

# **Statistical Data Analysis**

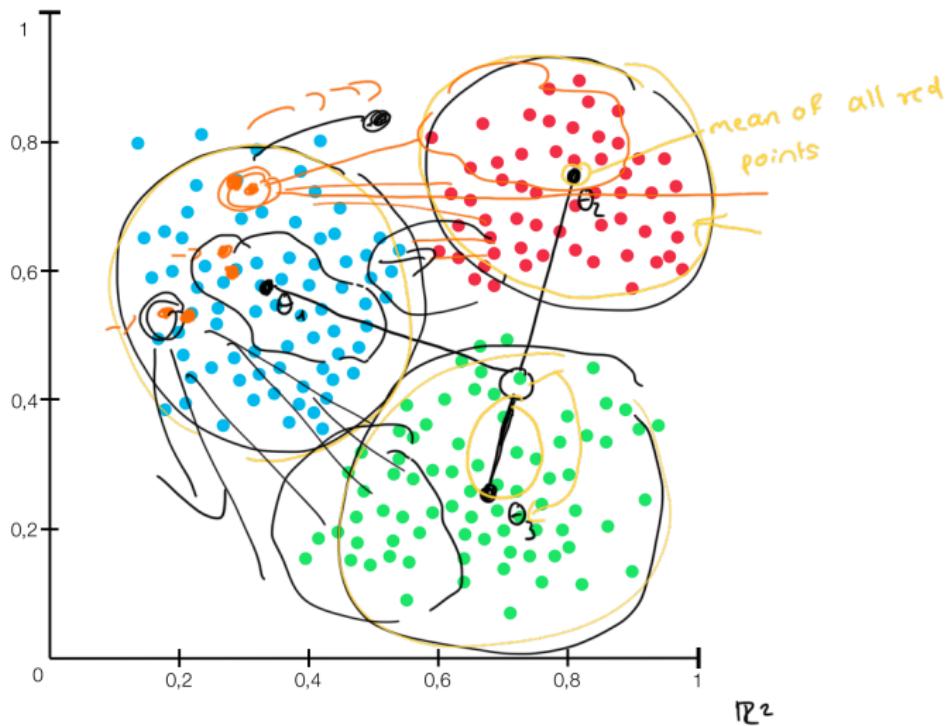
---

Dr. Jana de Wiljes

04.01.2022

Universität Potsdam

# Clustering



# K-means clustering

Input:

- Number of Clusters  $K$
- Set of points  $\{x_1, \dots, x_M\}$  in vector space that need to be classified

$$x_i \in \mathbb{R}^n$$

$n$  large  
 $M$  large

Output:

- Sets  $\mathcal{M}_k$  of the clusters

1. Initialize the centre of the cluster  $\theta_1, \dots, \theta_K \in \mathbb{R}^n$  randomly

2. Repeat till a stopping criterion is fulfilled {

for all  $k = 1 : K$

$$\mathcal{M}_k := \{\}$$

for all  $m = 1 : M$

$$j = \arg \min_h \|\theta_h - x_m\|_2^2$$

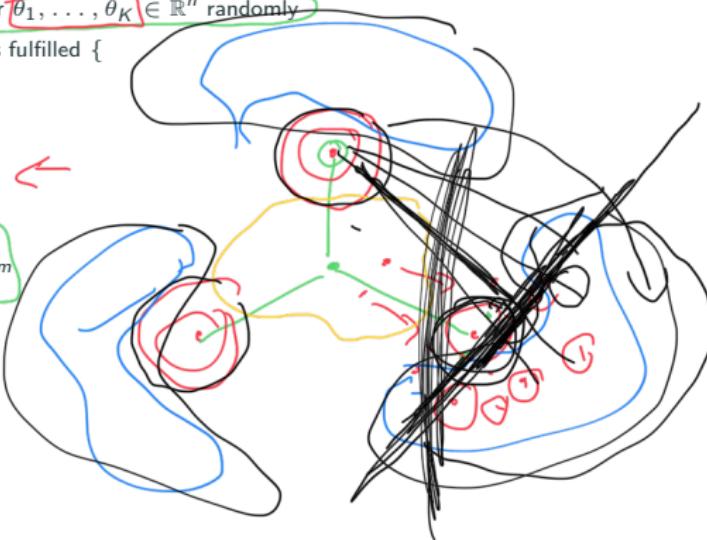
$$\mathcal{M}_j = \mathcal{M}_j \cup \{x_m\}$$

for all  $k = 1 : K$

$$\theta_k = \frac{1}{|\mathcal{M}_k|} \sum_{x_m \in \mathcal{M}_k} x_m$$

3. return  $\theta_1, \dots, \theta_K$

$$\theta_1, \theta_K$$



# Initialisation

- Random Partition Method

- Forgy Initialization

- kmeans++

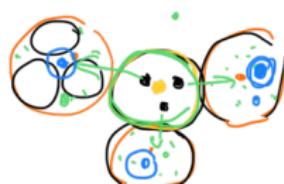
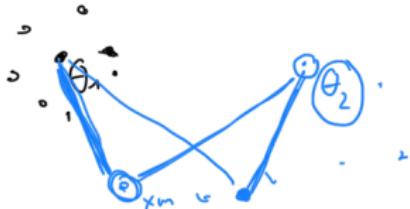
1. choose  $\theta_1$  uniformly at random from set of points
2. Choose new center  $\theta_i$  with probability

$$\frac{D(x_m)^2}{\sum_{x_l} D(x_l)^2}$$

$$D(x_m) \approx D(x_i)$$
  
  
(1)

where  $D(x_m)$  denotes the shortest distance from data point  $x_m$  to the closest center we have already chosen

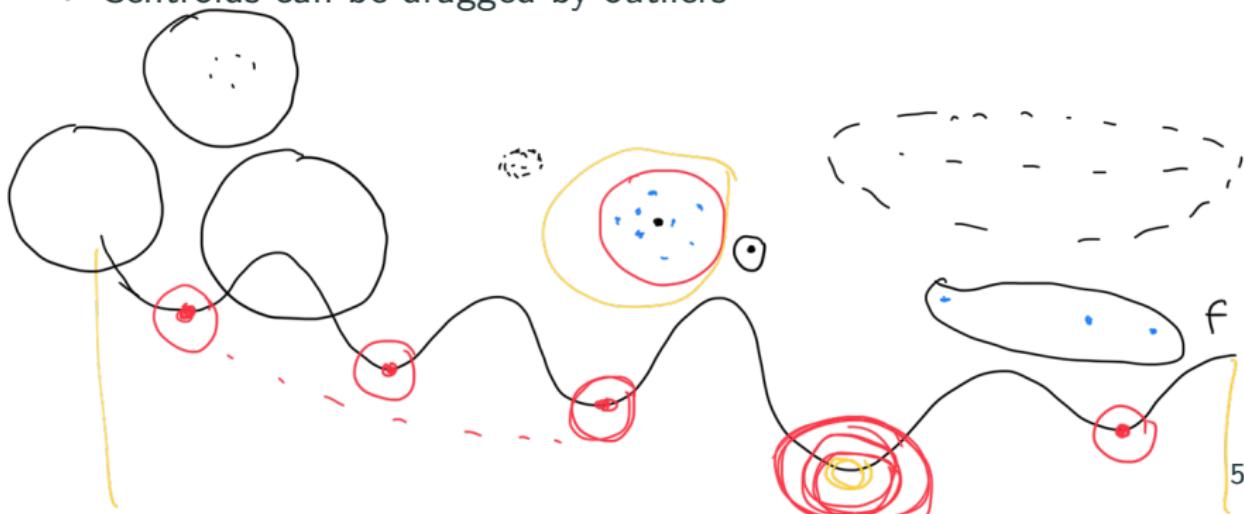
3. Repeat Step 2 until we have all K centers



# K-means clustering

## Disadvantages

- true number of clusters K unknown (requires tuning)
- K-means algorithm depends on the chosen initial values ↵
- Clustering data of varying sizes and density
- Centroids can be dragged by outliers



# Using the Triangle Inequality to Accelerate k-Means

Goal : Accelerate Kmeans with the help of the triangle inequality  $d(x, z) \leq d(x, y) + d(y, z)$  for any 3 vectors and any distance metric  $d$



Basic idea: avoid comparisons by having upper and lower bounds

Let  $x \in V$  and  $\theta_1, \theta_2 \in V$  centers. Need to know  $d(x, \theta_1) \geq d(x, \theta_2)$

Auxiliary Lemma 1: Let  $x \in V$ , let  $\theta_1, \theta_2$  be centers

if  $d(\theta_2, \theta_1) \geq 2d(x, \theta_2)$  then  $d(x, \theta_1) \geq d(x, \theta_2)$

Proof: we know that  $d(\theta_2, \theta_1) \leq d(\theta_2, x) + d(x, \theta_1) \Rightarrow [d(\theta_2, \theta_1) - d(x, \theta_2)] \leq d(x, \theta_1)$   
 $d(\theta_2, \theta_1) - d(x, \theta_2) = [2d(x, \theta_2) - d(x, \theta_2)] = d(x, \theta_2)$  now just consider left hand side  
so  $d(x, \theta_2) \leq d(x, \theta_1)$

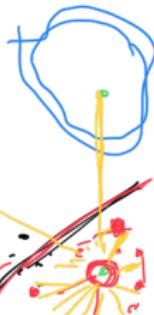
Lemma 2: Let  $x$  be a vector and let  $\theta_2$  and  $\theta_1$  be centers

Then  $d(x, \theta_1) \geq \max\{0, d(x, \theta_2) - d(\theta_2, \theta_1)\}$

Proof: we know that  $d(x, \theta_2) \leq d(x, \theta_1) + d(\theta_2, \theta_1)$  so

$d(x, \theta_1) \geq d(x, \theta_2) - d(\theta_2, \theta_1)$  Also  $d(x, \theta_1) \geq 0$

How to use this?  $x$  is currently assigned to  $\theta_2$  and  $\theta_1$  is another center  
 $\frac{1}{2} [d(\theta_2, \theta_1)] \geq d(x, \theta_2)$  then  $d(x, \theta_1) \geq d(x, \theta_2)$



not necessary to compute  $d(x, \theta_1)$

⚠ might not know  $d(x, \theta_2)$  exactly, but do know upper bound  $u$ :  $u \geq d(x, \theta_2)$

# Using the Triangle Inequality to Accelerate k-Means

~ need to compute:  $d(x, \theta_1)$  and  $d(x, \theta_2)$  only if  $u > \frac{1}{2} d(\theta_1, \theta_2)$   
if  $u = \frac{1}{2} \min_{\theta'} (\theta_1, \theta')$  where the minimum over all  $\theta' \neq \theta_1, \theta_2$   
then the point  $x$  must remain with the cluster center it has been assigned to.

Lemma 2 (use): Let  $x$  be any data point, let  $\theta$  be any center, let  $\theta'$  be the previous version of  $\theta$

• Suppose that in the previous iteration we know a lower bound  $\ell'$  such that  $d(x, \theta') \geq \ell'$

• Then we consider a lower bound  $\ell$  for current iteration:

$$\begin{aligned} d(x, \theta) &\geq \max \{0, \underline{d(x, \theta')} - d(\theta, \theta')\} \\ &\geq \max \{0, \underline{\ell'} - d(\theta, \theta')\} = \ell \end{aligned}$$

• suppose  $u(x) \geq d(x, \theta)$

$\theta$  current cluster of  $x$

• Suppose  $P(x, \theta') \leq d(x, \theta')$

lower bound on the distance between  $x$  and different center  $\theta'$

if  $u(x) < P(x, \theta')$  then

necess. to compute

~ not  $d(x, \theta')$

$$\boxed{d(x, \theta)} \leq \boxed{u(x)} \leq \boxed{P(x, \theta')} \leq \boxed{d(x, \theta')}$$

## Using the Triangle Inequality to Accelerate k-Means

---

# Using the Triangle Inequality to Accelerate k-Means

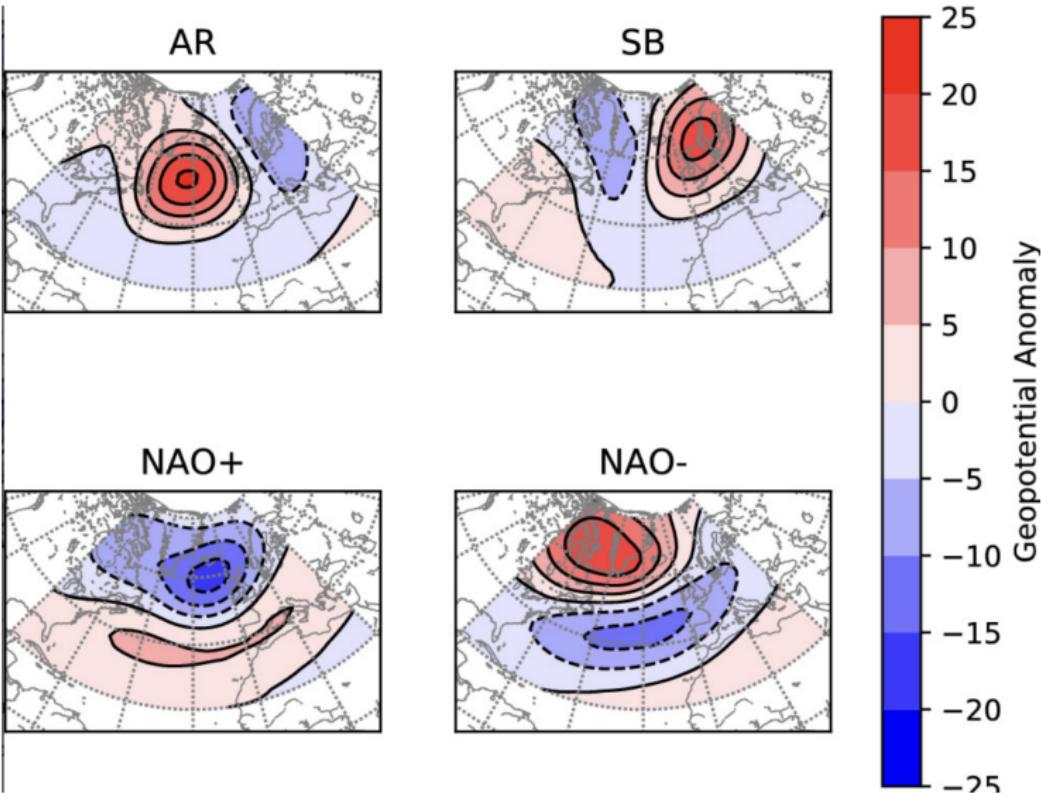
Algorithm:

1. Initialize the centre of the cluster  $\theta_1, \dots, \theta_K \in \mathbb{R}^n$  randomly ↵
2. Set lower bounds to  $l(x_m, \theta_i) = 0$  for all  $\theta_i$  and  $x_m$
3. Assign each  $x_m$  to its closest initial center  $\theta(x_m) = \arg \min_h \|\theta_h - x_m\|_2^2$  (avoid redundant calculations using Lemma 1)  
 $d(x_m, \theta)$  while  $\theta$  is the current center of  $x_m$
4. Each time  $\|\theta_h - x_m\|_2^2$  is computed, set  $l(x_m, \theta_h) = \|\theta_h - x_m\|_2^2$
5. Assign upper bounds  $u(x_m) = \min_i \|\theta_i - x_m\|_2^2$
6. Repeat till a stopping criterion is fulfilled {
  - 6.1 for all  $\theta_i$  and  $\theta_j$ , compute  $\|\theta_i - \theta_j\|_2^2$ . For all centers  $\theta_i$ , compute  $s(\theta_i) = \frac{1}{2} \min_j \|\theta_i - \theta_j\|_2^2$
  - 6.2 Identify all points  $x_m$  such that  $u(x_m) \leq s(\theta(x_m))$ .
  - 6.3 for all centers  $\theta_i$  for all remaining points  $x_m$  check
    - $\theta_i \neq \theta(x_m)$  and
    - $u(x_m) > l(x_m, \theta_i)$  and
    - $u(x_m) > \frac{1}{2} \|\theta(x_m) - \theta_i\|_2^2$If conditions  $r(x_m) = \text{true}$  are true compute  $\|x_m - \theta(x_m)\|$  and assign  $r(x_m) = \text{false}$ . Otherwise  $\|x_m - \theta(x_m)\|_2^2 = u(x_m)$ .
  - 6.4 if  $\|x_m - \theta(x_m)\|_2^2 > l(x_m, \theta_i)$  or  $\|x_m - \theta(x_m)\|_2^2 > \frac{1}{2} \|\theta(x_m) - \theta_i\|_2^2$  then
    - compute  $\|(x_m - \theta_i)\|_2^2$
    - if  $\|(x_m - \theta_i)\|_2^2 < \|x_m - \theta(x_m)\|_2^2$  then assign  $\theta(x_m) = \theta_i$
7. for all centers  $\theta_i$ , let  $m(\theta_i)$  be the mean of the points assigned to  $\theta_i$
8. for all points  $x_m$  and for all centers  $\theta_i$  assign  $l(x_m, \theta_i) = \max\{l(x_m, \theta_i) - \|\theta_i - m(\theta_i)\|_2^2, 0\}$
9. for all points  $x_m$ , assign  $u(x_m) = u(x_m) + \|m(\theta(x_m)) - \theta(x_m)\|$  and  $r(x_m) = \text{true}$
10. replace each center  $\theta_i$  with  $m(\theta_i)$
11. return  $\theta_1, \dots, \theta_K$

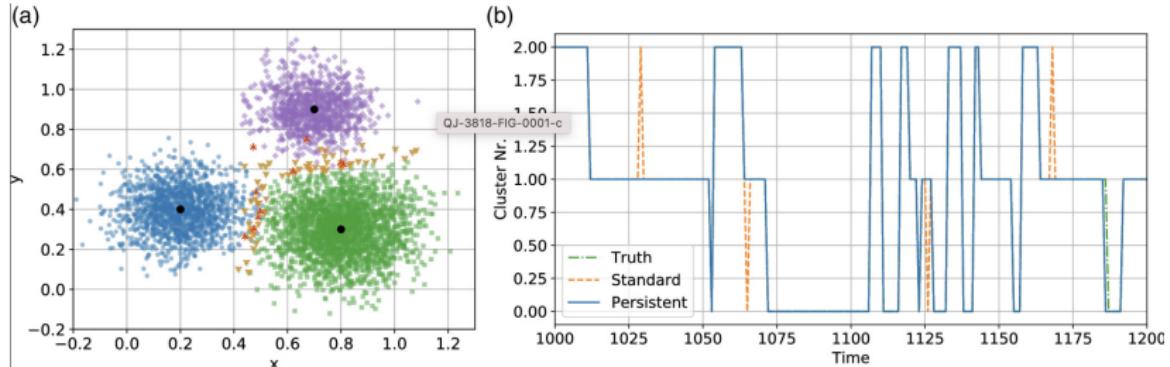
## **Example: pattern recognition for atmospheric circulation regimes**

---

# Regime

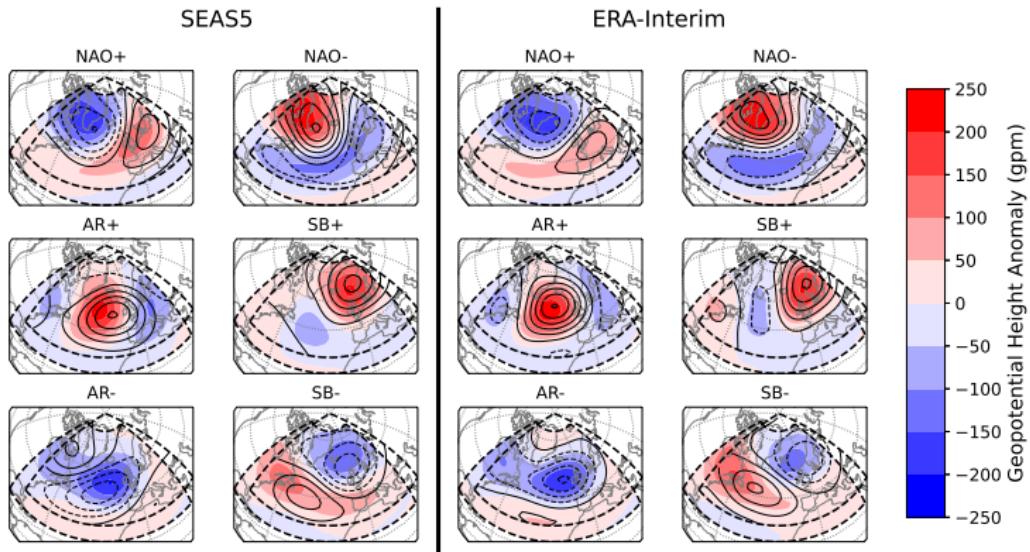


# Time persistency constraint

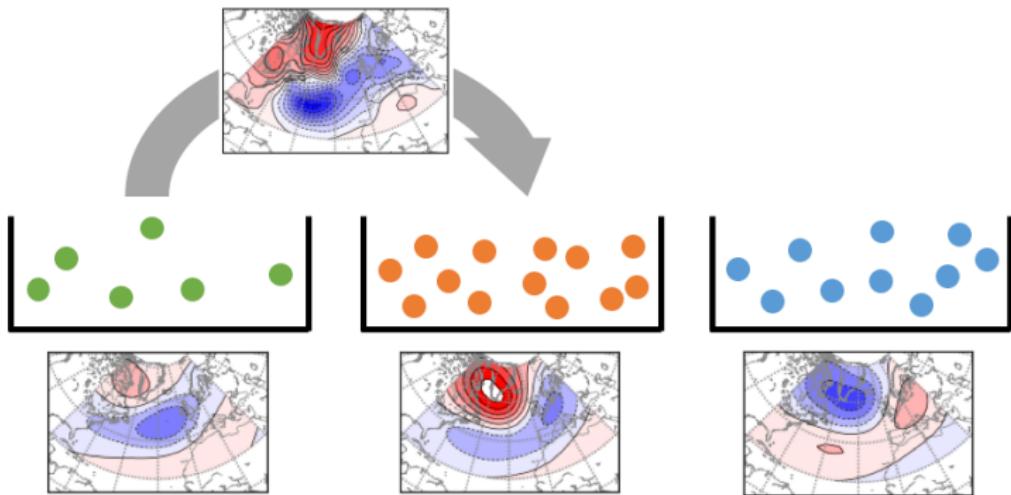


$$\sum_{t=1}^{T-1} |\gamma_k(t+1) - \gamma_k(t)| \leq N_C \quad \forall k$$

# *k*-means clustering for different domains



## $k$ -means clustering for different domains



# Optimisation problem

$$\mathbf{L}(\Theta, \Gamma) = \sum_{t=0}^T \sum_{n=1}^N \sum_{i=1}^k \gamma_i(t, n) \|x_{t,n} - \theta_i\|^2$$

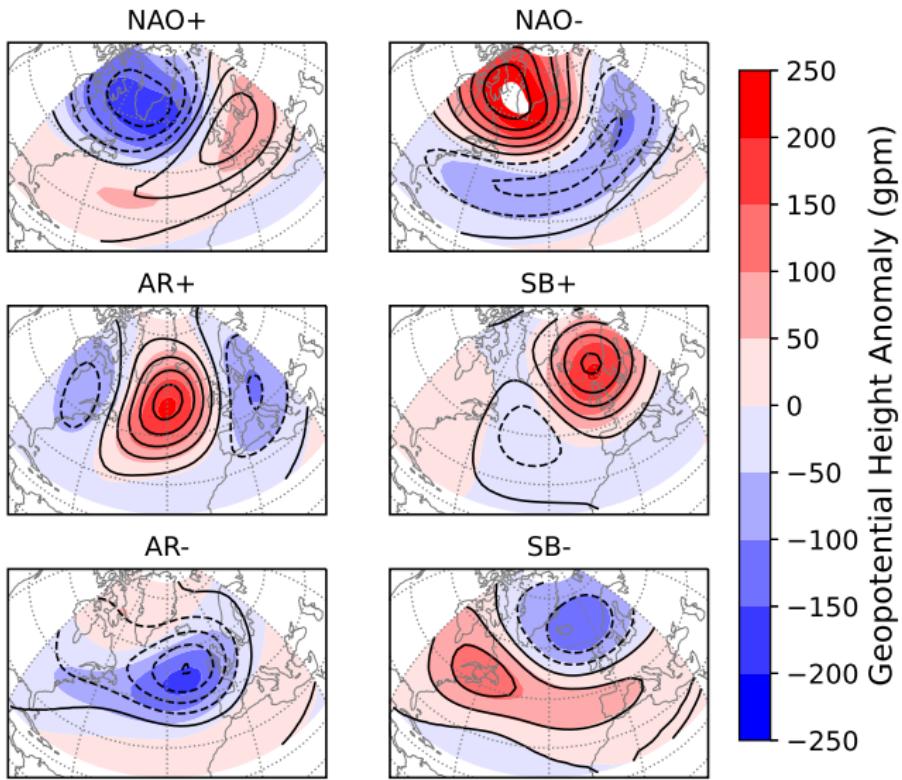
with

$$\sum_{i=1}^k \gamma_i(t, n) = 1, \quad \forall t \in [0, T], \quad \forall n \in [1, N].$$

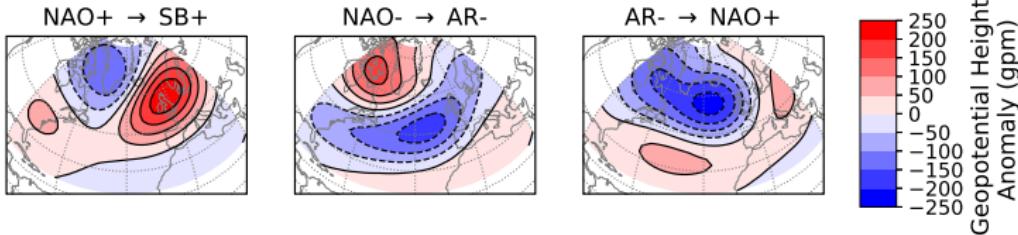
and

$$\sum_{i=1}^k \sum_{n_1, n_2} |\gamma_i(t, n_1) - \gamma_i(t, n_2)| \leq \phi \cdot C_{\text{eq}}, \quad \forall t \in [0, T],$$

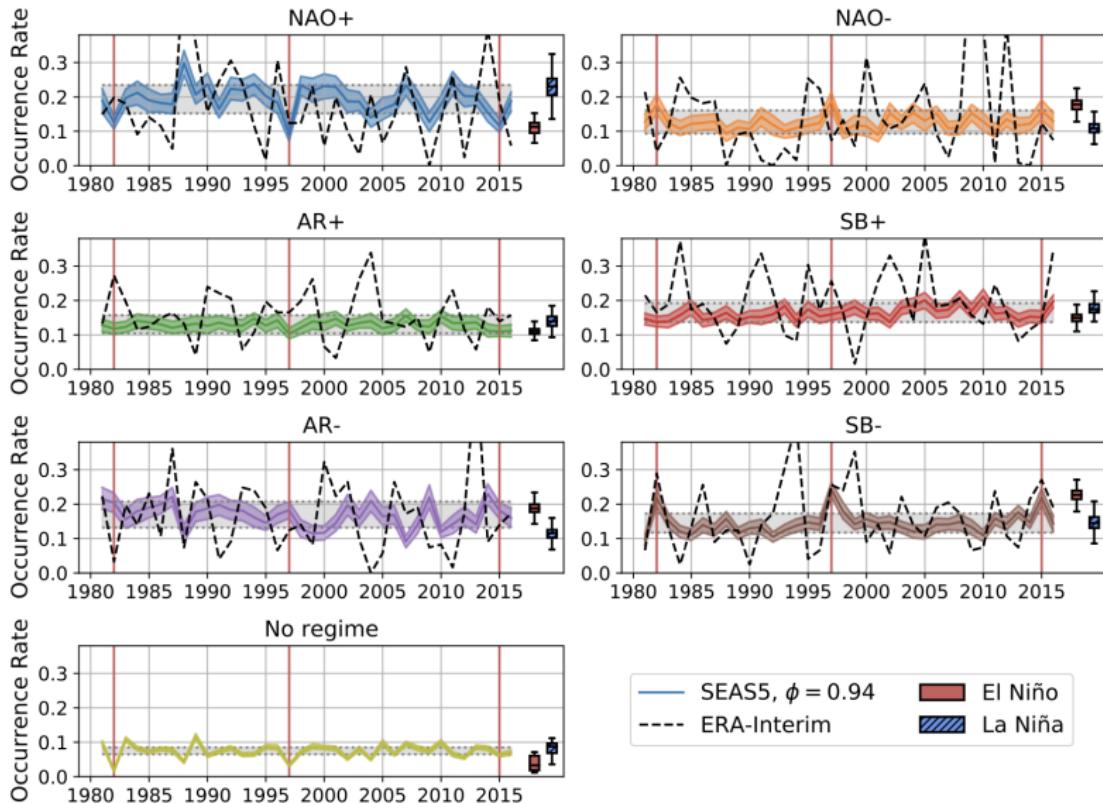
# Ensemble persistency constraint



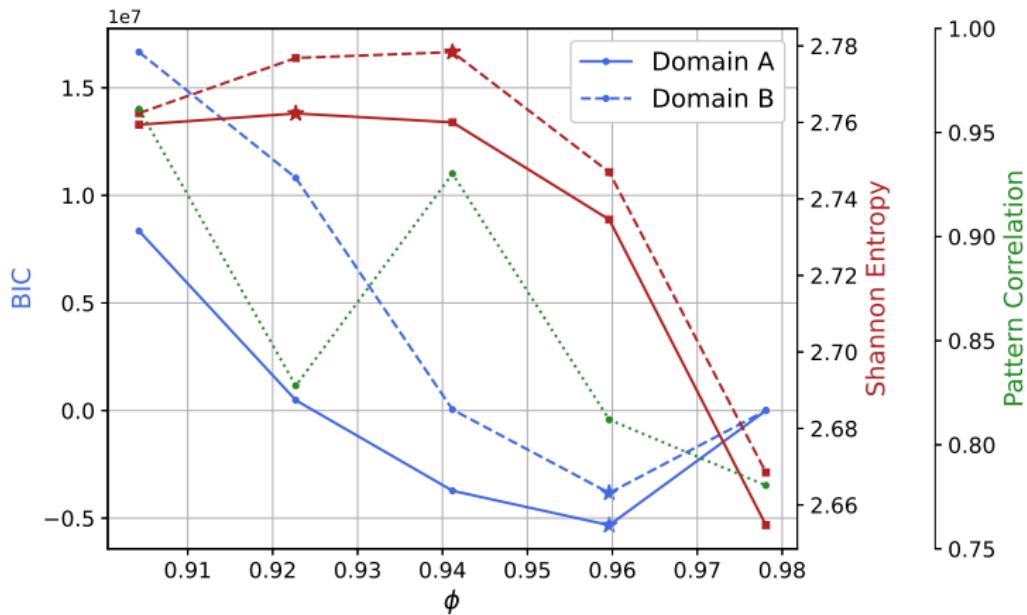
# Ensemble persistency constraint



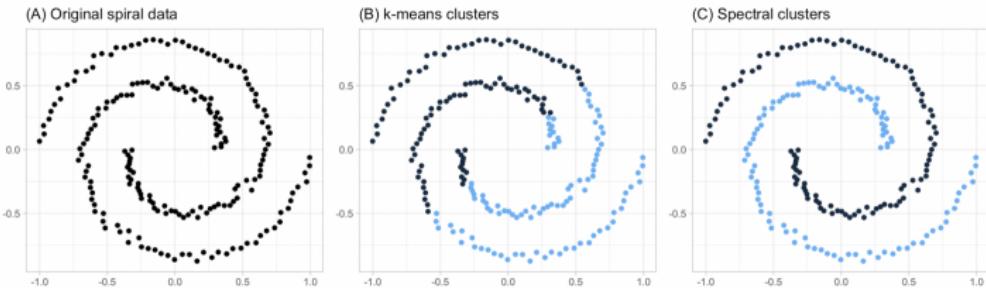
# Occurrence rates



# Optimal $\phi$



# K-Means vs Spectral Clustering



## Definition

Let  $V$  be a  $K$ -Vector space,  $f: V \rightarrow V$  an Endomorphismus,  $\lambda \in K$ . The scalar  $\lambda$  is called **Eigenvalue** of  $f$ , if there is a vector  $v \in V, v \neq 0$ , so that

$$f(v) = \lambda \cdot v.$$

The vector  $v$  is called **Eigenvector** of  $f$  an Eigenvalue  $\lambda$ .

**Note:** An Eigenvalue  $\lambda$  can be  $0 \in K$ , but an Eigenvector is always  $\neq 0$ .

## Theorem

---

### Theorem

Let  $V$  be a  $K$ -vector space,  $n = \dim V < \infty$  and  $f: V \rightarrow V$  an Endomorphismus. The following two are equivalent:

1.  $V$  has a basis of Eigenvectors of  $f$ .
2. There is a Basis  $\mathcal{B}$  of  $V$ , so that

$$M_{\mathcal{B}}^{\mathcal{B}}(f) = \begin{pmatrix} \lambda_1 & & 0 \\ & \ddots & \\ 0 & & \lambda_n \end{pmatrix} \text{ with } \lambda_i \in K.$$

# Characteristic Polynom

## Definition

Let  $A \in K^{n \times n}$  and  $\lambda \in K$  arbitrary. Then

$$\text{Eig}(A, \lambda) := \{v \in K^n \mid Av = \lambda v\}$$

is called the **Eigenspace** of  $A$  with respect to  $\lambda$ .

$$\chi_A(t) := \det(A - tE) \in K[t]$$

is called the **charakteristisches Polynom** of  $A$ .

**Remark:** For a matrix  $A \in K^{n \times n}$  the following holds:

$$\lambda \in K \text{ is an Eigenvalue of } A \Leftrightarrow \text{Eig}(A, \lambda) \neq 0.$$

## Theorem

---

Let  $A \in K^{n \times n}$  and  $\lambda \in K$ . Then:

$$\lambda \text{ is an Eigenvalue of } A \Leftrightarrow \lambda \text{ is a root of } \chi_A(t).$$

# Multiplicity

## Definition

Let  $P(t) \in K[t]$  be a Polynom.  $P(t)$  can be decomposed over  $K$  in **Linear factors** if and only if there are  $\lambda_1, \dots, \lambda_n \in K, c \in K$ , so that

$$P(t) = c \cdot (t - \lambda_1) \cdots (t - \lambda_n) = c \cdot \prod_{j=1}^r (t - \lambda'_j)^{m_j},$$

where  $m_j \in \mathbb{N}$  and  $\lambda'_1, \dots, \lambda'_r \in \{\lambda_1, \dots, \lambda_n\}$  are pairwise different.  $m_j$  is called the **Multiplicity** of the root  $\lambda'_j$ . It holds that

$$\sum_{j=1}^r m_j = n.$$

## Example

---

## Example

---

## Example

---