# The Good, the Bad and the Ugly of unsupervised learning
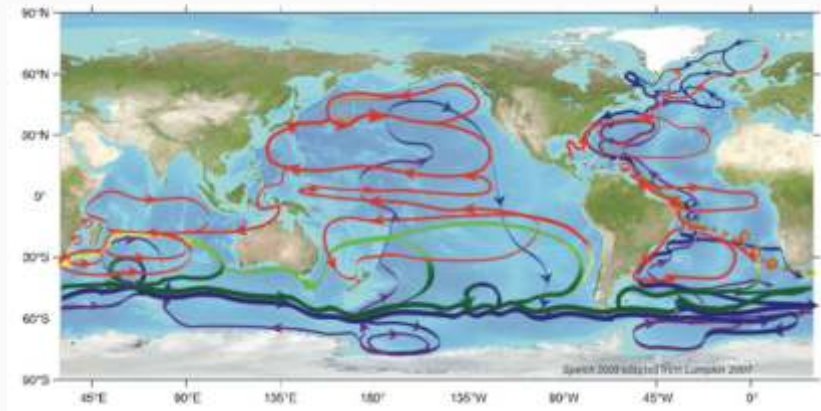
---

Maike Sonnewald[1,2]

July, 2019

[1]MIT & [2]Harvard

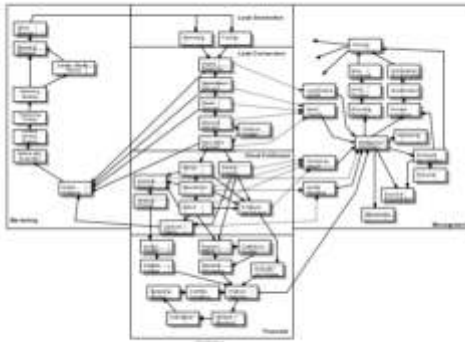**$CO_2$ release is the fastest in the past 55 million years...**



Movie by NASA

**Take home:**
> **Application of modern statistics/ML and traditional modeling/theory work can drive progress**
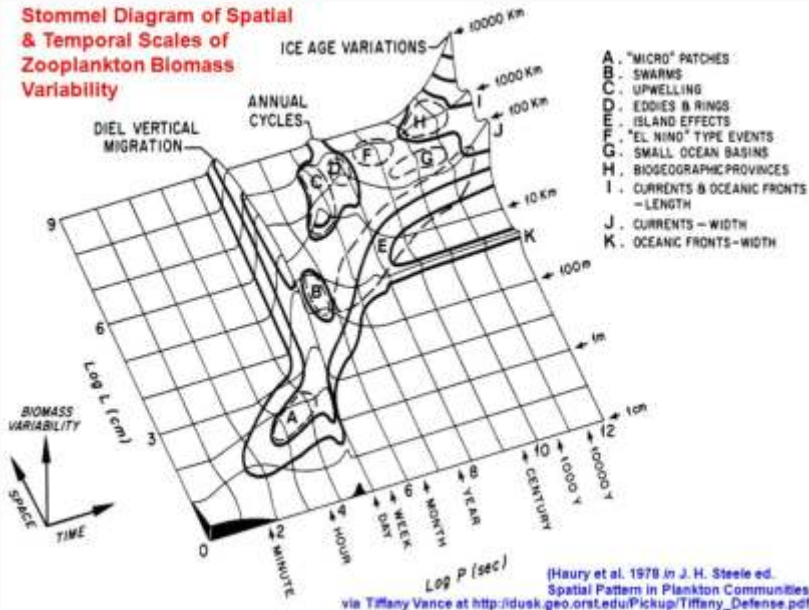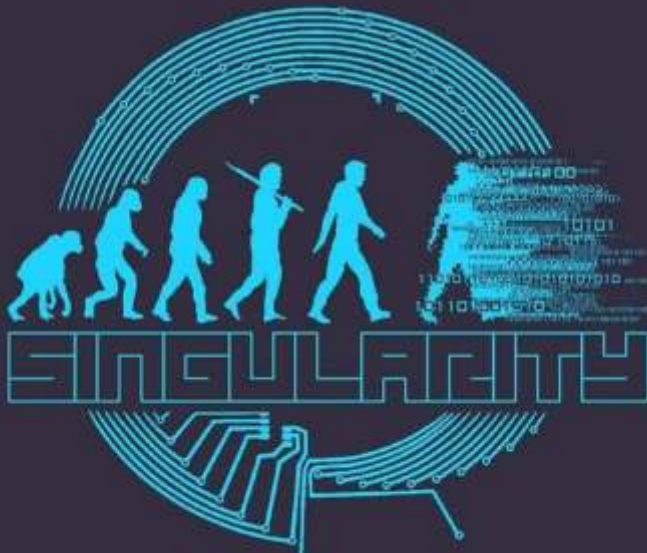
slideshare.net

Including **all** components of e.g.
an ESM at adequate resolution is unfeasible.

## How complicated is complicated enough?



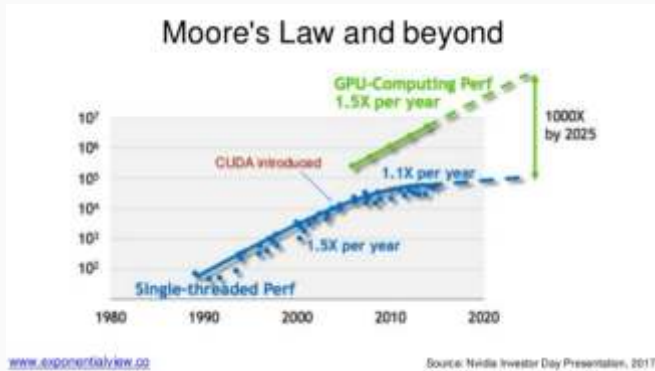**Stommel Diagram of Spatial & Temporal Scales of Zooplankton Biomass Variability**

A . "MICRO" PATCHES
B . SWARMS
C . UPWELLING
D . EDDIES & RINGS
E . ISLAND EFFECTS
F . "EL NINO" TYPE EVENTS
G . SMALL OCEAN BASINS
H . BIOGEOGRAPHIC PROVINCES
I . CURRENTS & OCEANIC FRONTS − LENGTH
J . CURRENTS − WIDTH
K . OCEANIC FRONTS − WIDTH

(Haury et al. 1978 in J. H. Steele ed. Spatial Pattern in Plankton Communities via Tiffany Vance at http://dusk.geo.orst.edu/Pickup/Tiffany_Defense.pdf)

3

Moore's Law: Jet-Pack!

**Increased computational power?**



Moore Vs Rock: Nano materials/quantum computing to the resque..?

eos.org

**The dream:**
ML can be used to lead the charge both in terms of supervised and unsupervised learning.

Can ML help?

- Make complicated data complex
- Allows insight to parameterize and simplify
- Create synergy between models, theory and observations



**ML towards the goal of science/geoscience:**
**Precise and accurate understanding of the natural world.**

"*...automate those parts that can be perfectly automated*"

Occam's razor in Machine Learning

"*...automate those parts that can be perfectly automated*"

Occam's razor in Machine Learning

- Modern data science is greatly increasing the efficiency of conventional research
- Find patterns to accelerate exploration of physics
- Highlight emergence of complex interactions



COMPLEX ≠ COMPLICATED

"Complexity: *I know it when I see it...*"

## Complexity: Find underlying "rules"
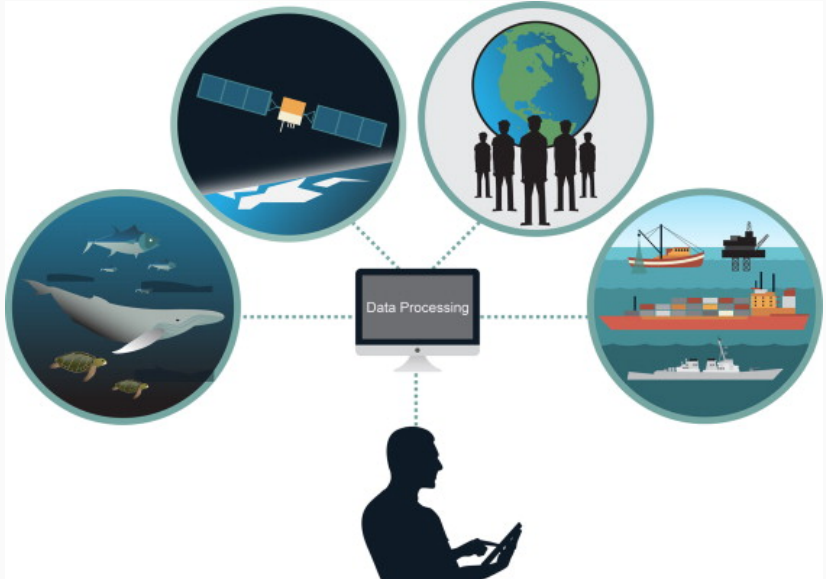
"Complexity: *I know it when I see it...*"



**BOIDS** Craig Reynolds (1986)

- separation: steer to avoid crowding local flockmates
- alignment: steer towards the average heading of local flockmates
- cohesion: steer to move towards the average position (center of mass) of local flockmates

Maxwell 2015, Kim Martini and Sonnewald et al 2013

Maxwell 2015, Kim Martini and Sonnewald et al 2013

# Unsupervised Learning: Find clusters

## The persuasive power of numbers

*"There are three kinds of lies: lies, damned lies, and statistics."*
Benjamin Disraeli (British prime minister)



- Bad statistics can bolster weak arguments
- Weak results can seem legitimate.
- False positives are bad.
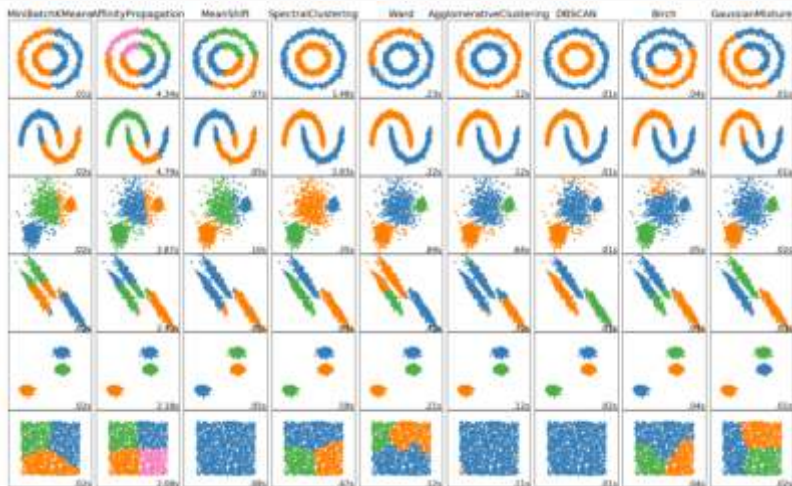
## Lies damned lies and ..ML?



**Do:**
Clean data, visualize data, highlight interactions as appropriate, choose an appropriate model, check its parameters are appropriate, account for stochastic elements.

**Don't:**
Trust data blindly, make assume variance topography, choose model blindly, if so: Brute force the statistical robustness.

**Keep It Simple Stupid**



scikit-learn.org

Libraries e.g. scikit learn, dask_ml

16

## Examples outline

- 1:      The good ...Global dynamical regimes
  - Insight into global dynamical regimes!
  - Keeping it simple
- 2:      The bad ... Global eco-provinces
  - Finding patterns in complicated data?
  - Avoiding false positives
- 3:      The robustly ugly ... Global eco-provinces
  - Complex solution to complicated data
  - Interdisciplinarity to the rescue...



17

# 1: The Good ...Global dynamical regimes
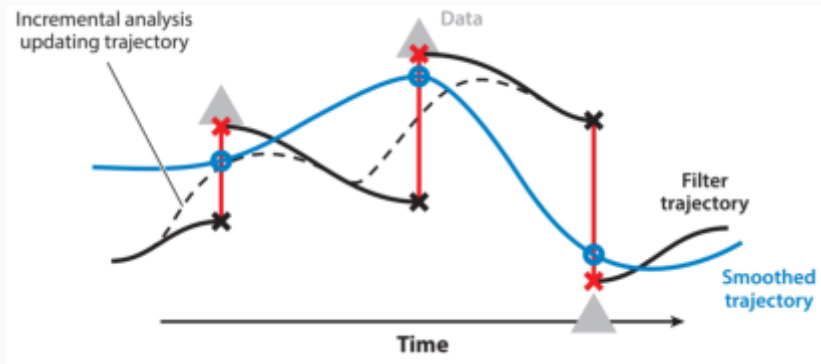
**Exploration** - Uncover new physical phenomena?

- Identify driving features of key currents

- Analyze Big Data for climate fast

e.g. Climate Model Inter-comparison Project

**Model Development** - Establish largely laminar regions?

$\rightarrow$ Save computational cost

- Make inferences of key transitions

$\rightarrow$ e.g. baroclinicity

- Focus parameterization

# Estimating the Circulation and Climate of the Ocean: 1992-2013



ECCO provides ocean state consistent with known physics and observations

## Keeping it simple: Barotropic Vorticity equation

Momentum equations:

$$\partial_t \mathbf{u} + f\mathbf{k} \times \mathbf{u} = -\frac{1}{\rho_0}\nabla p + \frac{1}{\rho_0}\partial_z \tau + \mathbf{a} + \mathbf{b}, \partial_z p = -g\rho, \nabla \cdot \mathbf{v} = 0.$$

-Depth integrate, take curl

Barotropic Vorticity:

$$0 = \overbrace{\nabla \cdot (f\mathbf{U})}^{\text{Advection}} - \underbrace{\nabla \times (p_b \nabla H)}_{\text{Bottom Pressure Troque}} + \overbrace{\nabla \times \tau}^{\text{Wind and Bottom stress}} - \underbrace{\nabla \times \mathbf{A}}_{\text{Non-linear Torque}} + \overbrace{\nabla \times \mathbf{B}}^{\text{Lat. Visc.}}$$
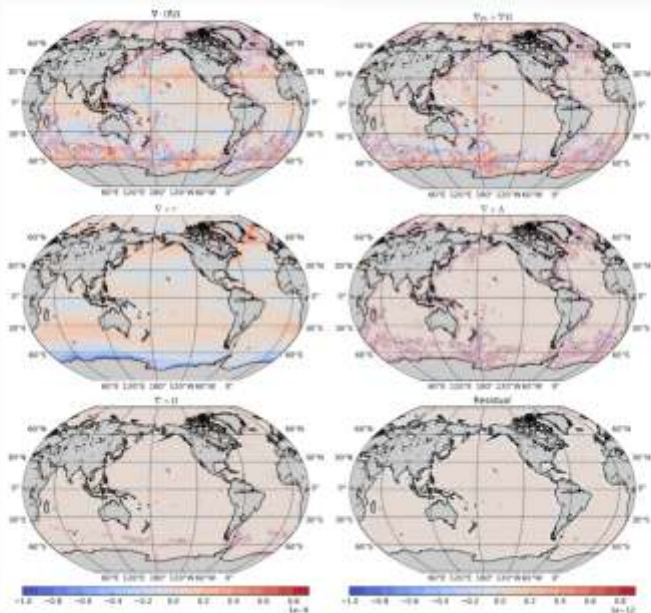
e.g. Sverdrup balance:

Wind stress curl and advection balance locally

$$\nabla \cdot (f\mathbf{U}) = \nabla \times \tau$$

## Keeping it simple: Barotropic Vorticity equation

Momentum equations:

$$\partial_t \mathbf{u} + f\mathbf{k} \times \mathbf{u} = -\frac{1}{\rho_0}\nabla p + \frac{1}{\rho_0}\partial_z \tau + \mathbf{a} + \mathbf{b}, \partial_z p = -g\rho, \nabla \cdot \mathbf{v} = 0.$$

-Depth integrate, take curl

Barotropic Vorticity:

$$0 = \overbrace{\nabla \cdot (f\mathbf{U})}^{\text{Advection}} - \underbrace{\nabla \times (p_b \nabla H)}_{\text{Bottom Pressure Torque}} + \overbrace{\nabla \times \tau}^{\text{Wind and Bottom stress}} - \underbrace{\nabla \times \mathbf{A}}_{\text{Non-linear Torque}} + \overbrace{\nabla \times \mathbf{B}}^{\text{Lat. Visc.}}$$

e.g. Sverdrup balance:

Wind stress curl and advection balance locally

$$\nabla \cdot (f\mathbf{U}) = \nabla \times \tau$$

*..Is this real? Is every location unique?*

20

What do we do when we don't have the answers a priori?



number of clusters    number of cases    centroid for cluster $j$

case $i$

objective function $\leftarrow$ $J = \sum_{j=1}^{k} \sum_{i=1}^{n} \left\| x_i^{(j)} - c_j \right\|^2$
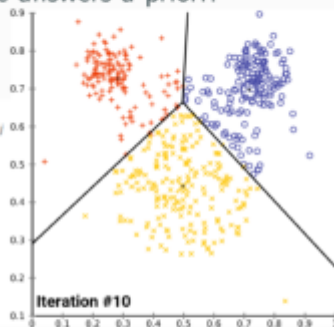
Distance function

commons.wikimedia.org

- **Assign**: Cluster w mean minimizing least squared Euclidean dist.
- **Itterate**: Calculate the new means.
- NB! NP-hard. Not global. Need to treat data. Sensitive to:
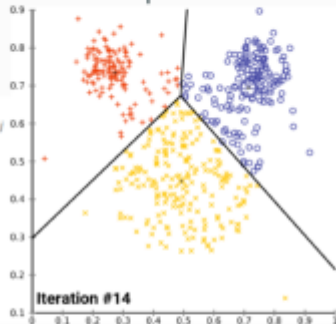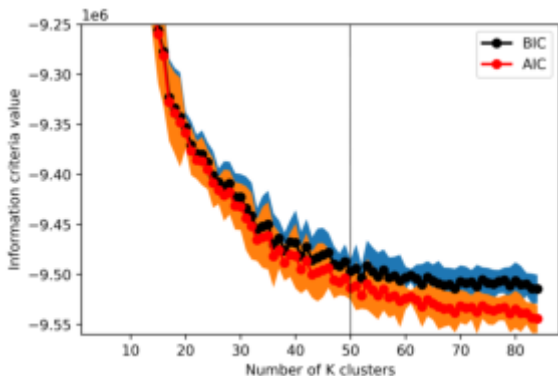K and initialization.

What do we do when we don't have the answers a priori?



$$J = \sum_{j=1}^{k} \sum_{i=1}^{n} \left\| x_i^{(j)} - c_j \right\|^2$$

objective function, number of clusters, number of cases, case $i$, centroid for cluster $j$, Distance function

commons.wikimedia.org

- Assign: Cluster w mean minimizing least squared Euclidean dist.
- Itterate: Calculate the new means.
- NB! NP-hard. Not global. Need to treat data. Sensitive to:
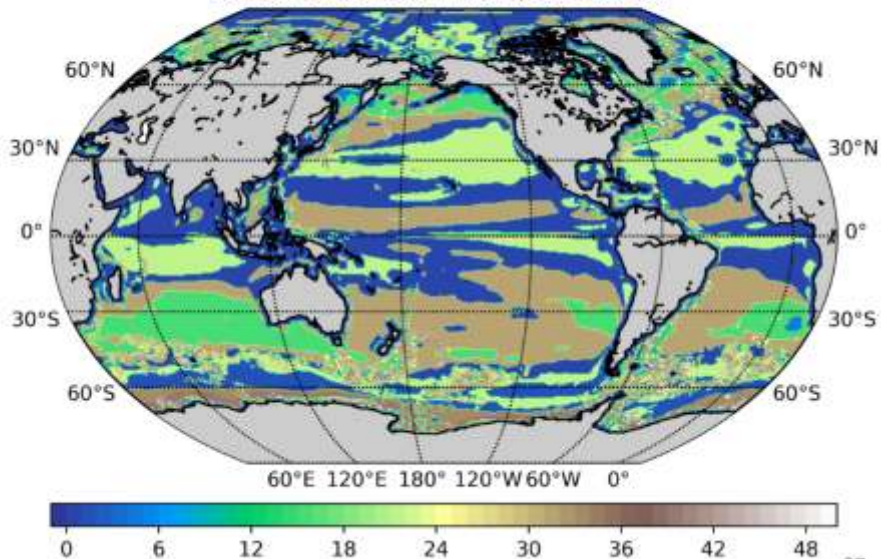K and initialization.

23

What do we do when we don't have the answers a priori?



number of clusters     number of cases

case $i$

centroid for cluster $j$

objective function $\leftarrow J = \sum_{j=1}^{k} \sum_{i=1}^{n} \left\| x_i^{(j)} - c_j \right\|^2$

Distance function

Iteration #10

commons.wikimedia.org

- Assign: Cluster w mean minimizing least squared Euclidean dist.
- Itterate: Calculate the new means.
- NB! NP-hard. Not global. Need to treat data. Sensitive to: K and initialization.

24

What do we do when we don't have the answers a priori?



$$J = \sum_{j=1}^{k} \sum_{i=1}^{n} \left\| x_i^{(j)} - c_j \right\|^2$$

commons.wikimedia.org

- Assign: Cluster w mean minimizing least squared Euclidean dist.
- Itterate: Calculate the new means.
- NB! NP-hard. Not global. Need to treat data. Sensitive to: K and initialization.

25

## Choosing K: Information criteria



$$\mathrm{BIC} = K \ln(n) - 2\ln(\mathcal{L}),$$

where $n$ is the number of datapoints and $\mathcal{L}$ is the likelihood:

$$\mathcal{L} = \Pi_{i=1}^{N} \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(\zeta_i - \hat{\zeta}_i)^2}{2\sigma^2}\right).$$

# Dynamical regions



K-Means with 50 clusters, scaled data

# Percentage of area accounted for

# Sverdrup balance: Subtropical Gyre

## 11.2%



Present globally!

Only here does the linearity approximation break down.

Ocean globally organizes into distinct 6 regimes
→ Use this to analyze Big Data **fast**

- Objectively describe 3D ocean
  using k-means
- Regions comply with theory
- Use IC and check for degeneracy



**Sonnewald, Wunsch and Heimbach (2019)** and
**Sonnewald et al. (Nature Com. resubmission)**

# 2: The Bad    ...Global ecological provinces

It's complicated

Dutkiewicz et al., 2015

## Biogoechemical ecosystem model: Empirical modeling approach



| Trait | Size-based | Functionality-based |
|---|---|---|
| Maximum growth rate | Yes | Yes |
| Nutrient uptake half saturation | Yes | No |
| Sinking rate | Yes | No |
| Light Absorption | Yes | Yes |
| Palatability | Yes | Yes |
| Grazing rate | Yes | No |
| Other mortality | No | No |
| Carbon quota | Yes | No |
| Stoichiometry | No | Yes |

Dutkiewicz et al., 2015

# 51 species (biomass) and 4 nutrients (concentration)

Biomass: K-Means with 250K, z=0, all year

36

AIC penalizes more complex models, i.e., models with additional

# cbiomes K-Means: Summer/Winter

The data is clearly not like the BV
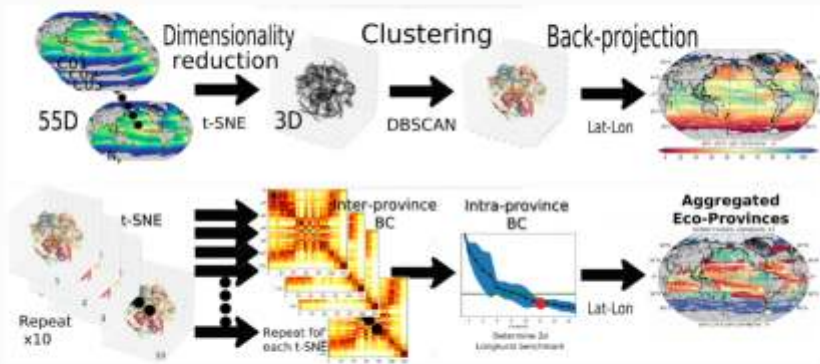$\rightarrow$ We do **not** have round ND distributions

- The yearly "provinces" looked reasonable
- Winter/Summer split was odd
- AIC+BIC confirmed suspicion

# 3: The (robustly) Ugly    ...Global ecological provinces

Given a set of N high-dimensional objects $x_1, ..., x_N$, the t-Statistic Neighbourhood Embedding minimize Kullbach-Leibner distance between the likelyhood of association between a low dimentional rendition and the high dimentional data.

- If $x_i$ it the i-th object in the N dim space and $y_i$ is the i-th object in the low-dim space:

$$p_{j|i} = \frac{\exp(-\|\mathbf{x}_i - \mathbf{x}_j\|^2 / 2\sigma_i^2)}{\sum_{k \neq i} \exp(-\|\mathbf{x}_i - \mathbf{x}_k\|^2 / 2\sigma_i^2)},$$

and the same for a reduced dimensional set:

$$q_{ij} = \frac{(1 + \|\mathbf{y}_i - \mathbf{y}_j\|^2)^{-1}}{\sum_{k \neq l} (1 + \|\mathbf{y}_k - \mathbf{y}_l\|^2)^{-1}}.$$

This is done as: $KL(P\|Q) = \sum_{i \neq j} p_{ij} \log \frac{p_{ij}}{q_{ij}}$
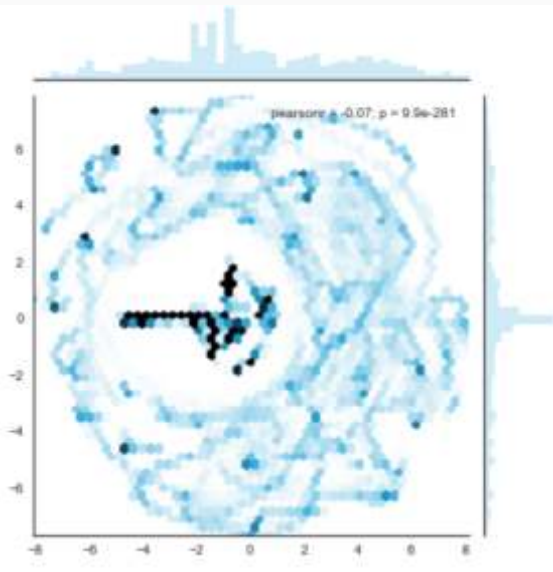
tSNE dimensions

stackexchange.org
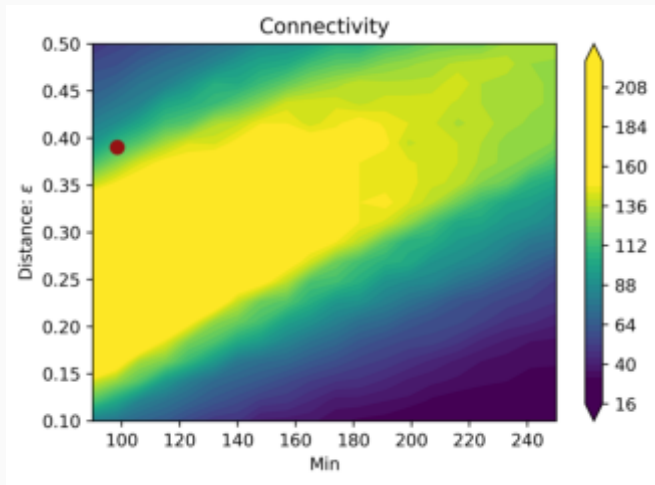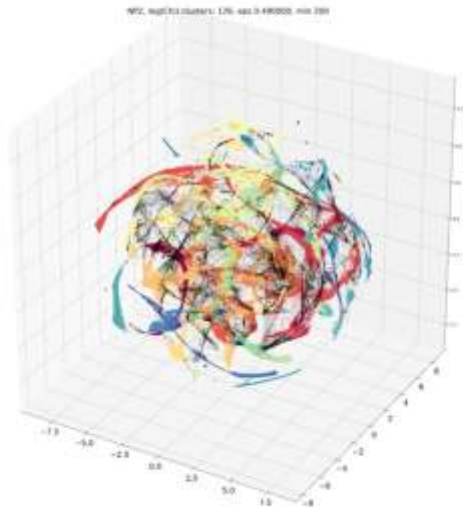
- Set: Eps and MinPts.
- Note: Not stochastic.
- Global. Need to preprocess data.

# Unsupervised learning: DBSCAN



2D "elbow" check in connectedness
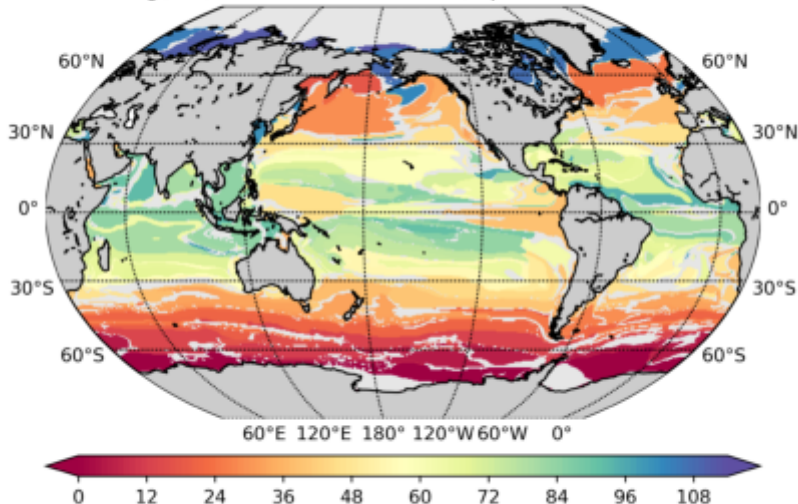
# Biogeography: Clustering Eco-Provinces



NPZ log(Chl) DBSCAN clusters: 115, eps 0.390000, min 100
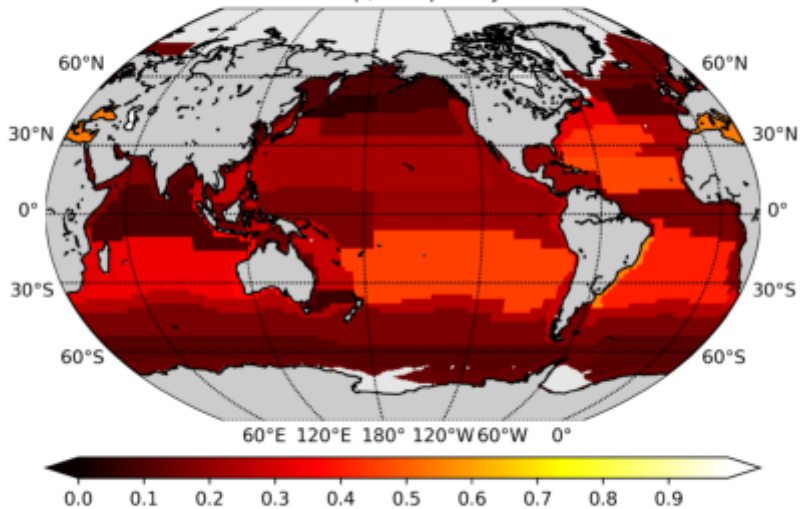
Do we need all provinces? Are we doing better than Longhurst?

Mean Map, complexity 131

## NPZ complexity: Bray-Curtis dissimilarity

How similar are the identified clusters?

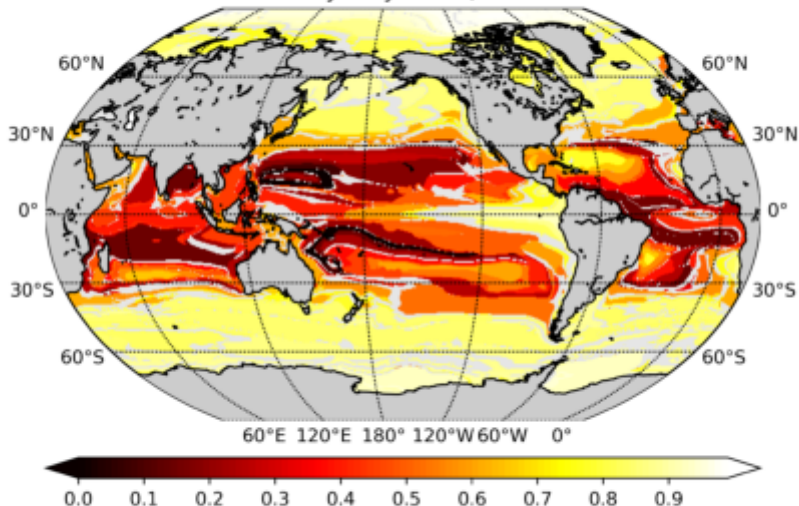$$BC_{ij} = 1 - \frac{2C_{ij}}{S_i + S_j}$$

$C_{ij}$ is the minimum value present where similarities exist:

$$C_{ij} = \sum_{c51}^{c01} min(\mathrm{biomass}_i, \mathrm{biomass}_j)$$

$S_i$ is the total across plankton: $S_i = \sum_{c51}^{c01}(\mathrm{biomass}_i)$

Dissimilarity Bray-Curtis, area: 1.7%
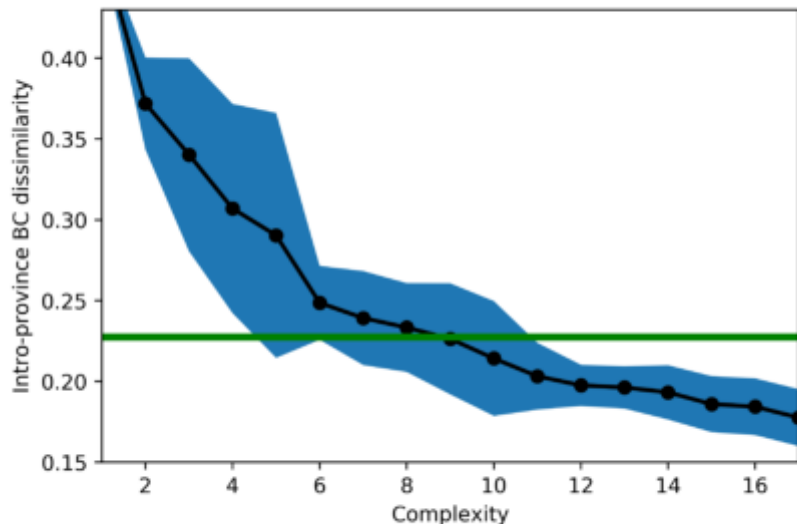
Many of eco-provinces are very similar

- Each eco-province connected to every other
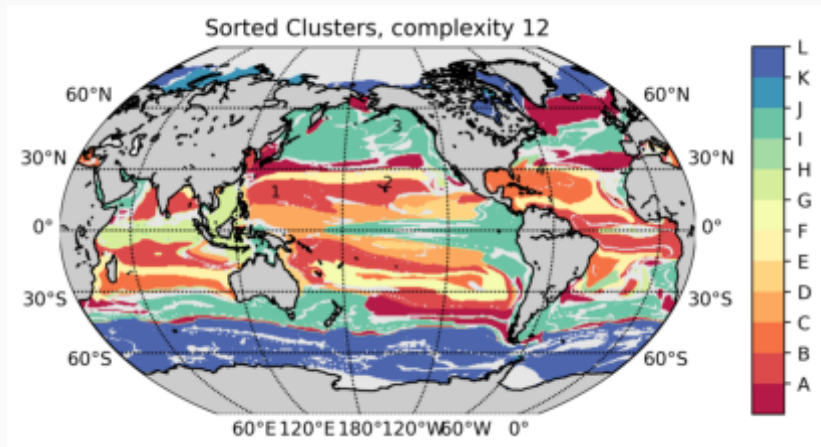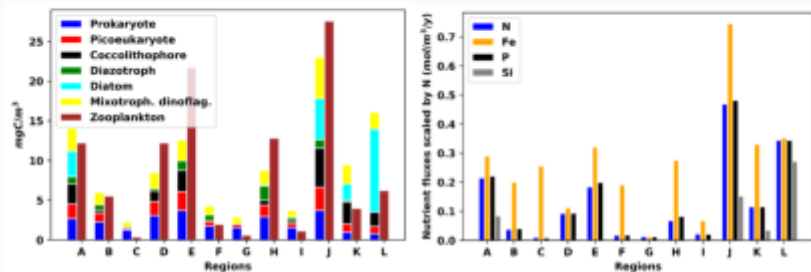- Facebook allegory
- Graph theory allows sorting

A minimum complexity of 12 is recommended
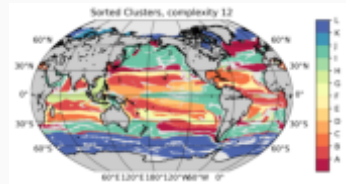
Sorted Clusters, complexity 12

- Similar biomass/chl but different community structure
- Biomass is poor predictor of zooplankton: Trophic cascades?

Visualizing the data in 3D allowed model selection
$\rightarrow$ Knowledge of topology avoids brute-force

- Robust eco-provinces
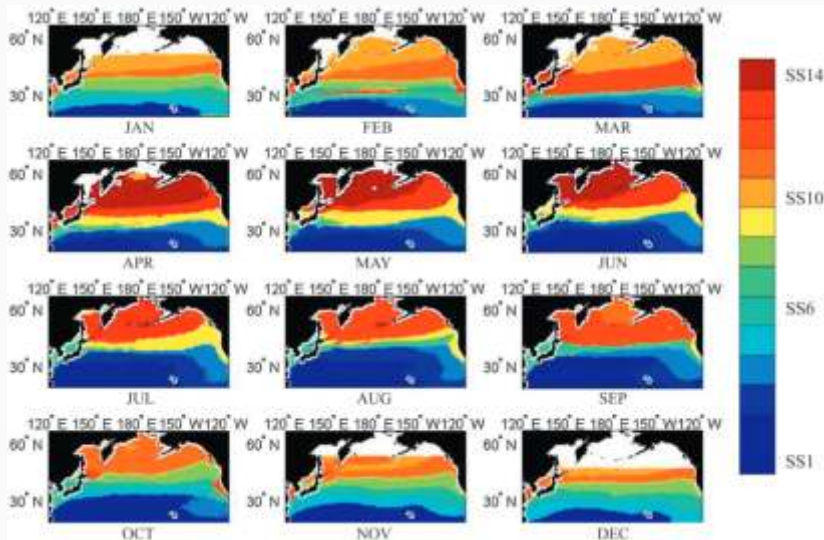- Improvements over Longhurst
- Aggregation allowed wider application



Sorted Clusters, complexity 12

How complicated is complicated "enough"?

How complicated is complicated "enough"? What does
conservation work require?
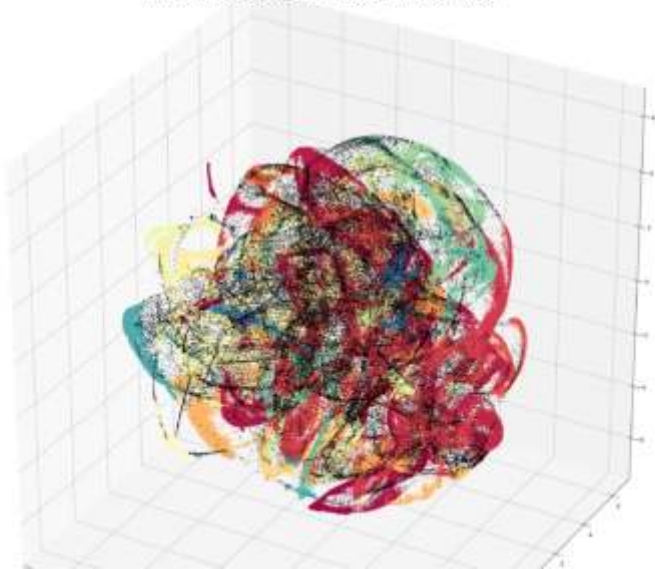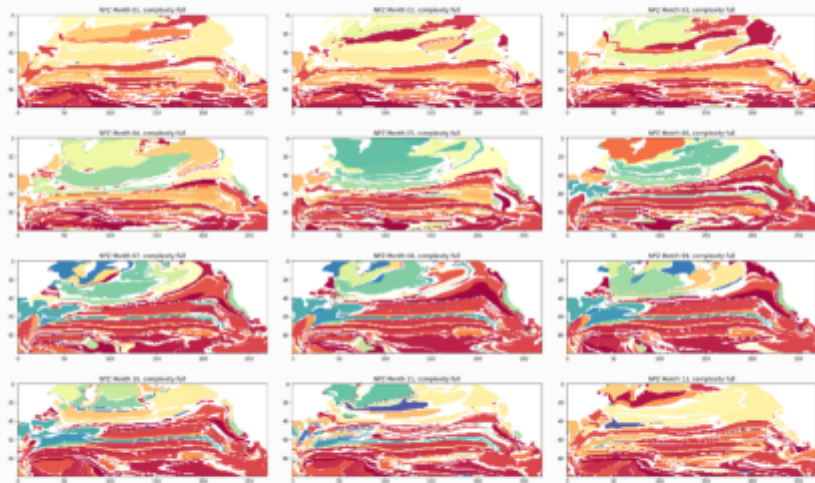
Kavenaugh et al. 2016

SOM+hierarchical clustering: SST, PAR and Chl a

NPZ Pacific, log(biomass) clusters: 159, eps 0.255000, min 80

Similarities when looking at ecology!
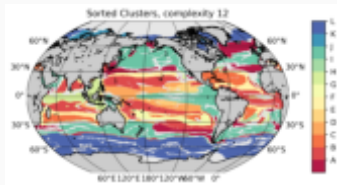
# Summary

## Summary

The singularity may not be coming; Make complicated models more complex?

- To allow unsupervised ML to help:

    - Preprocess data as appropriate
    - Start with simple model
    - Carefully estimate parameters
    - Assess with system insight

- System insight to simplify the complicated components:

    - Ecosystem models?
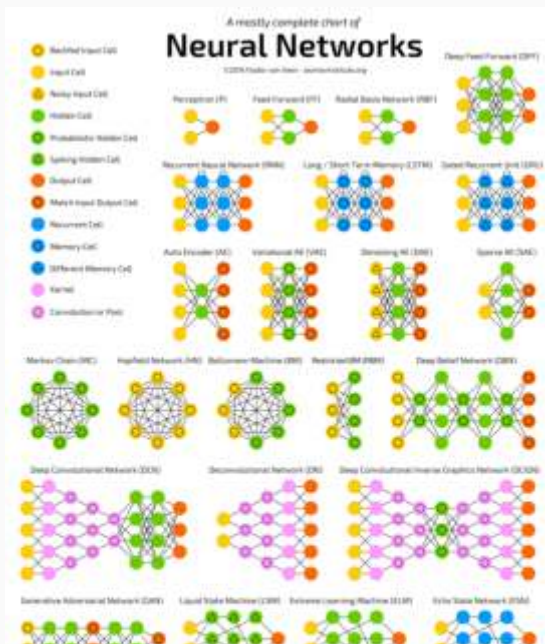    - Bathymetric interactions
    - 3D ocean insight

Modeling+Theory+Observations = Understanding of natural world
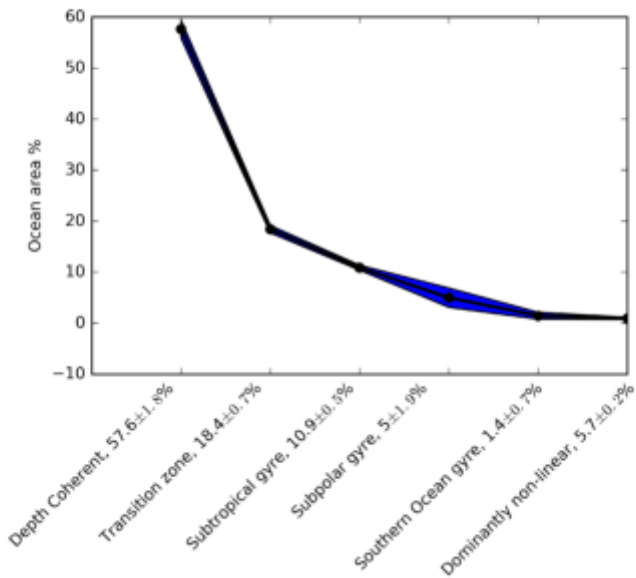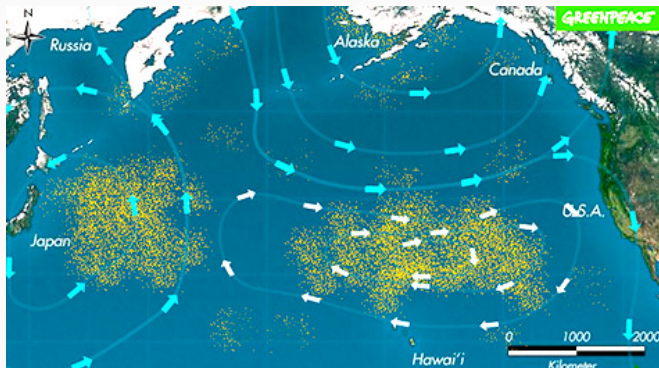
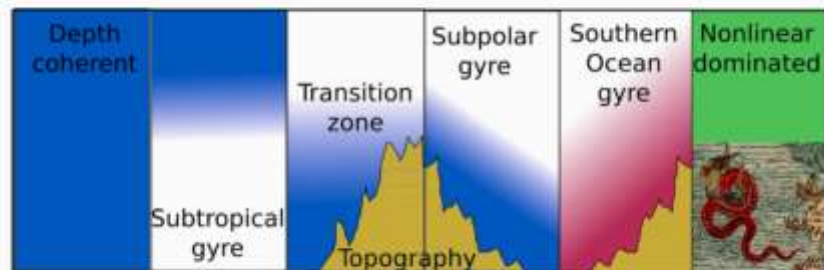Sorted Clusters, complexity 12

Thank you!

# Percentage of area accounted for

Garbage
follows currents. Ecology is not this simple...
Good to know if you want to clean it up with food webs in mind!

## Implications: Gyre ⇔ overturning?

kMeans_50_NAtlFit.png

## Non-linear terms

$$\nabla \times \mathbf{A} = \nabla \times \left[ \int_{-H}^{\eta} \nabla \cdot (\mathbf{uu}) \mathrm{d}z \right] + [w\zeta]_{z=H}^{z=\eta} + [\nabla w \times \mathbf{u}]_{z=H}^{z=\eta}, \quad (1)$$

where **uu** is a second order tensor. The RHS of equation (1) represents the curl of the vertically integrated momentum flux divergence, the non-linear contribution to vortex tube stretching and the conversion of vertical shear to barotropic vorticity.