

CLIMATEWINS  
PROJECT SUMMARY



# CLIMATE WINS

Keanu Gomes March, 2024  
*CareerFoundry Student*



# I. PROJECT OBJECTIVE AND HYPOTHESES

[Tableau Dashboard](#)

## GOAL OBJECTIVE

Apply optimization algorithms, supervised and unsupervised machine learning techniques to predict the consequences of climate change as a data analyst at the European non profit organization, ClimateWins.



## HYPOTHESES THAT CAN BE PROPOSED FROM THIS DATA:

1. Which algorithm predicts pleasant weather days best?
2. Will warmer temperatures correlate positively with the occurrence of pleasant weather days?
3. Does higher global radiation correspond to increased temperatures in cities?

# II. DATA ETHICS

## DATA SOURCE



<https://www.ecad.eu>

## BIAS TYPES



### Selection Bias

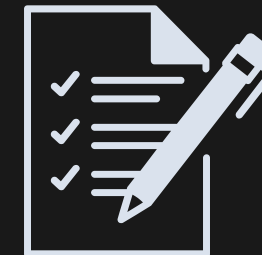
Only 18 out of 26321 weather stations were chosen as sample data

## DATA ACCURACY



The data selected for this analysis comes from reliable and trustable sourcing, as is it from 87 participants from verified meteorological stations across Europe totaling 26321 weather stations and 13 characteristics to be analyzed.

## DATA DIMENSIONS



22,951 rows x 170 columns



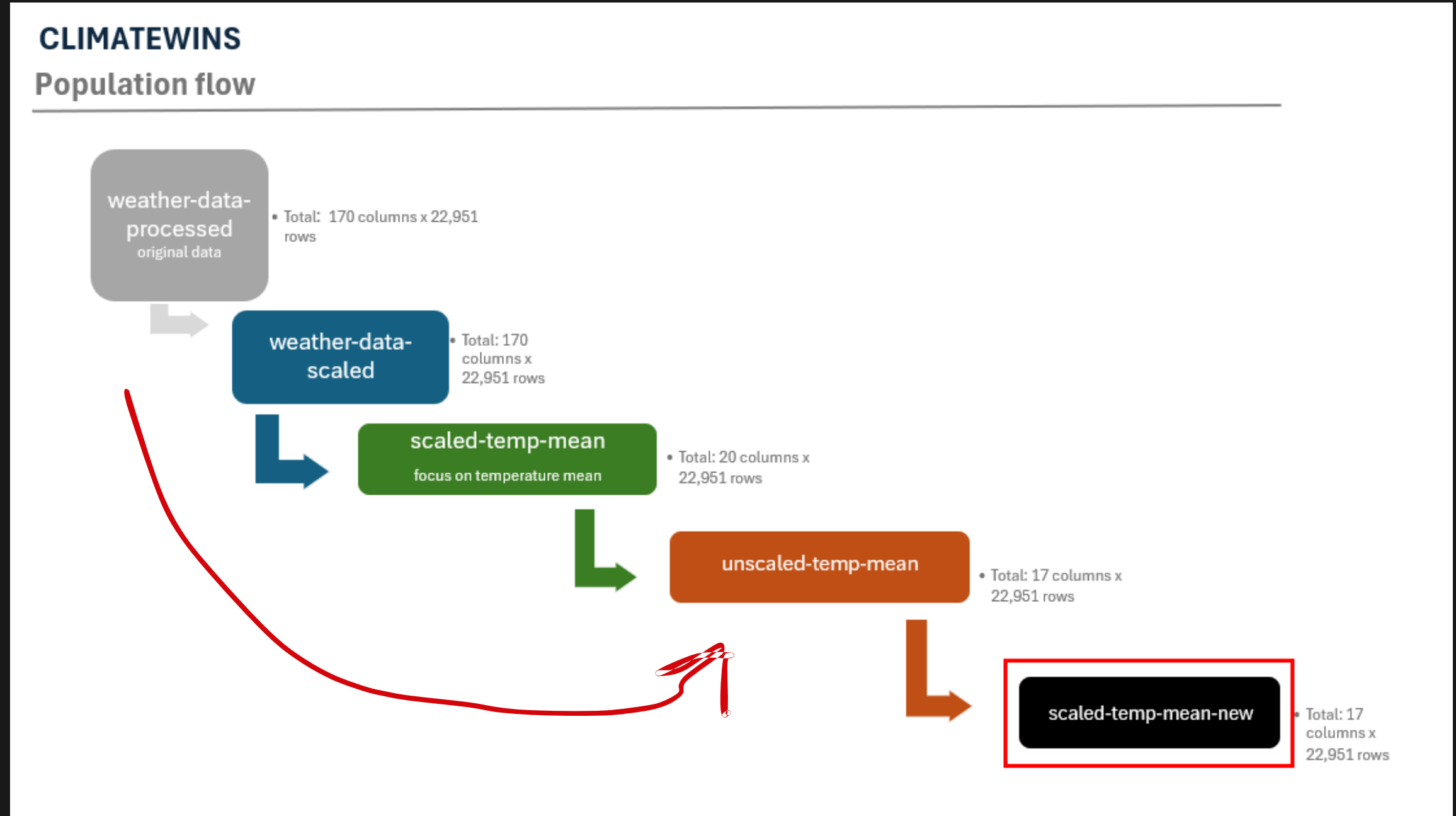
**18**

Total weather stations

# III. DATA WORKFLOW

## TEMPERATURE (MEAN) FOCUSED ANALYSIS

- Displayed on the right, is the population flow of my analysis through the (already) processed & cleaned data received from the weather stations.
- In order to feed the data into our supervised learning algorithms, we must removed non-pertinent columns & scale the data in order to normalize it for a more accurate analysis.

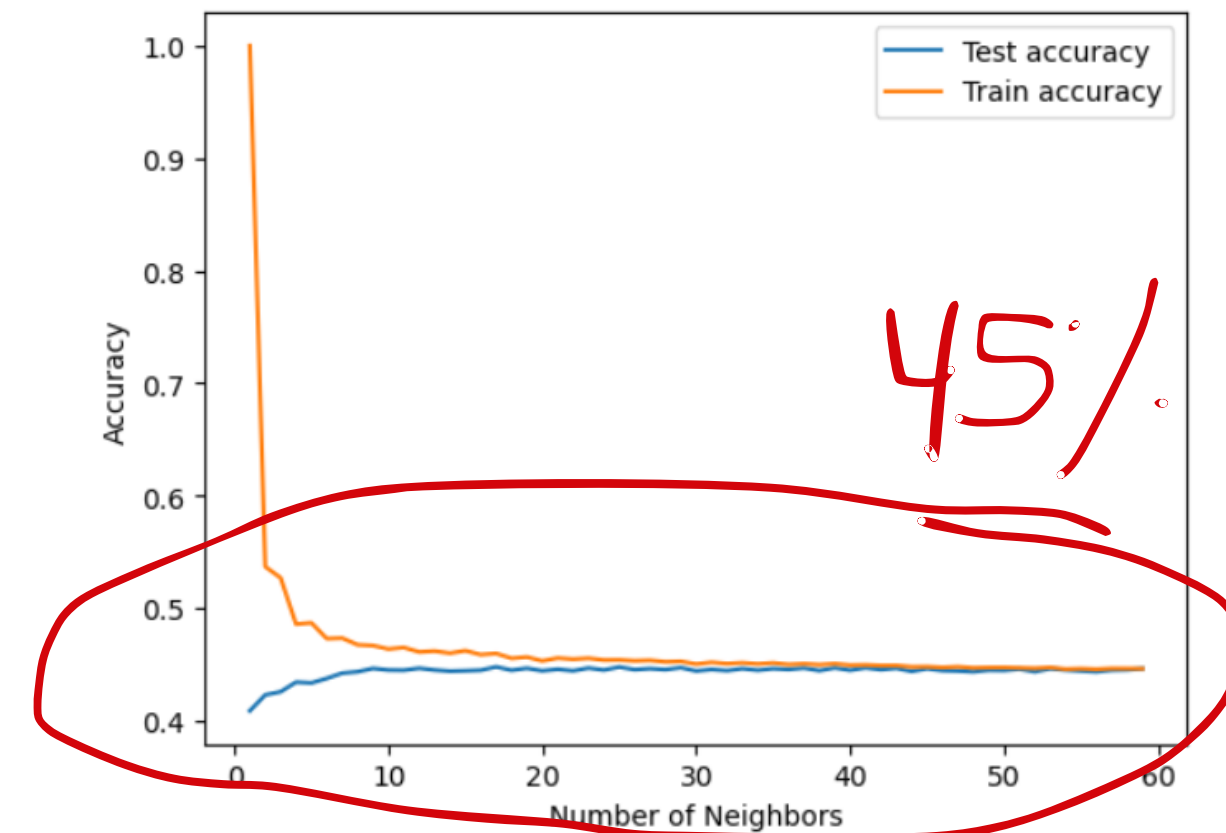
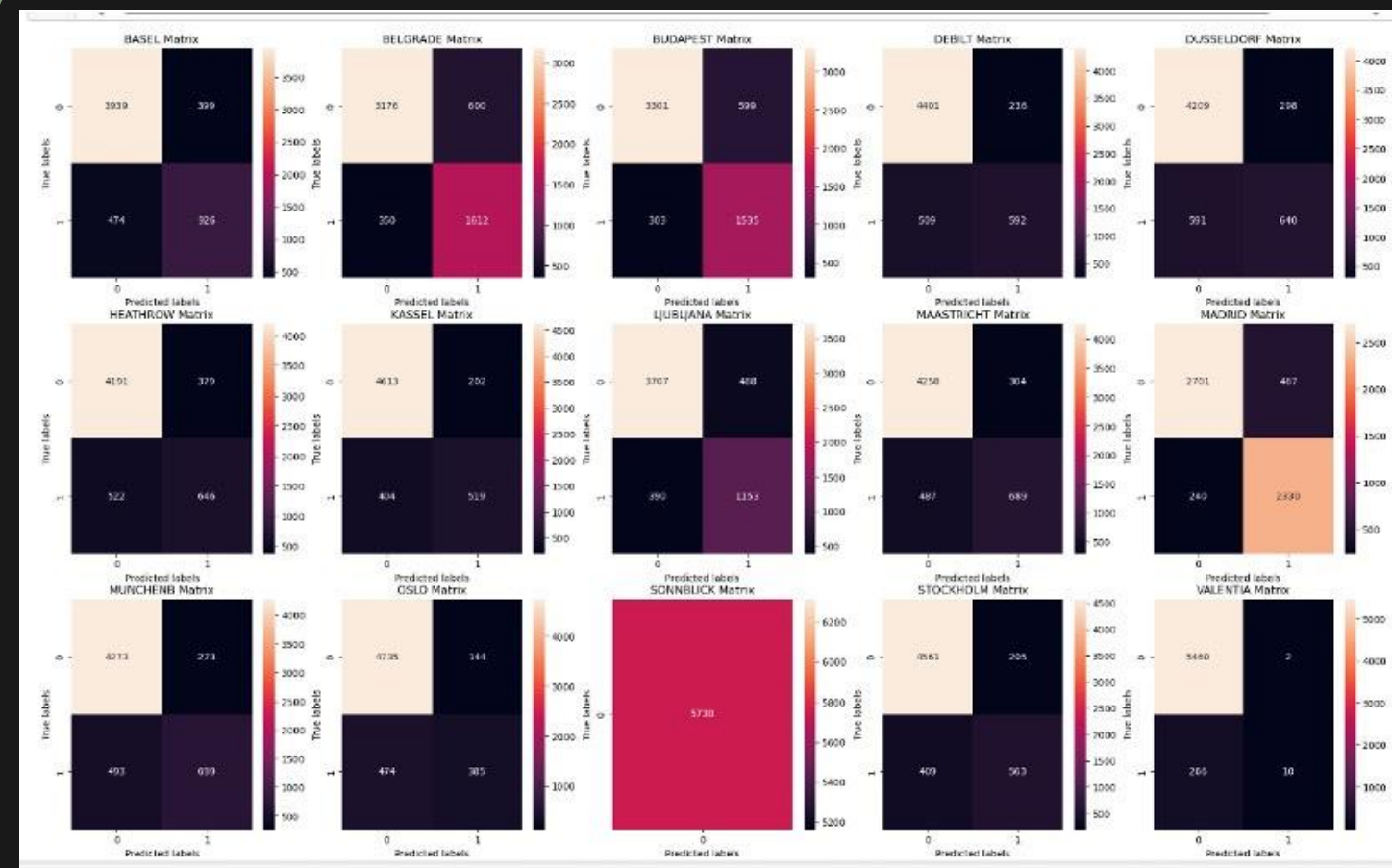




# V. ACCURACY SCORES

## KNN

### {K-NEAREST NEIGHBORS}



- Overall Testing Accuracy: approx. 0.45 or 45%
- Individual Station Accuracy: 0.82 - 0.95 or 82-95%
- Interpretation: Potential overfitting; lower overall accuracy compared to individual station scores.

# V.I. ACCURACY SCORES

{DECISION TREE}

```
In [12]: ▶ 1 #What is the testing accuracy score? Using the cross validation method
          2 y_pred = weather_dt.predict(X_test)
          3 print('Test accuracy score: ', accuracy_score(y_test, y_pred))
          4 multilabel_confusion_matrix(y_test, y_pred)
```

Test accuracy score: 0.4051934471941443

```
Out[12]: array([[3735, 603],
                [ 555, 845]],
              [[3143, 633],
                [ 622, 1340]],
```

40%

- Overall Testing Accuracy: approx. 0.405 or 40%
- Individual Station Accuracy: 0.82 - 0.95 or 82-95%
- Interpretation: Potential overfitting; lower overall accuracy compared to individual station scores.

# V.II. ACCURACY SCORES

ANN

{ARTIFICIAL NEURAL NETWORK}

```
2 mlp = MLPClassifier(hidden_layer_sizes=(20, 10, 10), max_iter=1000, tol=0.0001) #increasing hidden layers
3 #Fit the data to the model
4 mlp.fit(X_train, y_train)
```

Out[31]:

```
MLPClassifier
MLPClassifier(hidden_layer_sizes=(20, 10, 10), max_iter=1000)
```

In [32]:

```
1 y_pred = mlp.predict(X_train)
2 print(accuracy_score(y_pred, y_train))
3 y_pred_test = mlp.predict(X_test)
4 print(accuracy_score(y_pred_test, y_test))
```

```
0.45044155240529865
0.4527710003485535
```

45%

- Overall Testing Accuracy: approx. 0.452 or 45%
- Individual Station Accuracy: 0.82 - 0.95 or 82-95%
- Interpretation: Potential overfitting; lower overall accuracy compared to individual station scores.



## KNN/ANN/ OR DECISION TREE?

VI. How accurately do the algorithms predict pleasant and non-pleasant days per weather station?

### VALENTIA PREDICTION METRICS for 60 neighbors

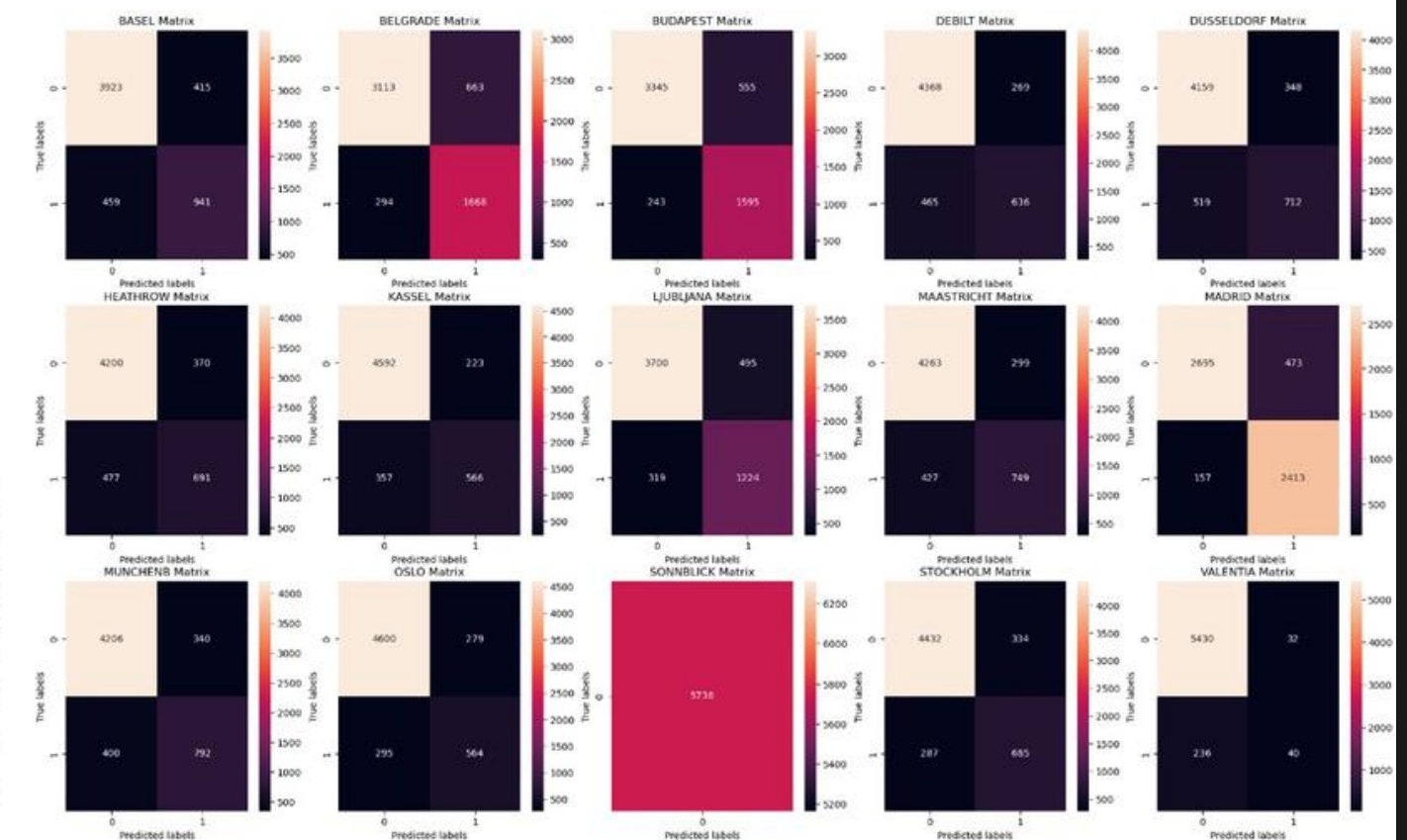
1. Accuracy: 95.34%
2. Precision: 99.96%
3. Recall (Sensitivity): 95.37%
4. F1 Score: 97.61%

## Confusion Matrix Scores (pleasant vs non-pleasant weather)

```
27  
28 print("Accuracy scores for each group:")  
29 for i, accuracy in enumerate(accuracy_scores):  
30     print(f"Group {i + 1}: {accuracy:.4f}")
```

Accuracy scores for each group:

Group 1: 0.8520  
Group 2: 0.8294  
Group 3: 0.8513  
Group 4: 0.8717  
Group 5: 0.8536  
Group 6: 0.8487  
Group 7: 0.9001  
Group 8: 0.8526  
Group 9: 0.8761  
Group 10: 0.8890  
Group 11: 0.8736  
Group 12: 0.8996  
Group 13: 1.0000  
Group 14: 0.8961  
Group 15: 0.9540



- VALENTIA seems to have the least false positives and negatives, & the highest number of true positives out of every station and algorithm used, this indicates that it may be the most accurate at the individual level.

VII. Which supervised learning algorithm types will be most effective for our hypotheses?



## PRIMARY RECOMMENDATIONS

### CONSISTENT TREND:

**40-45%**

Overall Accuracy

**82-100%**

Individual Station Accuracy

### VALENTIA STANDS OUT:

**95%**

Achieves high accuracy scores consistently around

### RECOMMENDATIONS:

- Investigate data quality and potential biases.
- Conduct feature importance analysis to leverage Valentia's strengths.
- Continue model refinement for improved accuracy.

### SUMMARY:

- Engage stakeholders to discuss implications and actions.
- All models show potential overfitting with lower overall accuracy compared to individual station scores.
- Further analysis and model refinement are necessary to address overfitting and improve generalization performance.



**THANKS FOR FOLLOWING ALONG!**

Any Questions?

Please contact me below at  
[keanudatatech@gmail.com](mailto:keanudatatech@gmail.com)

or visit:

<https://keanudatatech.github.io/portfolio>