Real-World Applications of Machine Learning: Evaluating Hyperparameters
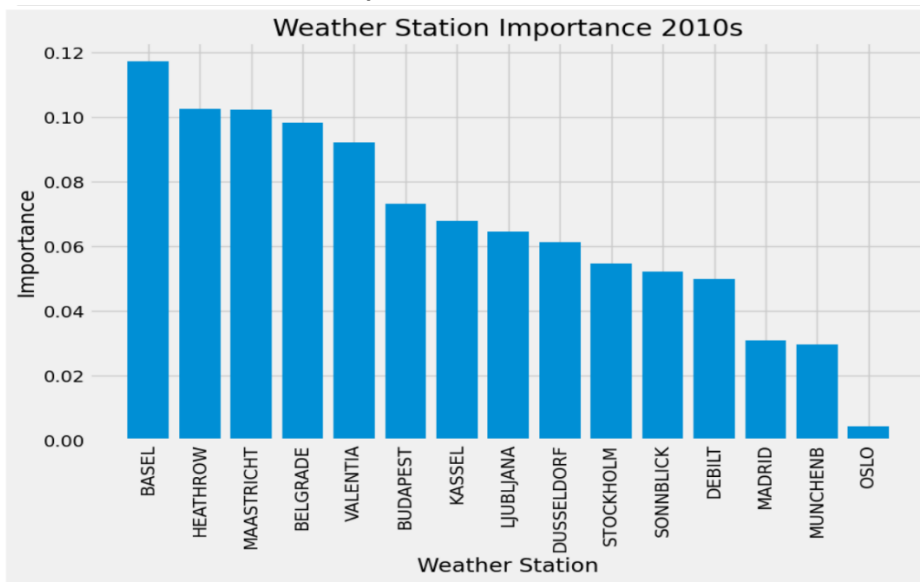
**Part 1 – Random Forest Model**

Before optimization, the Random Forest model achieved an accuracy of **59.1%** when predicting weather conditions across all stations from 2010 to 2019. After hyperparameter optimization, the accuracy remained almost the same, showing a marginal drop to **58.7%**, indicating that the optimization primarily affected the model's structure rather than its predictive performance.

For the single station, *Maastricht*, which included data spanning all years, the model achieved a perfect accuracy of **100%**, both before and after optimization. This suggests that the dataset for Maastricht is inherently separable, making it relatively easier for the model to make precise predictions.

Optimization brought notable structural improvements. The optimized model's decision trees were simpler, more balanced, and less prone to overfitting, focusing on key features like temp_max and precipitation. These features consistently emerged as the most influential predictors across stations, both before and after optimization. Other variables, such as sunshine and humidity, played secondary roles with slight adjustments to their weights.

In terms of station-level importance, the optimized model redistributed focus among the stations. While Düsseldorf and Maastricht retained their high significance, smaller shifts reflected a more balanced reliance on features across different locations.

**All Weather Stations Before Optimization**

**All Weather Stations After Optimization**



Weather Station Importance 2010s — Optimized

**All Weather Top 15 Feature Importances After Optimized**



Top 15 Feature Importances (Optimized RF)

## Maastricht Before Optimization (All Years)



Model for MAASTRICHT Accuracy: 1.0

Random Forest Tree for MAASTRICHT

MAASTRICHT_global_radiation <= 1.795
gini = 0.329
samples = 10915
value = [13646, 3566]

MAASTRICHT_temp_max <= 17.95
gini = 0.132
samples = 8172
value = [11987, 913]

MAASTRICHT_temp_max <= 17.95
gini = 0.473
samples = 2743
value = [1659, 2653]

gini = 0.0
samples = 6678
value = [10526, 0]

MAASTRICHT_precipitation <= 0.005
gini = 0.473
samples = 1494
value = [1461, 913]

gini = 0.0
samples = 549
value = [875, 0]

MAASTRICHT_humidity <= 0.755
gini = 0.352
samples = 2194
value = [784, 2653]

MAASTRICHT_cloud_cover <= 6.5
gini = 0.157
samples = 630
value = [86, 913]

gini = 0.0
samples = 864
value = [1375, 0]

MAASTRICHT_pressure <= 1.014
gini = 0.267
samples = 1777
value = [440, 2334]

MAASTRICHT_precipitation <= 0.005
gini = 0.499
samples = 417
value = [344.0, 319.0]

(...)   (...)   (...)   (...)   (...)   (...)

MAASTRICHT Feature Importances

**Maastricht After Optimization (All Years)**

Optimized RF - Maastricht (Entire Timeline) - Tree #0

```
                    MAASTRICHT_temp_max <= 17.95
                         gini = 0.329
                        samples = 10915
                     value = [13646, 3566]
                    /                        \
          gini = 0.0               MAASTRICHT_precipitation <= 0.005
       samples = 7227                     gini = 0.474
      value = [11401, 0]                 samples = 3688
                                      value = [2245, 3566]
                                     /                    \
                  MAASTRICHT_sunshine <= 0.95          gini = 0.0
                         gini = 0.046               samples = 1351
                        samples = 2337             value = [2159, 0]
                      value = [86, 3566]
                     /                \
           gini = 0.0              gini = 0.0
         samples = 54            samples = 2283
        value = [86, 0]         value = [0, 3566]
```
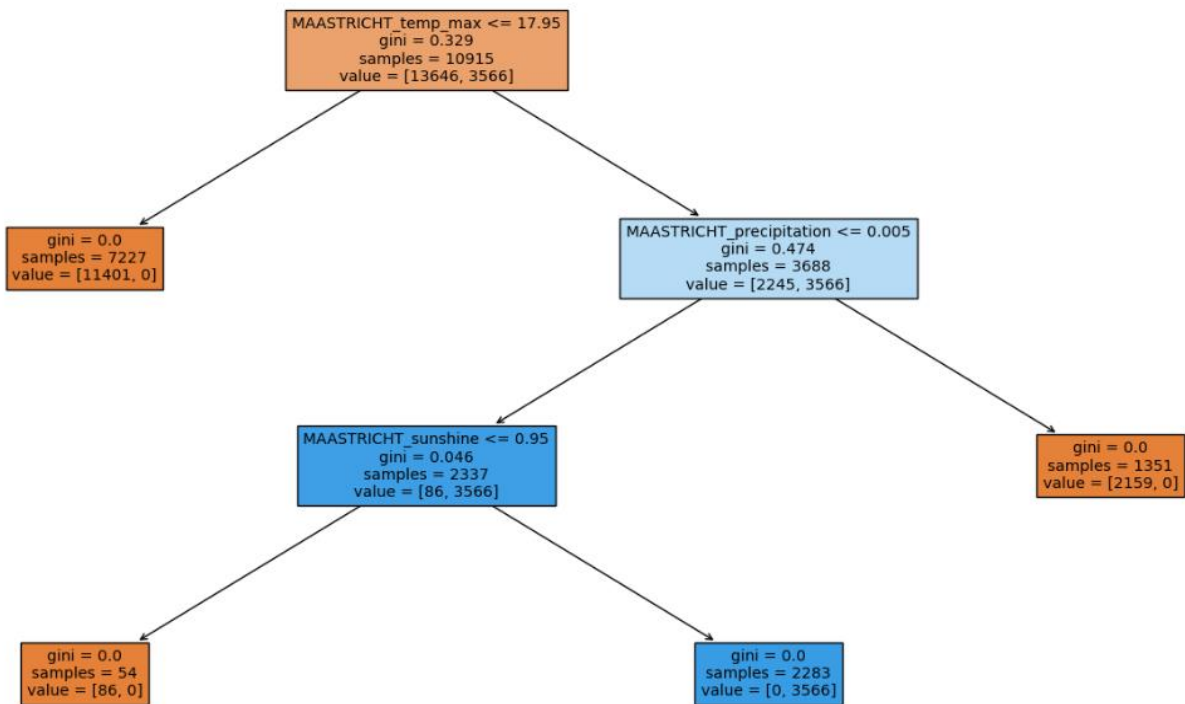


Maastricht Feature Importances - Entire Timeline (Optimized)

## Part 2 – Deep Learning (CNN Model)

The CNN model showed a clear improvement after hyperparameter optimization, with test accuracy rising from **73.9%** to **81.1%** and a corresponding decrease in loss from **0.7537** to **0.6609**. This improvement demonstrates the value of fine-tuning parameters like the number of filters, kernel size, dropout rate, and learning rate.

The optimized CNN had the following configuration:

- **Filters**: 83
- **Kernel Size**: 1
- **Dropout**: 0.48
- **Learning Rate**: 0.016
- **Batch Size**: 116
- **Epochs**: 16

This setup allowed the model to better capture non-linear relationships and patterns in the data, outperforming the Random Forest model in predictive capability. However, the CNN model is computationally more demanding and less interpretable than Random Forest, which limits its applicability for scenarios requiring explainability.

### Confusion Matrix (Stations)

| True Label \ Predicted | BASEL | BELGRADE | BUDAPEST | DEBILT | DUSSELDORF | HEATHROW | KASSEL | LJUBLJANA | MAASTRICHT | MADRID | MUNCHENB | OSLO | SONNBLICK | STOCKHOLM | VALENTIA |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| BASEL | 2800 | 138 | 6 | 0 | 0 | 0 | 0 | 0 | 0 | 11 | 0 | 0 | 0 | 0 | |
| BELGRADE | 353 | 519 | 2 | 0 | 0 | 0 | 0 | 0 | 0 | 5 | 0 | 0 | 0 | 0 | |
| BUDAPEST | 47 | 25 | 82 | 5 | 0 | 0 | 0 | 0 | 0 | 3 | 0 | 0 | 0 | 0 | |
| DEBILT | 22 | 2 | 19 | 15 | 0 | 3 | 0 | 0 | 0 | 3 | 0 | 0 | 0 | 0 | |
| DUSSELDORF | 11 | 0 | 2 | 3 | 0 | 4 | 0 | 0 | 0 | 5 | 0 | 0 | 0 | 0 | |
| HEATHROW | 18 | 1 | 2 | 0 | 0 | 13 | 0 | 5 | 0 | 28 | 0 | 0 | 0 | 0 | |
| KASSEL | 4 | 1 | 1 | 0 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | |
| LJUBLJANA | 19 | 2 | 3 | 0 | 0 | 0 | 0 | 4 | 0 | 18 | 0 | 0 | 0 | 0 | |
| MAASTRICHT | 4 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 2 | 0 | 0 | 0 | 0 | |
| MADRID | 47 | 13 | 10 | 0 | 0 | 1 | 0 | 1 | 0 | 288 | 0 | 0 | 0 | 0 | |
| MUNCHENB | 8 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | |
| OSLO | 4 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | |
| SONNBLICK | 3 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | |
| STOCKHOLM | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | |
| VALENTIA | | | | | | | | | | | | | | | |

**Part 3 – Iteration**

To further refine the models and improve predictions for the Air Ambulance company, breaking the dataset into smaller, more focused segments is recommended. Dividing the data by location, time intervals, or weather features can provide valuable insights:

1. **By Location**: Segmenting the data by individual weather stations or grouping stations with similar patterns (e.g., coastal vs. inland) can help capture localized trends and improve model precision.

2. **By Time Intervals**: Analyzing data seasonally, monthly, or annually could highlight temporal patterns and provide insights into seasonal variations critical for flight safety.

3. **By Weather Features**: Narrowing the focus to specific variables like temperature, precipitation, or sunshine can help pinpoint which conditions are most predictive of safe flying days.

In terms of model selection, both Random Forest and CNN have distinct advantages. Random Forest is a great starting point for its interpretability and efficiency, particularly for single-station predictions where it achieved 100% accuracy. However, the CNN model is better suited for identifying complex relationships and temporal trends, making it an excellent complementary tool for broader, more advanced analyses.

For the Air Ambulance company, prioritizing weather variables like temp_max and precipitation is essential, as these consistently proved to be the most critical factors in predicting weather conditions. By focusing on these key predictors, the models can offer practical, reliable insights for planning helicopter flights.