



# ClimateWins: Predicting Weather Conditions

Isavannah Reyes

# TABLE OF CONTENTS

## **01** Objective

Objective and Hypotheses

## **02** Data and Optimization

Data features and gradient descent findings

## **03** Supervised Models

KNN, Decision Tree and ANN findings

## **04** Conclusion

Findings and future project steps

# Objective



**Objective:** Predict weather in mainland Europe by using historical weather data and machine learning methods.



**Why?** Increasing temperature and weather events has become more common in recent times. ClimateWins hopes to be able to predict these changes to help predict effects of climate change.



**Method:** Machine learning will allow ClimateWins to leverage pattern recognition often found in weather data to help predict future weather patterns.

# Hypothesis

## Hypotheses 1

If temperatures have historically risen, then machine learning will predict a continual rise in temperature.

## Hypothesis 2

Temperature in southern regions are more likely to be predicted as good weather.

## Hypothesis 3

Machine learning can identify climate change. For example, temperatures in Northern regions that were historically predicted as bad weather will more likely be predicted as good weather if temperatures become warmer.

# Data

- Contains temperature data from 18 different weather stations across Europe
  - Historical data from late 1800s to 2022
  - Daily data such as temperature min, temperature mean, temperature max
  - Additional readings for wind speed, snow, global radiation etc.
- Potential bias
  - 18 different weather stations of a geographically diverse continent may not fully cover the scope of weather patterns within Europe or outside of Europe if we wanted to expand the model
  - Pleasant and Unpleasant labels is subjective. Many people may one or the other. Defining this will likely not align with a small percentage of the population's own definitions of pleasant and unpleasant.
  - Weather stations such as, Sonnblick, have imbalanced data distribution. Sonnblick only has unpleasant weather labels. This skews predictions and overfits on unpleasant data predictions causing a bias towards unpleasant conditions.



Those who like rainy weather may call this good weather!

# Methodology

## Methodology – Supervised Modeling (Train and Test Data created)

- Compare F1 scores to account for data imbalance



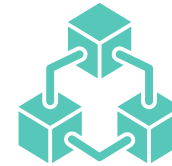
### KNN

- Non-parametric classifier that takes on supervised learning.
- Used to classify good weather and bad weather by reducing distance between similar groups
- Tested on both scaled and unscaled data



### Decision Tree

- Classification tree used to define good and poor weather.
- Function: DecisionTreeClassifier
- Gini criterion used to measure the quality of split
- Minimum number of samples required to split the node is 2



### ANN

- Input data goes through various layers to classify bad and good weather.
- Many layer adjustments but final model had three layers:
  - Layer 1: 100
  - Layer 2: 50
  - Layer 3: 25

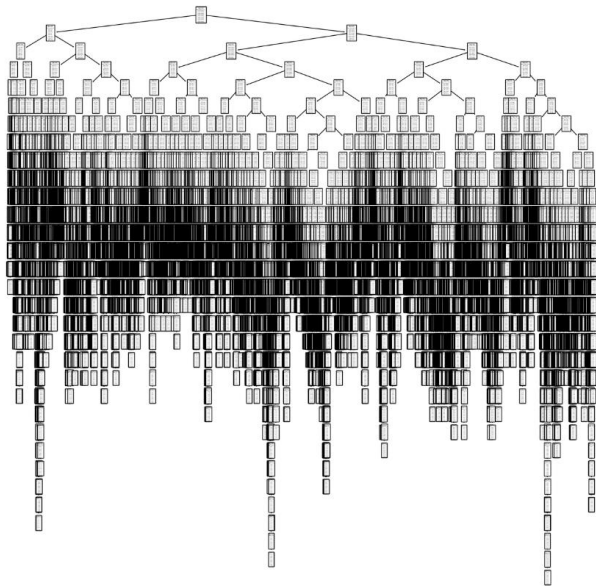
# KNN (Scaled)

- Train data scores higher than Test data both on average and per individual city
- Sonnblick only has 'unpleasant weather days' leading to Nan score

Station		F1 Score (Train Data)	Station		F1 Score (Test Data)
Basel		0.85	Basel		0.6925
Belgrade		0.87	Belgrade		0.7536
Budapest		0.88	Budapest		0.7599
Debilt		0.84	Debilt		0.6968
Dusseldorf		0.84	Dusseldorf		0.6698
Heathrow		0.84	Heathrow		0.6599
Kassel		0.82	Kassel		0.6602
Ljbljana		0.86	Ljbljana		0.7227
Maastricht		0.85	Maastricht		0.6902
Madrid		0.92	Madrid		0.8576
Munchen		0.84	Munchen		0.6629
Oslo		0.83	Oslo		0.6211
Sonnblick		Nan	Sonnblick		Nan
Stockholm		0.83	Stockholm		0.6520
Valentia		0.74	Valentia		0.5349
		0.86			0.688

# Decision Tree

- Train data scores higher than Test data both on average and per individual city
- Train data scores are scored 'perfectly' at 1 compared to Test data scores
- Sonnblick only has 'unpleasant weather days' leading to Nan score



Station	F1 Score (Train)	Station	F1 Score (Test)
Basel	1	Basel	0.67
Belgrade	1	Belgrade	0.72
Budapest	1	Budapest	0.72
Debilt	1	Debilt	0.67
Dusseldorf	1	Dusseldorf	0.67
Heathrow	1	Heathrow	0.63
Kassel	1	Kassel	0.63
Ljbljana	1	Ljbljana	0.69
Maastricht	1	Maastricht	0.68
Madrid	1	Madrid	0.85
Munchen	1	Munchen	0.67
Oslo	1	Oslo	0.63
Sonnblick	Nan	Sonnblick	Nan
Stockholm	1	Stockholm	0.64
Valentia	1	Valentia	0.64
	1		0.65



# ANN

- Train data scores higher than Test data both on average and per individual city
- Train data scores are scored 'perfectly' at 1 compared to Test data scores
- Sonnblick only has 'unpleasant weather days' leading to Nan score

Station	F1 Score	Station	F1 Score
Basel	0.84	Basel	0.78
Belgrade	0.88	Belgrade	0.84
Budapest	0.88	Budapest	0.84
Debilt	0.83	Debilt	0.78
Dusseldorf	0.82	Dusseldorf	0.76
Heathrow	0.84	Heathrow	0.77
Kassel	0.86	Kassel	0.75
Ljbljana	0.86	Ljbljana	0.82
Maastricht	0.8-	Maastricht	0.8-
Madrid	0.91	Madrid	0.91
Munchen	0.77	Munchen	0.77
Oslo	0.75	Oslo	0.75
Sonnblick	Nan	Sonnblick	Nan
Stockholm	0.79	Stockholm	0.79
Valentia	0.66	Valentia	0.66
	0.75		0.75

# Summarized Results: Supervised Models



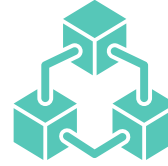
## KNN

- F1 scores for the scaled data mirror the other classifiers.
- Test data performs much better indicating potential overfitting ( $0.86 > 0.69$  avg F1 scores)



## Decision Tree

- Varies in performance between the different weather stations.
- The testing data scores perfect F1 scores while the test data scores an average of 0.65 which indicates overfitting



## ANN

- Does not converge on the test data without the use of further hyperparameters.
- **Highest average F1 Score: 0.75**

# Conclusion: Supervised Models

## Which model performed the best?

- KNN and decision tree performed similarly via F1 scores but the decision tree is overfit. ANN has the highest F1 score.
  - With a more complex black box (100, 50, 25) and higher iterations the ANN approach the f1 score is slightly higher than the KNN and decision tree classifiers.
  - ANN may be the best option to fine-tune going forward.
- A common issue is that the models are doing well on specific cities but not others. Using hyperparameters may help with this problem. Each model overfits on SONNBLICK data due to it only having one type of classified weather.

## Future Steps

- Fine-tuning ANN with hyperparameters
- If insufficient, other methods or even unsupervised models may be tested.
- Perfected model can be used to predict future weather events!



# Thank you!



QUESTIONS?

GITHUB: [ISAVANNAHR/CF\\_CLIMATEWINS](https://github.com/ISAVANNAHR/CF_CLIMATEWINS)