

MESSI Cluster Analysis

19.01.2023

Robert Wright

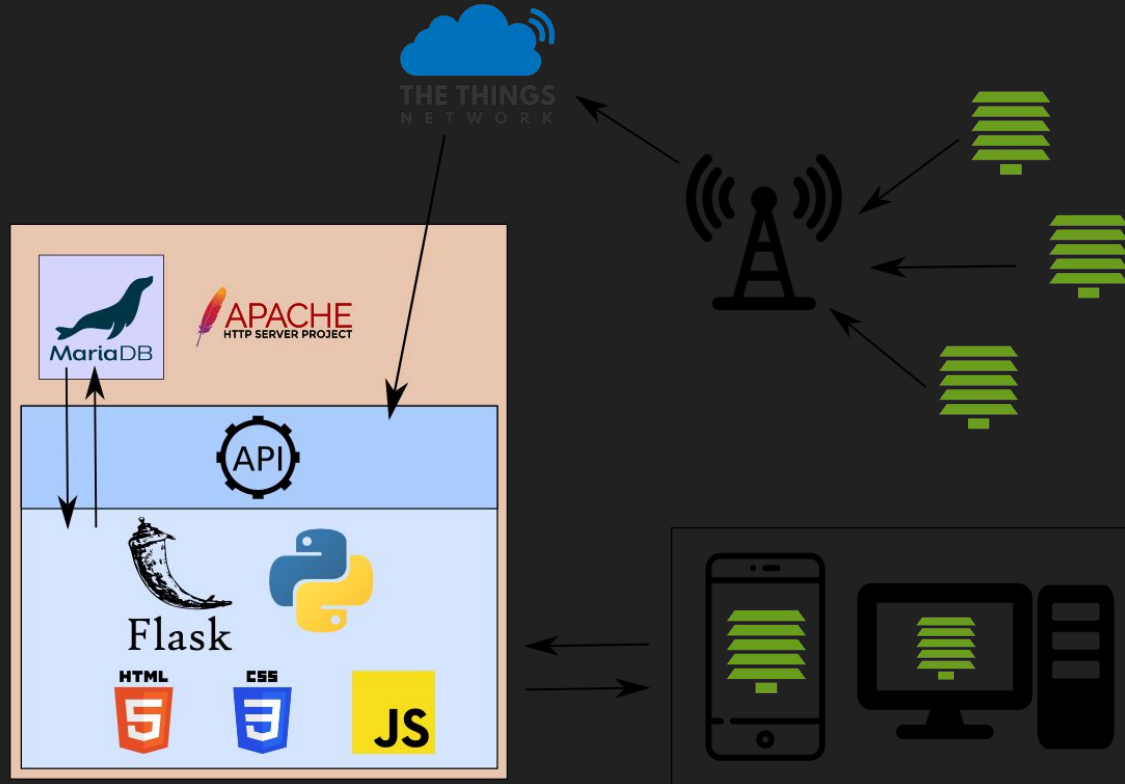
Weather & Climate Diagnosis

MESSIs

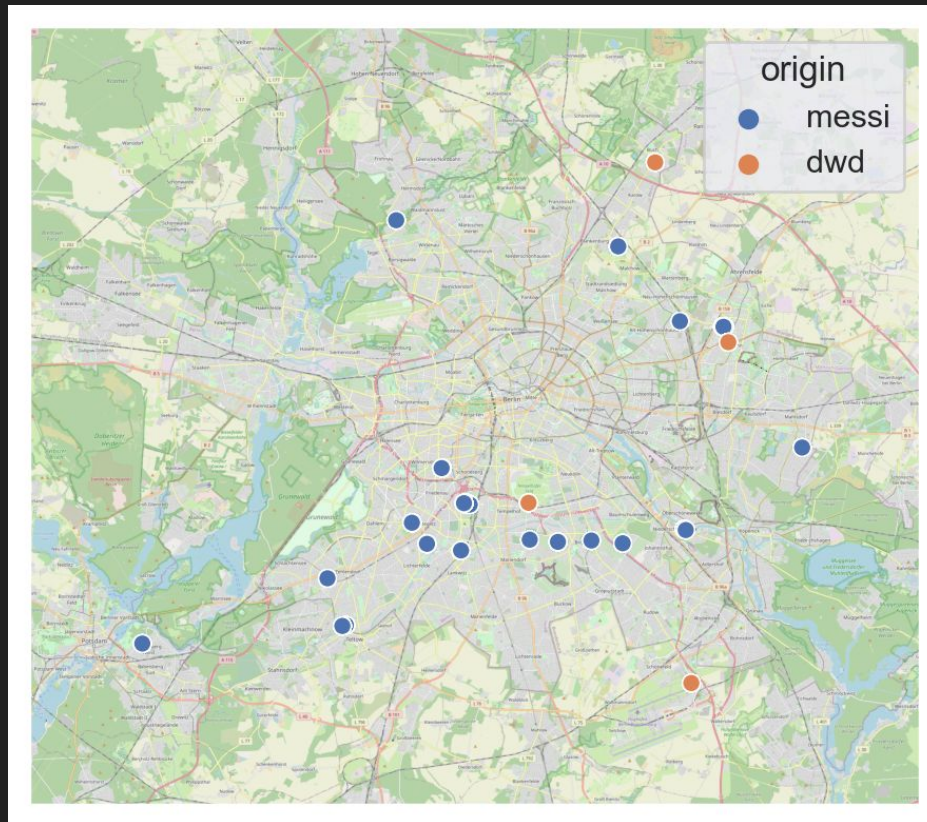
- = **M**ein **E**igenes **S**ub-**S**kalen Instrument
- Designed by FU & TU
- Citizen Science Project: [OpenUCO](#)



Data Assimilation

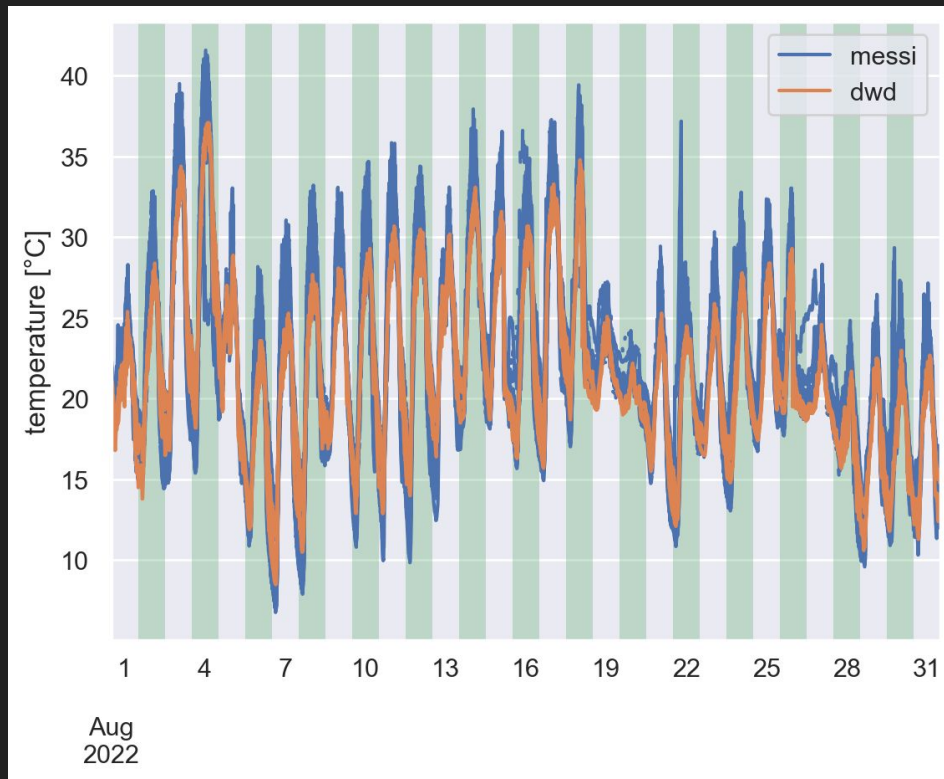


MESSI Locations



Data

- Temperature
- 10 min time interval
- August 2022

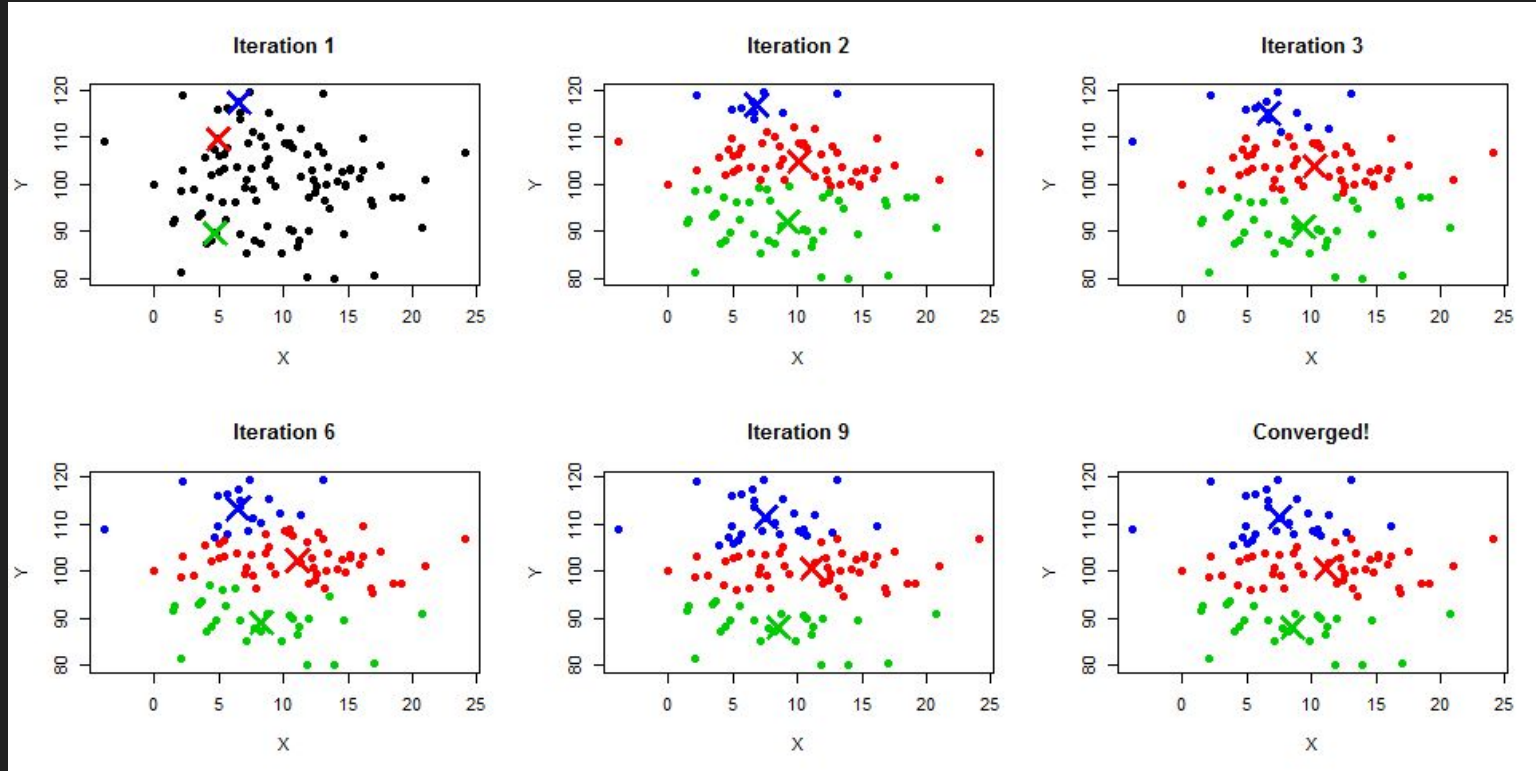


Clustering

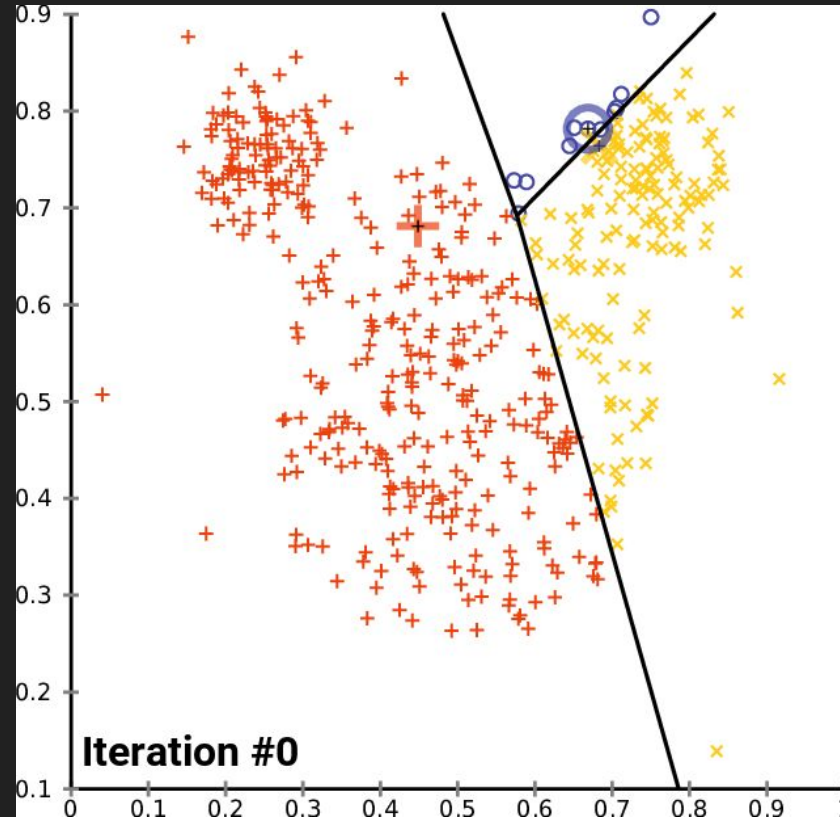
- kmeans
- Unsupervised learning = no target variable
- Separate samples into groups (*clusters*)
- # of clusters has to be specified
- Minimize within-cluster sum-of-squares (*inertia*):

$$SSE = \sum_{i=1}^k \sum_{\mathbf{x}_j \in S_i} \|\mathbf{x}_j - \boldsymbol{\mu}_i\|^2$$

kmeans



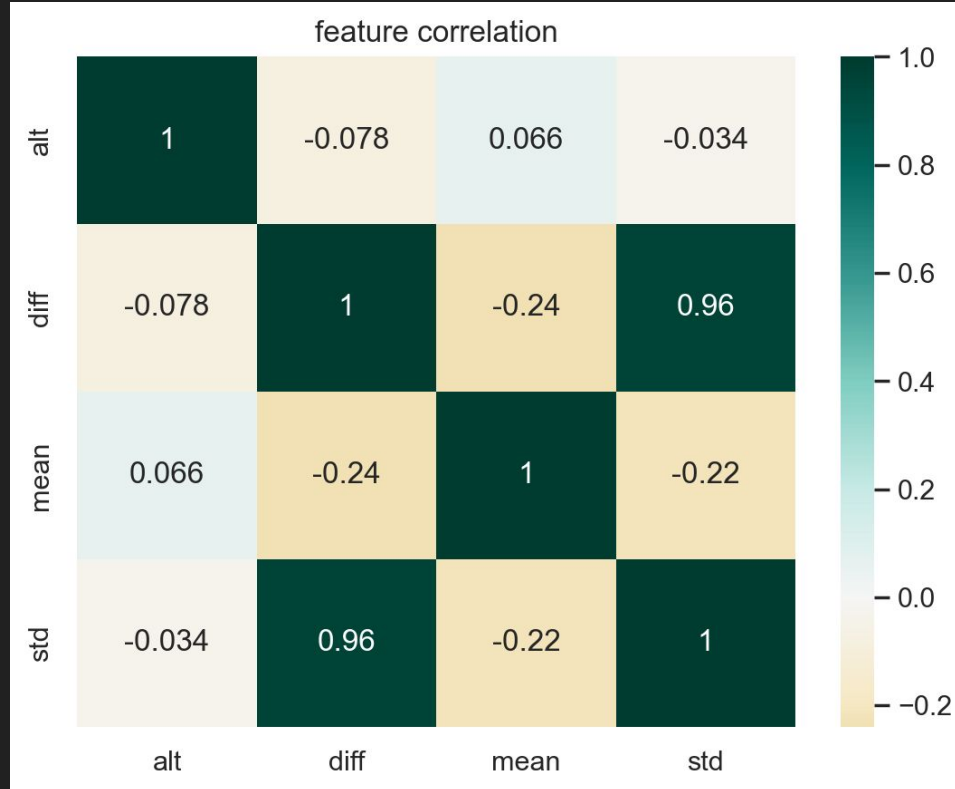
kmeans



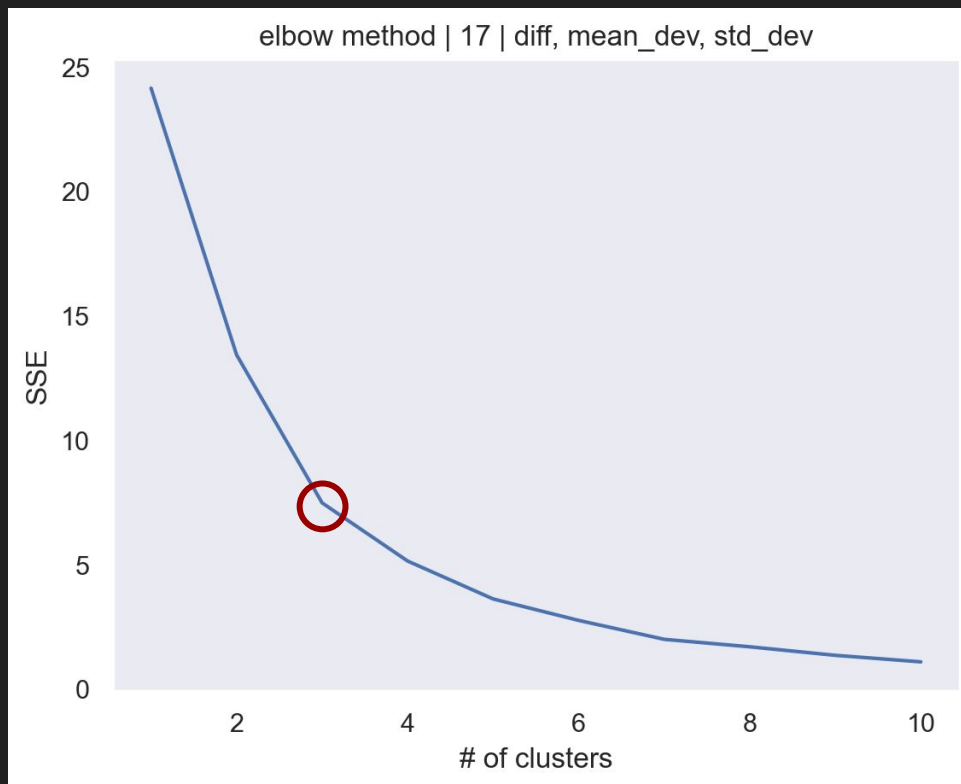
Feature Engineering

- Deviation from *dwd* daily mean: ΔT
- 1. Difference of maximum and minimum temperature: $\Delta T_{max} - \Delta T_{min}$
- 2. Mean: $\overline{\Delta T}$
- 3. Standard deviation: $\sigma \Delta T$

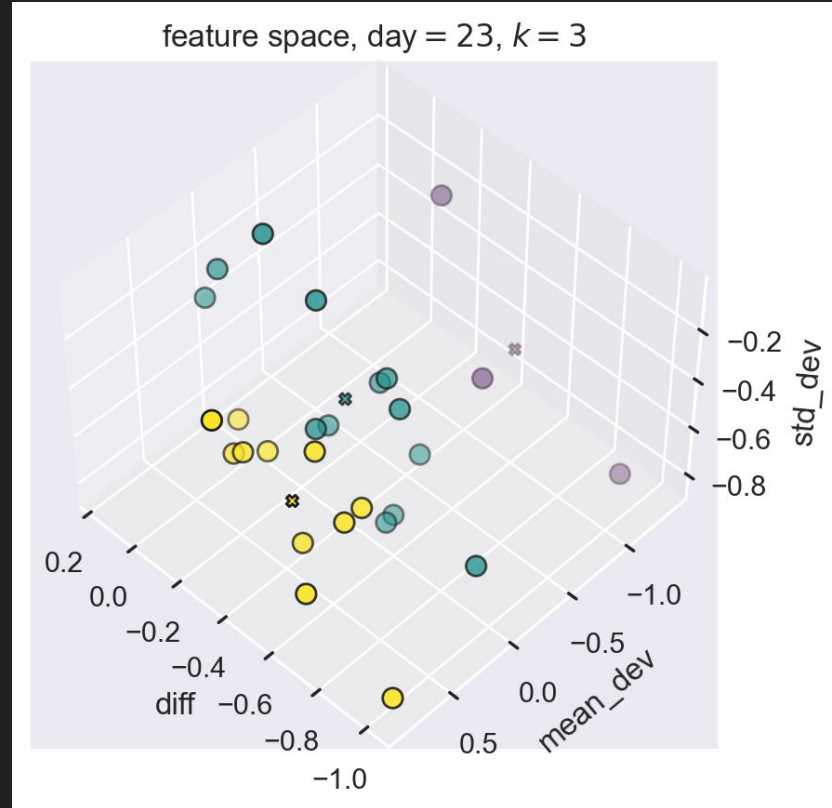
Feature Correlation



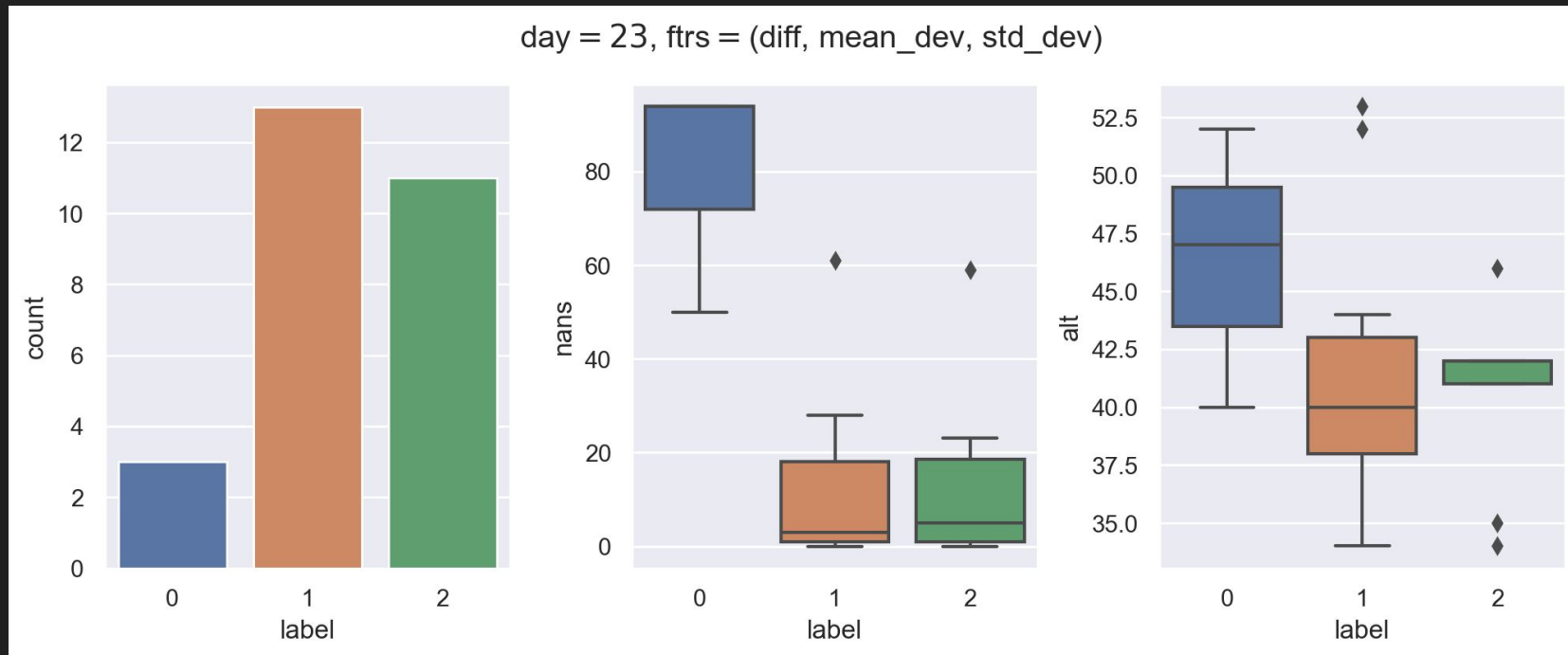
of Clusters



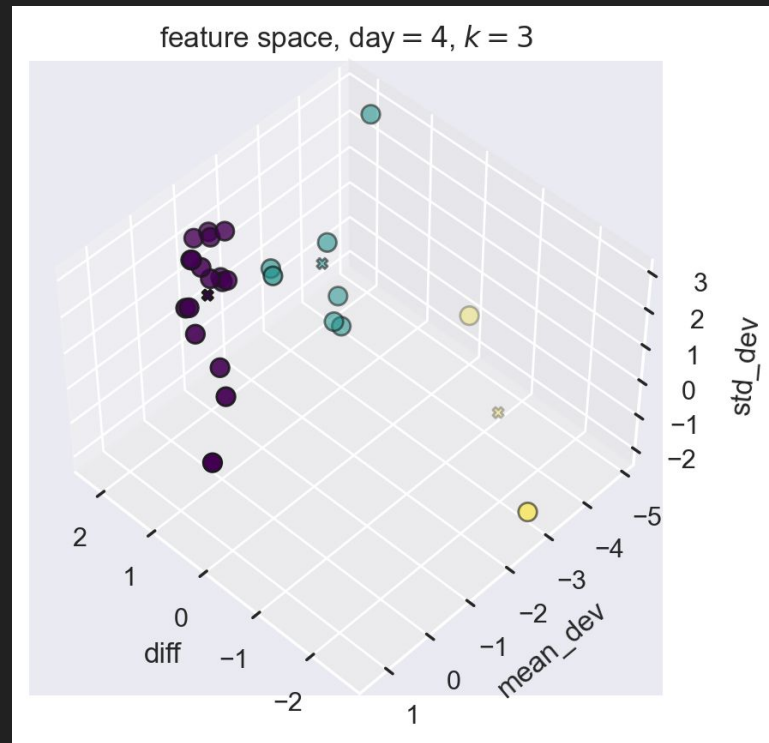
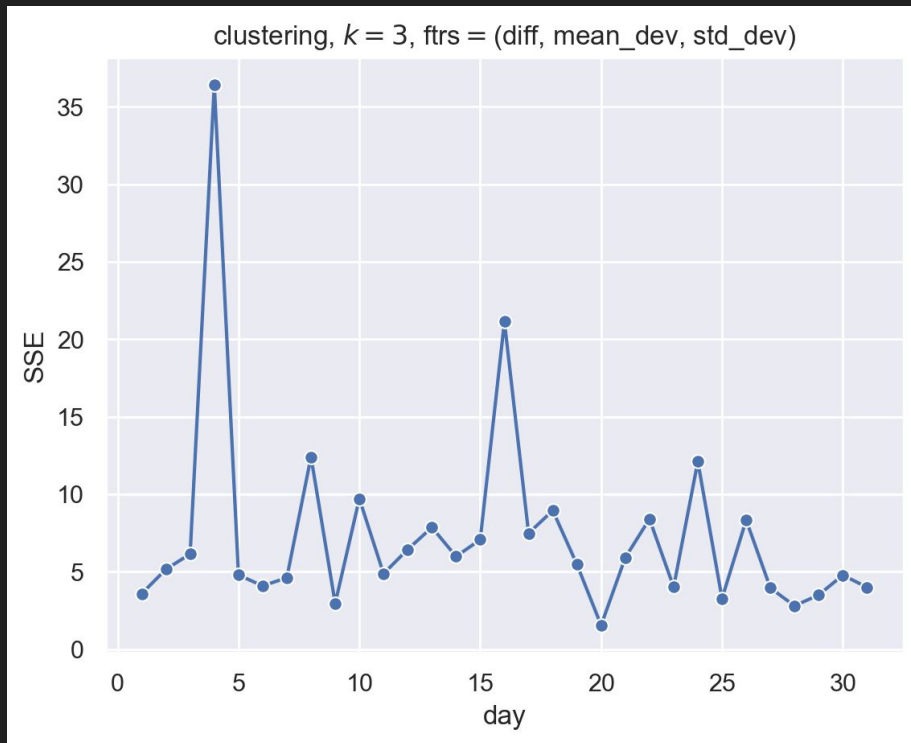
Clusters in Feature Space



Characteristics of Clusters

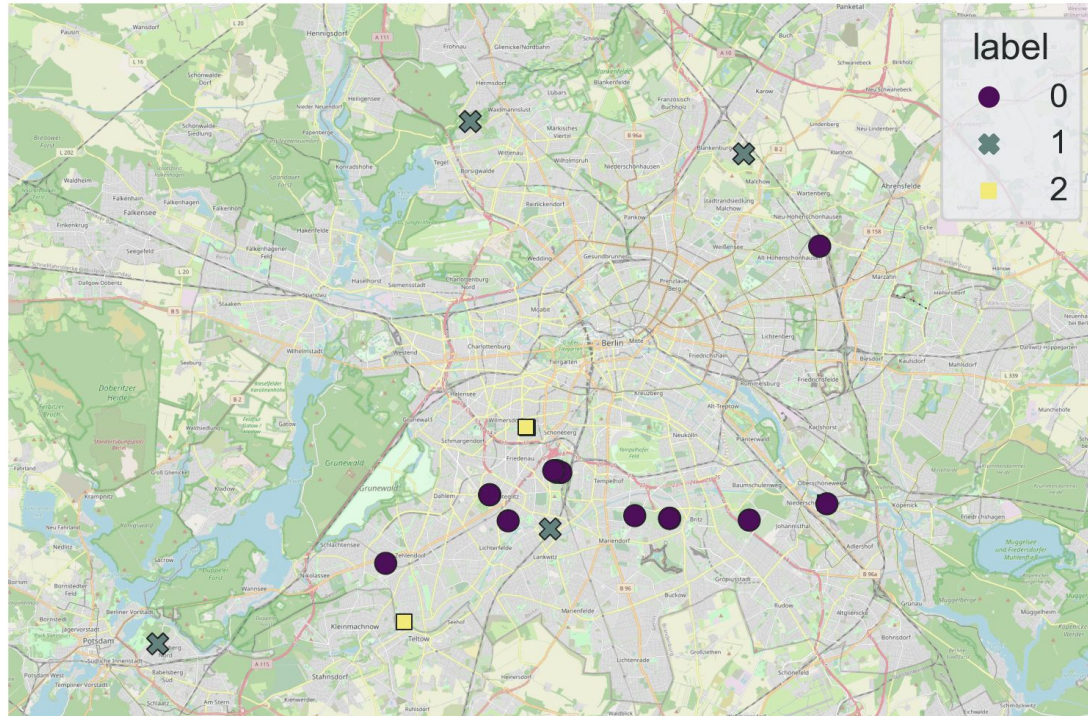


Inertia of Clusterings

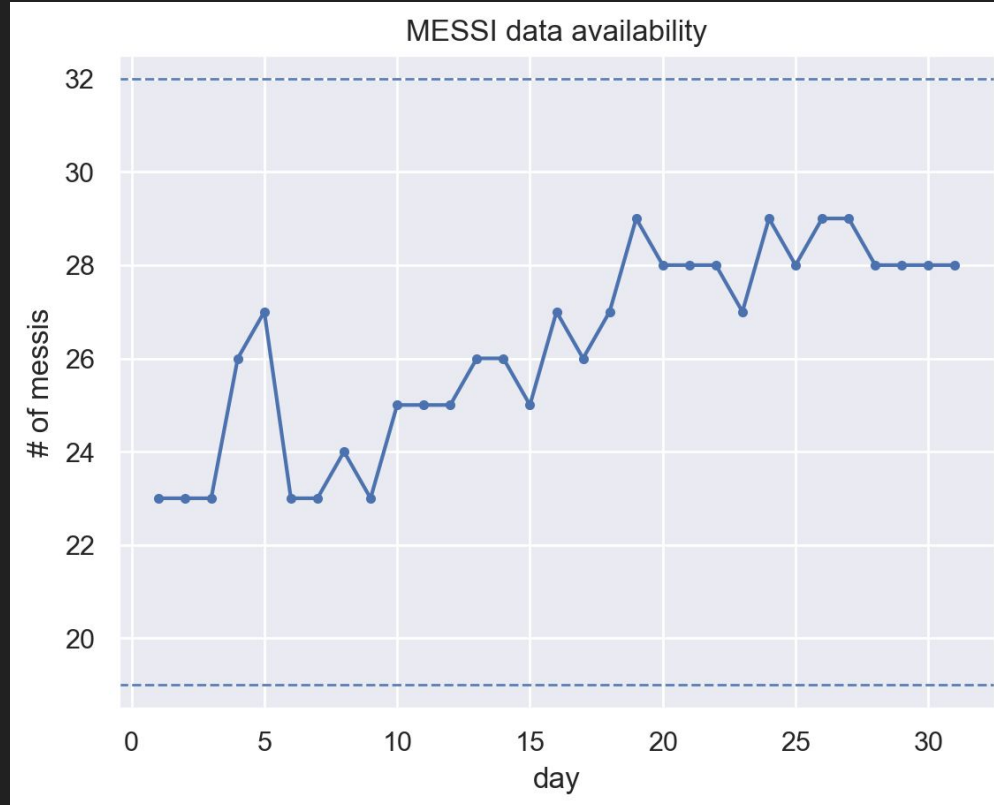


Clusters on Berlin map

day = 1



Changing # of MESSIs throughout August

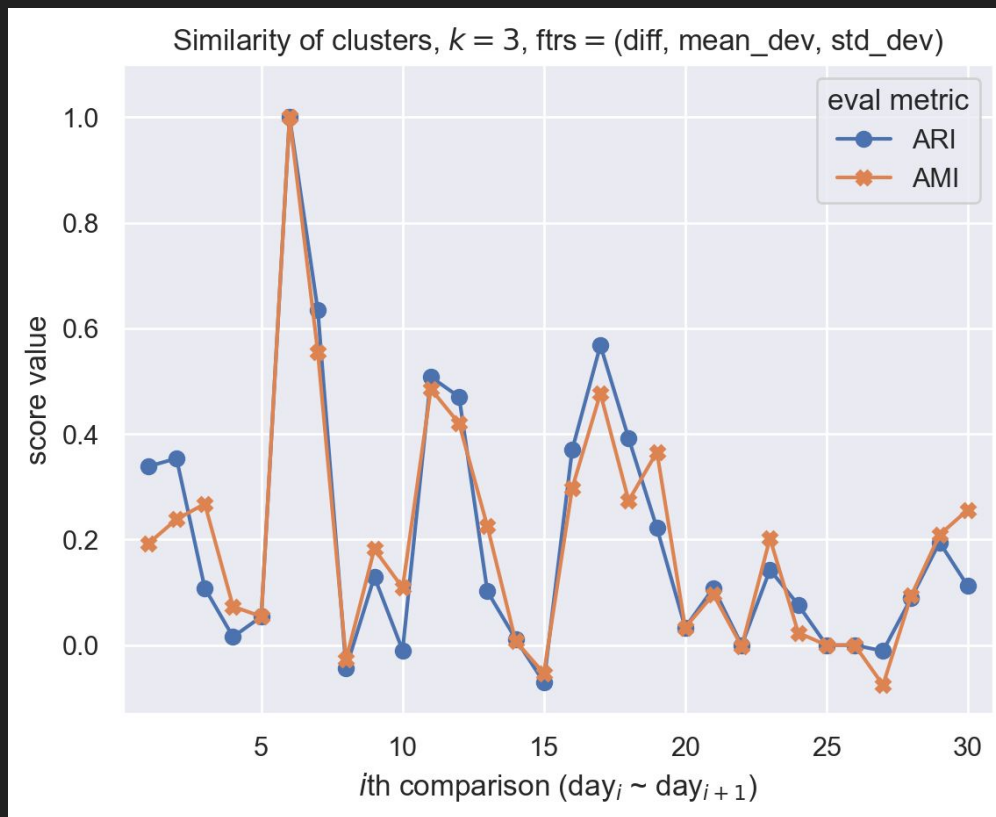


Rand Index

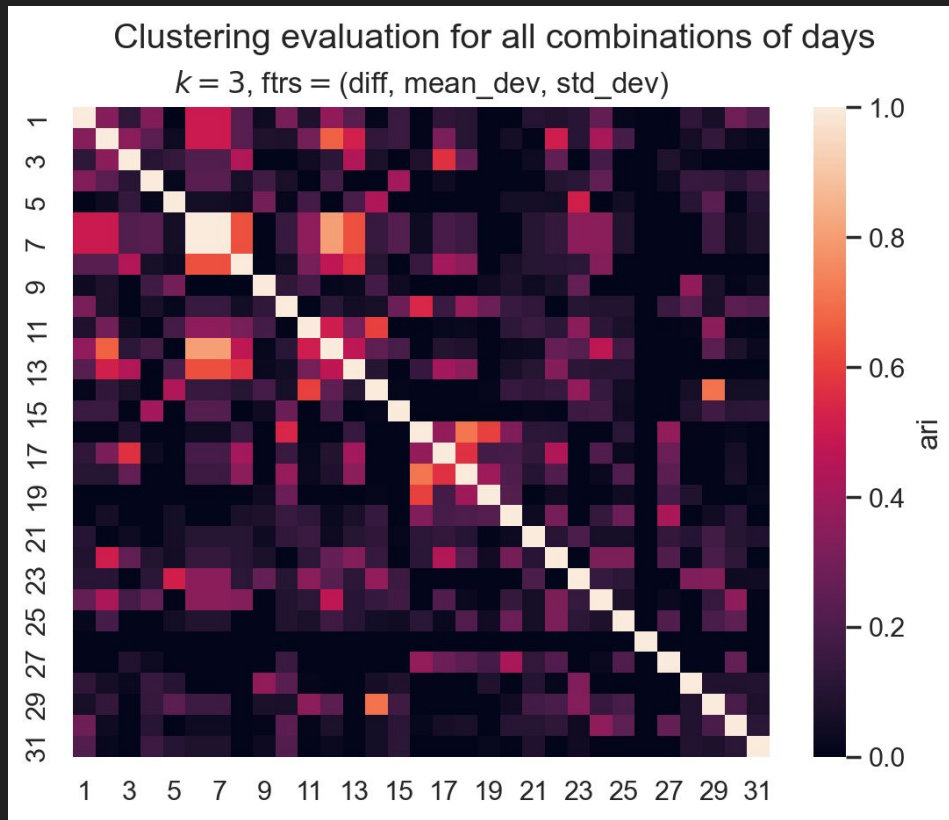
- Measures similarity of two cluster assignments
- Values between 0 (independent labellings) and 1 (perfect labelling)

$$RI = \frac{a + b}{C_2^{n_{samples}}}$$

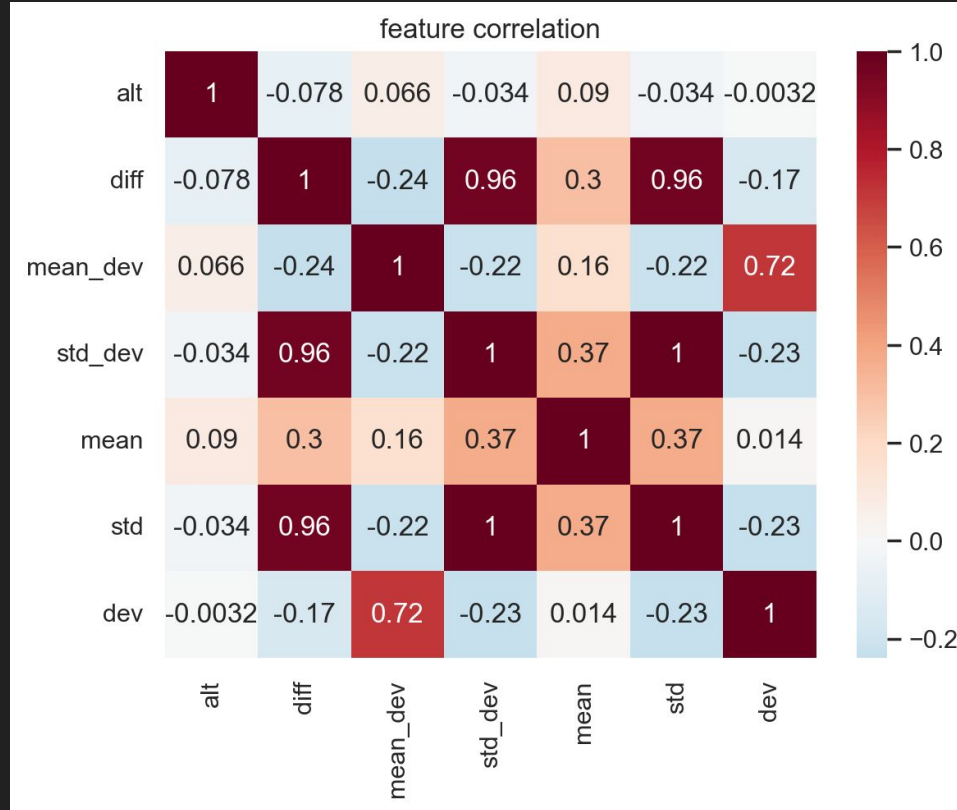
Cluster Similarity I



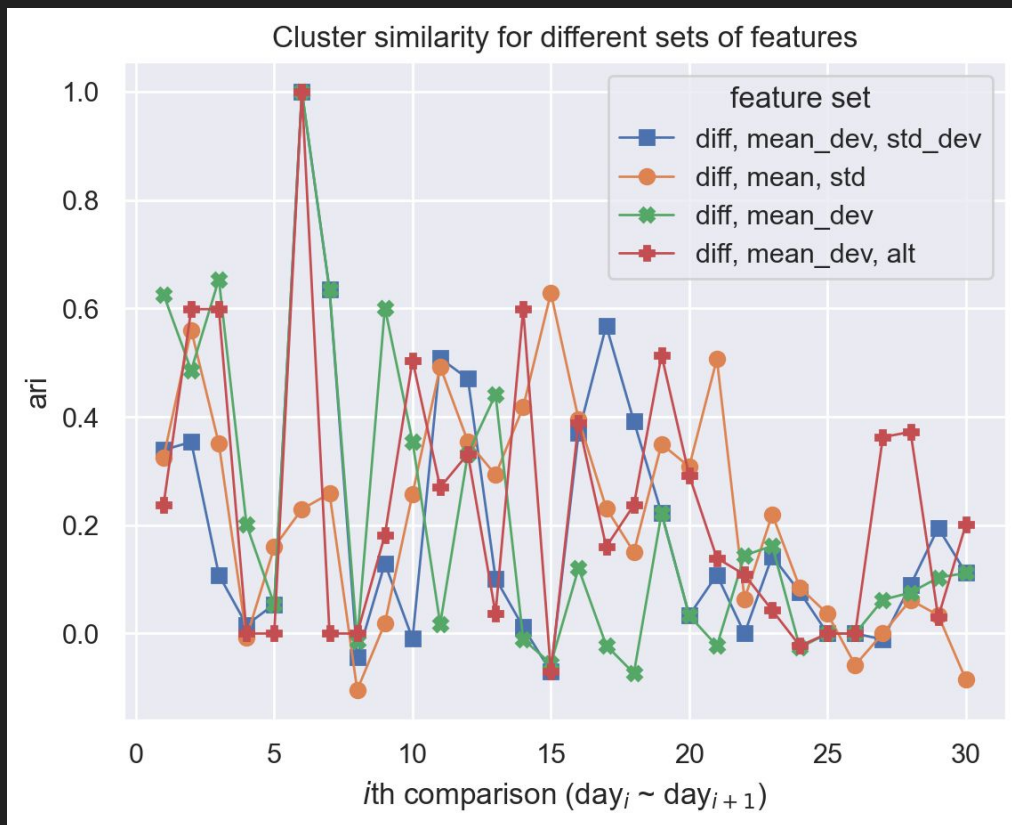
Cluster Similarity II



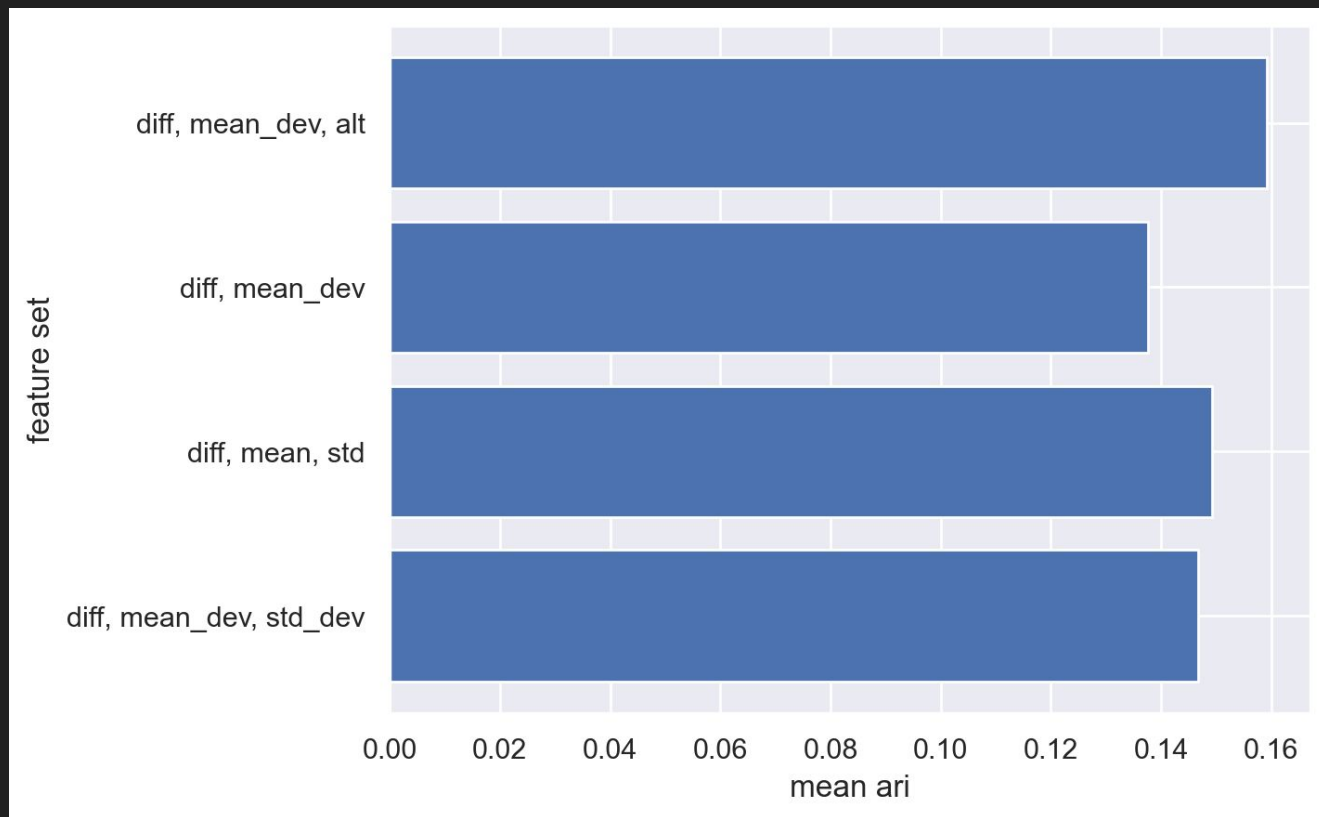
Different Sets of Features



Results I



Results II



Future Research

- Note: # of clusters is determined on a single day
- More MESSIs with continuous data / fill data gaps
 - Run clustering & performance evaluation on **same** data points
- Identify “meaning” of clusters with more meta-information

MESSI Cluster Analysis

19.01.2023

Robert Wright

Weather & Climate Diagnosis