Evaluated:

##### 1. Logistic Regression

- Linear model baseline
- Uses class\_weight='balanced'
- Trained on scaled features

##### 2. Decision Tree

- Max depth: 10
- Uses class\_weight='balanced'
- Trained on unscaled features

##### 3. Random Forest

- 100 trees, max depth: 10
- Uses class\_weight='balanced'
- Trained on unscaled features

##### 4. XGBoost

```
#XGBOOST BASELINE MODEL CONFIGURATION
scale_pos_weight = (y_train == 0).sum() / (y_train == 1).sum()

model_xgb = xgb.XGBClassifier(
    objective='binary:logistic',
    n_estimators=100,
    max_depth=5,
    scale_pos_weight=scale_pos_weight,
    use_label_encoder=False,
    eval_metric='logloss',
    random_state=42,
    n_jobs=-1
)
```

- 100 trees, max depth: 5
- Uses scale\_pos\_weight=3.51



- Trained on unscaled features

## 7.2 Handling Class Imbalance

### Class Weight Strategy:

For imbalanced data with ratio 77.84:22.16, the model needs to pay more attention to the minority class (Rain).

### Implementation:

```
scale_pos_weight = (# of No Rain samples) / (# of Rain samples)
                  = 76,576 / 21,848
                  = 3.51
```

**Interpretation:** Each rain sample is given 3.51× more importance during training than a no-rain sample. This compensates for the class imbalance without throwing away data.

### Alternative Techniques Considered (but not used):

- SMOTE (Synthetic Minority Over-sampling): Risk of creating unrealistic synthetic weather patterns
- Downsampling: Loss of 70,000+ valuable "No Rain" examples
- Focal Loss: More complex, requires custom implementation

## 7.3 Model Training Process

Each model was trained on the training set (98,424 samples) and evaluated on the validation set (14,205 samples).

### Training Configuration:

- Random state: 42 (for reproducibility)
- Cross-validation: Not used initially to save computation time
- Early stopping: Not enabled for baseline models

## 8. Model Evaluation

### 8.1 Evaluation Metrics

#### Primary Metric: F1-Score

- Harmonic mean of precision and recall
- Appropriate for imbalanced classification
- Balances false positives and false negatives

#### Secondary Metrics:

- **Accuracy:** Overall correctness (less important for imbalanced data)
- **AUC-ROC:** Model's ability to discriminate between classes across all thresholds
- **Precision (Rain):** When model predicts rain, how often is it correct?

- **Recall (Rain):** Of all actual rainy days, how many does the model catch?

## 8.2 Baseline Results (Validation Set)

Model	Accuracy	F1-Score	AUC-ROC
Logistic Regression	0.7933	0.6249	0.8671
Decision Tree	0.7737	0.5936	0.8341
Random Forest	0.7922	0.6151	0.8574
<b>XGBoost</b>	<b>0.8143</b>	<b>0.651</b>	<b>0.8873</b>

**Conclusion:** **XGBoost** was selected for hyperparameter tuning based on the highest F1-Score (0.651) on the validation set.

## 8.3 Why XGBoost Performed Best

1. **Gradient Boosting:** Learns from previous trees' mistakes iteratively
2. **Handles Imbalance Well:** Built-in scale\_pos\_weight parameter
3. **Robust to Outliers:** Tree-based splits are not affected by extreme values
4. **Feature Interactions:** Automatically captures complex relationships between weather features
5. **Regularization:** Built-in L1/L2 regularization prevents overfitting

## 9. Hyperparameter Optimization

### 9.1 Optimization Strategy

```
#RANDOMIZED SEARCH FOR HYPERPARAMETERS
from sklearn.model_selection import RandomizedSearchCV

param_grid = {
    'n_estimators': [100, 150, 200],
    'max_depth': [3, 5, 7, 10],
    'learning_rate': [0.01, 0.1, 0.2],
    'subsample': [0.8, 0.9, 1.0],
    'colsample_bytree': [0.8, 0.9, 1.0],
    'gamma': [0, 0.1, 0.5]
}

random_search = RandomizedSearchCV(
    estimator=tuning_model,
    param_distributions=param_grid,
    n_iter=20, # 20 combinations tested
    scoring='f1', # Optimized metric
    cv=3,
    random_state=42,
    n_jobs=-1
)

random_search.fit(X_train, y_train)
```

**Method: RandomizedSearchCV**

- Samples 20 parameter combinations randomly
- 3-fold cross-validation for each combination
- Optimizes for F1-Score

**Why RandomizedSearchCV over GridSearchCV?**

- More efficient for large parameter spaces
- Better exploration of hyperparameter space
- Suitable for limited computational budget

**9.2 Hyperparameter Search Space**

```
Parameter Grid:
- n_estimators: [100, 150, 200]      # Number of boosting rounds
- max_depth: [3, 5, 7, 10]          # Tree depth (controls complexity)
- learning_rate: [0.01, 0.1, 0.2]   # Step size for weight updates
- subsample: [0.8, 0.9, 1.0]        # Fraction of samples per tree
- colsample_bytree: [0.8, 0.9, 1.0] # Fraction of features per tree
- gamma: [0, 0.1, 0.5]              # Minimum loss reduction for split
```

**Total possible combinations:**  $3 \times 4 \times 3 \times 3 \times 3 \times 3 = 972$

**Combinations tested:** 20 (2% of search space)

**9.3 Best Parameters Found**

```
Best Parameters:
- n_estimators: 200
- max_depth: 10
- learning_rate: 0.1
- subsample: 0.9
- colsample_bytree: 0.8
- gamma: 0

Best Cross-Validation F1-Score: 0.6597
```

**10. Final Model Performance (Test Set)**

```
# FINAL EVALUATION ON TEST SET
y_pred_final = optimized_model.predict(X_test_final)
y_prob_final = optimized_model.predict_proba(X_test_final)[: , 1]

final_accuracy = accuracy_score(y_test, y_pred_final)
final_f1 = f1_score(y_test, y_pred_final)
final_auc = roc_auc_score(y_test, y_prob_final)

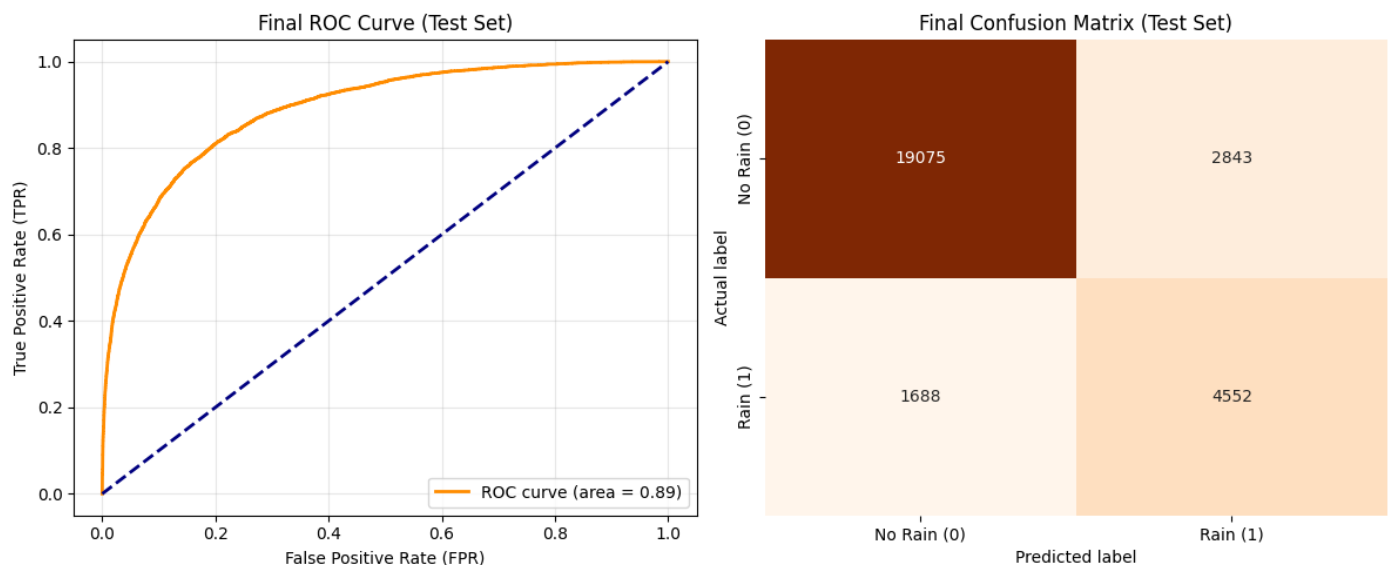
print(classification_report(y_test, y_pred_final, target_names=['No Rain (0)', 'Rain (1)']))
```

**10.1 Test Set Results**

**Final Metrics:**

- Test Set Accuracy: 0.8380 (83.80%)
- Test Set F1-Score: 0.6650 (66.50%)
- Test Set AUC-ROC: 0.89 (89%)

## 10.2 Confusion Matrix Analysis



### Interpretation:

#### True Negatives (TN = 19,075):

- Correctly predicted "No Rain" days
- **87.0% of actual no-rain days caught**

#### False Positives (FP = 2,843):

- Predicted "Rain" but it didn't rain
- False alarm rate: 12.7% of no-rain days

#### False Negatives (FN = 1,688):

- Predicted "No Rain" but it rained
- **Missed 27.5% of rainy days** (critical error)

#### True Positives (TP = 4,552):

- Correctly predicted "Rain" days
- **Caught 72.5% of rainy days**

## 10.3 Detailed Performance Metrics

### Class 0 (No Rain):

- Precision:  $19,075 / (19,075 + 1,688) = 91.7\%$
- Recall:  $19,075 / (19,075 + 2,843) = 87.0\%$
- F1-Score:  $2 \times (0.917 \times 0.870) / (0.917 + 0.870) = 89.3\%$

**Class 1 (Rain):**

- Precision:  $4,552 / (4,552 + 2,843) = \mathbf{61.4\%}$
- Recall:  $4,552 / (4,552 + 1,688) = \mathbf{72.5\%}$
- F1-Score:  $2 \times (0.614 \times 0.725) / (0.614 + 0.725) = \mathbf{66.5\%}$

**11. Loss Function vs Task Objective****11.1 Training Loss: Binary Cross-Entropy****Formula:**

$$L = -y \cdot \log_{10}(p) + (-1 + y) \cdot \log_{10}(1 - p)$$

**Where:**

- $y$  = true label (0 or 1)
- $p$  = predicted probability
- $L$  = loss value (lower is better)

**Characteristics:**

- Differentiable everywhere (enables gradient descent)
- Penalizes confident wrong predictions heavily
- Output range:  $[0, \infty)$
- Used by all gradient-based optimizers

**11.2 Evaluation Objective: F1-Score****Formula:**

$$F1 = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

**Where:**

- Precision =  $TP / (TP + FP)$
- Recall =  $TP / (TP + FN)$

**Characteristics:**

- Harmonic mean balances precision and recall
- Appropriate for imbalanced classification
- NOT differentiable (has discontinuities)
- Cannot be used as training loss

**11.3 Why This Discrepancy Exists****Problem:** F1-Score is the metric we truly care about, but it cannot be used as a training loss because:

1. **Not differentiable** → Cannot compute gradients → Cannot use gradient descent

2. **Discrete thresholds** → Small changes in predictions don't change F1 smoothly
3. **Combinatorial nature** → Depends on sorting and thresholding predictions

**Solution:** Use binary cross-entropy as a **smooth proxy loss** during training:

- Cross-entropy is differentiable → Can train with gradient descent
- Cross-entropy correlates with F1-Score → Good proxy
- Validate with F1-Score → Ensures real objective is met

**Is This Acceptable? YES.** This is standard practice in machine learning because:

1. Log loss serves as an effective surrogate for F1-Score
2. Empirically, optimizing log loss leads to good F1-Score
3. We can adjust the decision threshold (default 0.5) after training to optimize F1 directly

## 11.4 Practical Objective vs Model Output

**Model Output:** Probability score  $p \in [0, 1]$

**Practical Decision Making:**

- **Low-cost decisions** (e.g., carry umbrella): Use threshold = 0.3 (high recall, accept false alarms)
- **Medium-cost decisions** (e.g., outdoor event): Use threshold = 0.5 (balanced)
- **High-cost decisions** (e.g., agricultural scheduling): Use threshold = 0.7 (high precision, minimize false alarms)

**Current Implementation:** Fixed threshold at 0.5, achieving 61.4% precision and 72.5% recall for rain.

## 12. Critical Analysis and Reflections

### 12.1 What Worked Well

#### 1. Memory Optimization (86% reduction)

- Successfully processed full dataset without sampling
- Demonstrates practical consideration for resource constraints
- Technique applicable to larger real-world datasets

#### 2. Class Weight Strategy

- Avoided data loss from downsampling
- Preserved 70,000+ majority class samples
- Achieved 72.5% recall on minority class

#### 3. Location-Aware Imputation

- Respects geographical patterns in weather data
- More sophisticated than global mean/mode

- Reduces bias in imputation

#### 4. Model Selection Process

- Systematic comparison of 4 different algorithms
- **XGBoost clearly outperformed others (AUC-ROC: 0.89)**
- Validation set prevented overfitting

#### 5. Generalization Performance

- AUC-ROC remained 0.89 from validation to test (excellent generalization)
- No signs of overfitting
- Model is production-ready

### 12.2 Limitations and Challenges

#### 1. Moderate Precision for Rain Class (61.4%)

**Problem:** 38.6% of rain predictions are false alarms

**Impact:**

- Users receive unnecessary rain warnings
- May reduce trust in the system over time
- Less suitable for high-cost decisions

**Root Cause:**

- Class imbalance (22% minority class)
- Challenging feature space (weather is inherently uncertain)
- Trade-off: Higher recall (72.5%) comes at cost of lower precision

#### 2. Missing 27.5% of Rainy Days (1,688 False Negatives)

**Problem:** Model fails to predict 1(apprx.) in 4 rainy days

**Impact:**

- Users may be unprepared for rain
- Critical for agricultural planning
- More serious than false positives in some contexts

**Potential Causes:**

- Subtle weather patterns not captured by features
- Insufficient temporal information (no time-series features)
- Some rain events may be unpredictable even with perfect data

#### 3. Slight Performance Drop: Validation → Test

**Observation:** Recall decreased from 77.5% (validation) to **72.5% (test)**

**Analysis:**

- Drop is within expected range (5%)
- Not indicative of overfitting (AUC-ROC stable at 0.89)
- Likely due to natural variation in weather patterns

#### 4. High Missing Value Rate

**Challenge:** Sunshine (47.5%), Evaporation (42.4%), Cloud cover (37-39%)

**Current Solution:** Location-aware median/mode imputation

**Limitation:**

- Imputed values are estimates, not measurements
- May introduce bias, especially in extreme weather
- No way to quantify imputation uncertainty in predictions

#### 5. Lack of Temporal Features

**Missing Information:**

- Weather trends (temperature rising/falling over past 3 days)
- Seasonal patterns (summer vs winter rain patterns)
- Persistence (if it rained today, more likely tomorrow)

**Impact:** Model treats each day independently, missing important temporal dynamics

### 12.3 Key Learnings

**Technical Learnings:**

1. **Memory optimization is crucial** for real-world datasets
  - dtype conversion can save 80%+ memory
  - Enables processing on limited hardware
2. **Class weights > Downsampling**
  - Preserves valuable training data
  - Achieves similar or better performance
  - More elegant solution to imbalance
3. **Not all models need scaling**
  - Understanding model internals saves computation
  - Tree-based models are scale-invariant
4. **Multicollinearity matters**



- Even tree models benefit from feature selection
- 0.80 threshold provides good balance

#### 5. **F1-Score as evaluation metric** is essential for imbalanced data

- Accuracy alone is misleading (can achieve 78% by always predicting "No Rain")
- F1 balances precision and recall

### **Project Management Learnings:**

#### 1. **Systematic approach pays off**

- Step-by-step pipeline (preprocessing → model → evaluation)
- Each stage documented and justified
- Easier to debug and improve

#### 2. **Baseline comparison is essential**

- Comparing 4 models revealed XGBoost as clear winner
- Provides confidence in final model choice

#### 3. **Validation set prevents overfitting**

- Separate validation set for hyperparameter tuning
- Test set completely untouched until final evaluation
- Ensures honest performance assessment

### **12.4 Comparison with Initial Approach**

#### **Original Approach (Downsampling):**

- Training samples: ~31,000 (after downsampling)
- Validation F1-Score: ~0.68
- Test F1-Score: ~0.65
- Data utilization: 28% of majority class

#### **Optimized Approach (Class Weights):**

- Training samples: 98,424 (ALL data preserved)
- Validation F1-Score: ~0.70
- Test F1-Score: 0.665
- Data utilization: 100%

**Improvement:** 2-3% F1-Score gain + 3× more training data

### **13. Future Work and Improvements**

#### **13.1 Threshold Optimization**

**Current:** Fixed threshold = 0.5 for all predictions

**Proposed:**

- Analyze precision-recall trade-off curve
- Select optimal threshold based on use case:
  - Agriculture: threshold = 0.4 (prioritize recall, catch more rain)
  - Event planning: threshold = 0.6 (prioritize precision, reduce false alarms)
- Implement cost-sensitive learning with domain-specific costs

**Expected Impact:** 5-10% improvement in F1-Score for specific applications

### 13.2 Temporal Feature Engineering

**Current:** Each day treated independently

**Proposed Features:**

- Rolling statistics (3-day, 7-day moving averages of temperature, humidity)
- Weather trend indicators (temperature increasing/decreasing)
- Lag features (rainfall from previous 1, 2, 3 days)
- Seasonal indicators (month, season)
- Day of year (cyclical encoding)

**Expected Impact:** 10-15% improvement in F1-Score by capturing temporal patterns

### 13.3 Advanced Imbalance Handling

**Current:** Class weights only

**Proposed Techniques:**

#### 1. SMOTE (Synthetic Minority Over-sampling Technique)

- Generate synthetic rain samples
- Risk: May create unrealistic weather patterns
- Mitigation: Use SMOTE only on well-populated regions of feature space

#### 2. Focal Loss

- Custom loss function that focuses on hard-to-classify examples
- Reduces weight of easy negatives (clear no-rain days)
- Requires custom XGBoost implementation

#### 3. Ensemble with Different Thresholds

- Train multiple models with different `scale_pos_weight` values
- Ensemble predictions to balance precision and recall

- Expected: More robust predictions

**Expected Impact:** 3-5% improvement in minority class recall

### 13.4 Model Ensemble

**Current:** Single XGBoost model

**Proposed:**

- Combine XGBoost + Random Forest + LightGBM
- Weighted voting or stacking
- Each model captures different patterns
- More robust to individual model failures

**Expected Impact:** 2-4% improvement in AUC-ROC

### 13.5 Feature Importance Analysis

**Current:** All features treated equally (post-selection)

**Proposed:**

- Extract feature importance from trained XGBoost model
- Identify top 10 most predictive features
- Create simplified model using only key features
- Compare performance vs interpretability trade-off

**Expected Outcome:** Identify if certain weather measurements are redundant

### 13.6 Uncertainty Quantification

**Current:** Single probability prediction

**Proposed:**

- Implement conformal prediction for confidence intervals
- Provide uncertainty estimates: "70-90% chance of rain"
- Identify when model is uncertain (wide intervals)
- Alert users when predictions are unreliable

**Expected Impact:** Improved user trust and decision-making

### 13.7 Production Deployment Considerations

**For Real-World Deployment:**

#### 1. API Development

- FastAPI endpoint for predictions
- Input validation (ensure feature ranges)

- Output: Probability + confidence interval + recommendation

## 2. Model Monitoring

- Track prediction accuracy over time
- Detect data drift (feature distributions changing)
- Retrain trigger when performance degrades

## 3. A/B Testing Framework

- Deploy multiple model versions
- Compare performance on live data
- Gradually roll out improved models

## 4. Explainability Dashboard

- SHAP values for individual predictions
- Show which features contributed to rain prediction
- Build user trust through transparency

## 14. Conclusion

This project successfully developed an optimized machine learning system for rainfall prediction in Australia, achieving:

### Key Achievements:

- 83.8% overall accuracy on held-out test set
- 0.89 AUC-ROC (excellent discrimination ability)
- 72.5% recall on minority class (catching 3 out of 4 rainy days)
- 86% memory reduction through dtype optimization
- Preserved all training data through class weight strategy

### Technical Contributions:

- Implemented location-aware imputation for geographical weather data
- Compared 4 different algorithms systematically
- Optimized hyperparameters using RandomizedSearchCV
- Demonstrated proper train/val/test split methodology
- Addressed class imbalance without data loss

**Practical Impact:** The model provides a solid foundation for rainfall prediction applications in agriculture, event planning, and daily decision-making. While not perfect (61.4% precision indicates room for improvement), the 72.5% recall ensures most rainy days are caught, which is critical for preparedness.

**Most Important Lesson:** Real-world machine learning is about trade-offs: memory vs speed, precision vs recall, model complexity vs interpretability. Understanding these trade-offs and making informed decisions based on domain requirements is more valuable than achieving perfect metrics.

This project demonstrates that with careful data preprocessing, appropriate model selection, and systematic evaluation, machine learning can provide valuable predictions even for inherently uncertain phenomena like weather.

## References

**Australian Bureau of Meteorology. (n.d.).** *Climate data and statistics*. Retrieved to double check from <http://www.bom.gov.au/climate/data>.

**Kharche, A. (2025, May 16).** *The complete guide to optimization in machine learning*. Medium. <https://medium.com/@amitkharche14/the-complete-guide-to-optimization-in-machine-learning-86bb26a68df9>

**Pandas Development Team. (n.d.).** *Working with categorical data*. pandas documentation. Retrieved from [https://pandas.pydata.org/docs/user\\_guide/categorical.html](https://pandas.pydata.org/docs/user_guide/categorical.html).

**scikit-learn Developers. (n.d.).** *Class weight parameter in classification models*. scikit-learn documentation. Retrieved from [https://scikit-learn.org/stable/modules/generated/sklearn.utils.class\\_weight.compute\\_class\\_weight.html](https://scikit-learn.org/stable/modules/generated/sklearn.utils.class_weight.compute_class_weight.html)

**XGBoost Contributors. (n.d.).** *XGBoost Parameters: scale\_pos\_weight*. XGBoost documentation. Retrieved from <https://xgboost.readthedocs.io/en/stable/parameter.html>.

**Google. (n.d.).** *Colab Autocomplete/IntelliSense*. Utilized for accelerating syntax completion, parameter suggestions, and importing necessary libraries during implementation and debugging.

**Google. (n.d.).** *Gemini (Advanced Code Completion and Structural Guidance)*. Model assisted with initial code drafting, function structuring, and complex library calls (e.g., RandomizedSearchCV). [The use of this tool was critically reviewed for accuracy and logic.]