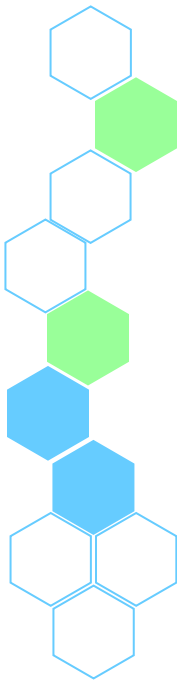
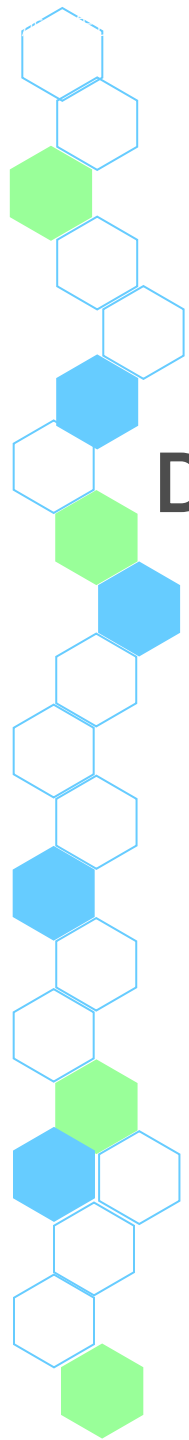
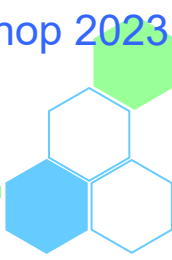


Regression Tree and Clustering for Distributions, and Homogeneous Structure of Population Characteristics

Mihoko Minami (Keio University)

Joint work with Cleridy E. Lennert-Cody of Inter-American Tropical Tuna Commission (IATTC)





Introduction

We often collect samples on characteristics of different observation units and wonder **Whether the characteristics of the observation units have similar distributional structure?**

We consider methods to find homogeneous subpopulations

- using **regression tree** and **clustering for distributions** approaches
- based on a **modified Jensen–Shannon divergence**

and present

- a testing procedure **for homogeneity of a cluster** and
- a hierarchical testing procedure to find **the minimal homogeneous/near-homogeneous tree structure** of the distributions of a population characteristic.

Motivational Example

Yellowfin tuna fork length data

- collected from the tuna catch of purse-seine vessels operated
- in the eastern Pacific Ocean during 2003 – 2007

A total of **797 samples** were available.

Each sample contains

- **the fork lengths (cm)** of about **50 yellowfin tunas**, and
- **the date and the location** of the fishing operations associated with the tuna catch

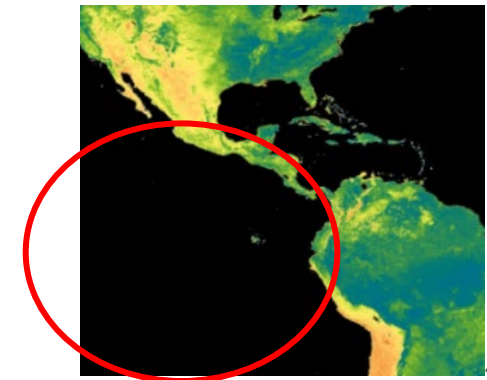
Yellowfin tuna (キハダマグロ)



Purse-seine set (巻網漁)



The eastern Pacific Ocean





The fork length data

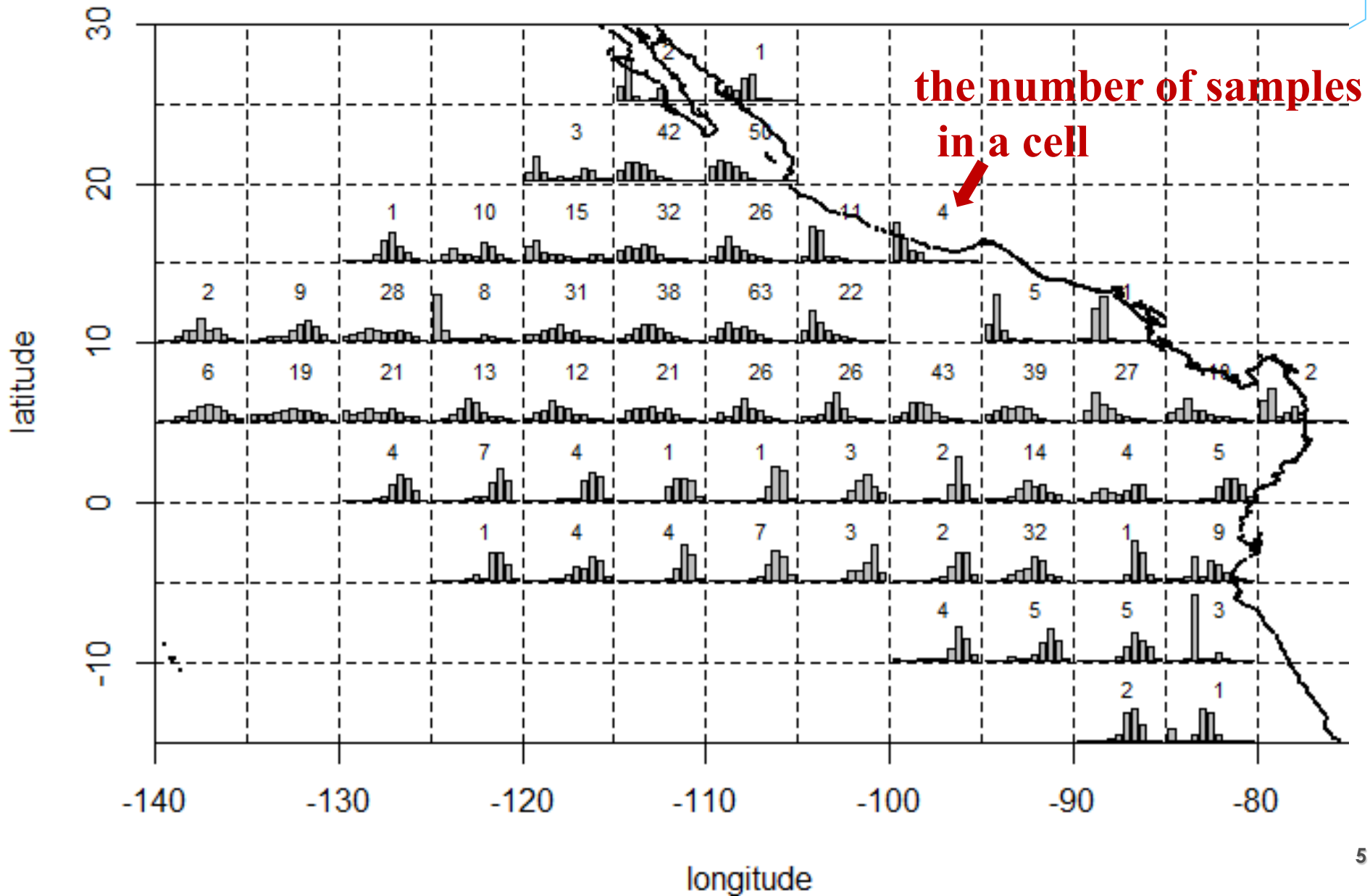
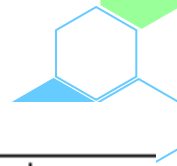
Data on the area and time period corresponding to the fishing operations are obtained from data recorded by onboard observers or from fishermen's logbooks.

The samples were collected by the port-sampling program of the Inter-American Tropical Tuna Commission (IATTC; www.iattc.org), which is the regional fishery management agency responsible for the conservation of tuna and other marine resources in the eastern Pacific Ocean

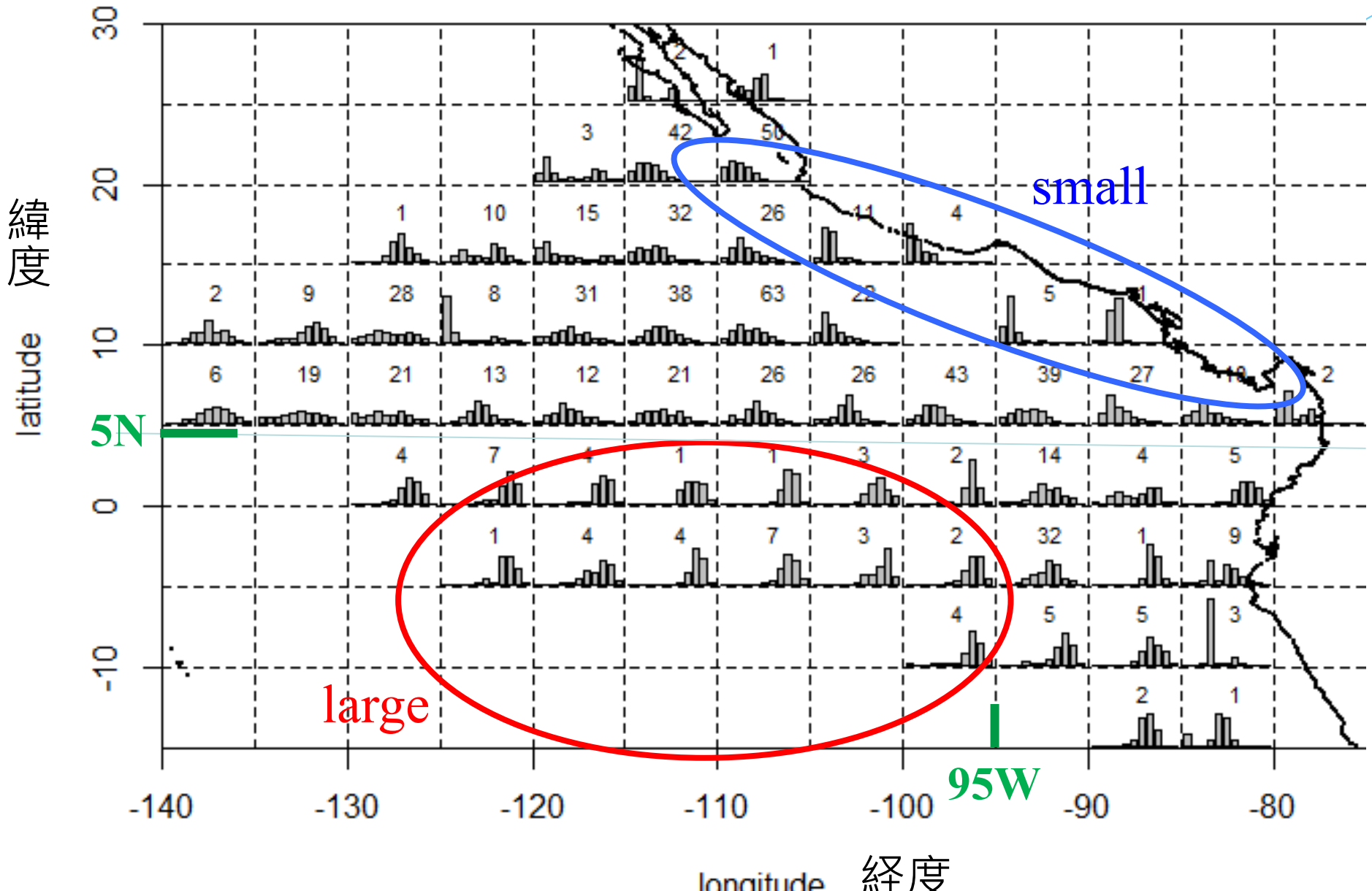
In our analysis, the fork length data were aggregated by location

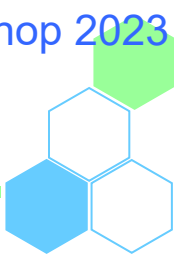
- over time **into 5° by 5° cells**
- so that 797 samples were **combined into 60 spatial cells.**

Histograms and numbers of samples in cells



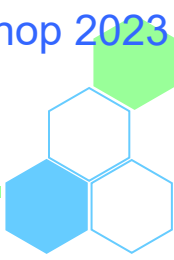
Histograms and numbers of samples in cells



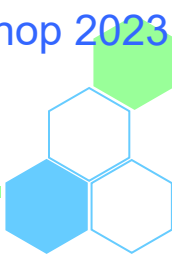


Outline of Talk

- ◆ Regression Tree and Hierarchical Clustering
- ◆ Modified Jensen–Shannon divergence, related distance and Impurity measures
- ◆ Regression Tree and Clustering for the Yellowfin tuna fork length data
 - ◆ With histograms
 - ◆ With density estimates
- ◆ Testing procedures for homogeneity and the minimum homogeneous tree structure
 - ◆ Near–homogeneous tree structure
- ◆ Summary and Future Work



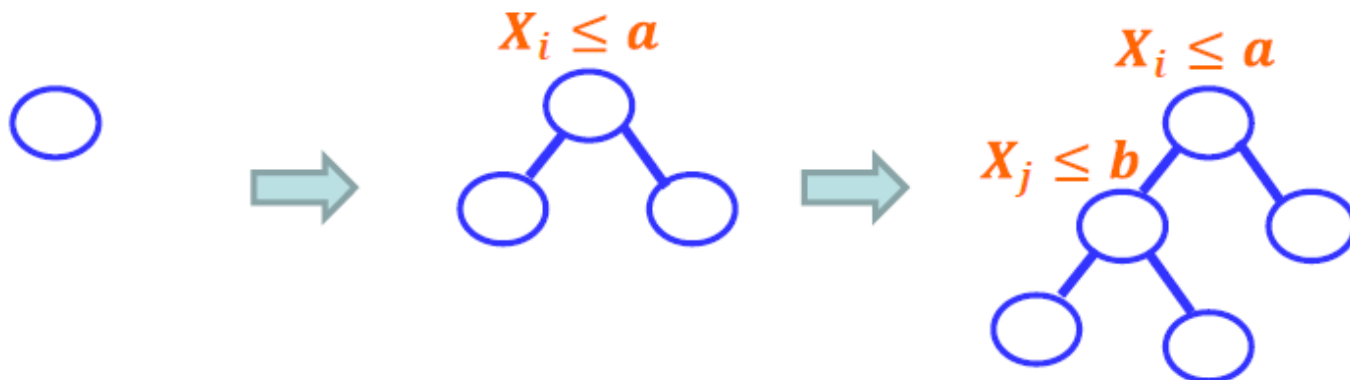
Regression Tree and Hierarchical Clustering

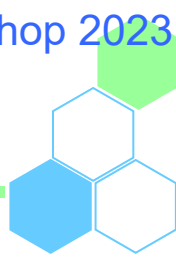


Regression Tree

Classification and regression tree (CART, Breiman et al. (1984))

- starts from a set of all units and
- repeatedly subdivides that set using **binary partitions** defined by **the values of an explanatory variable** selected to provide the greatest decrease of the values of a response variable in a measure of **impurity** until all divided nodes satisfy the termination rule.

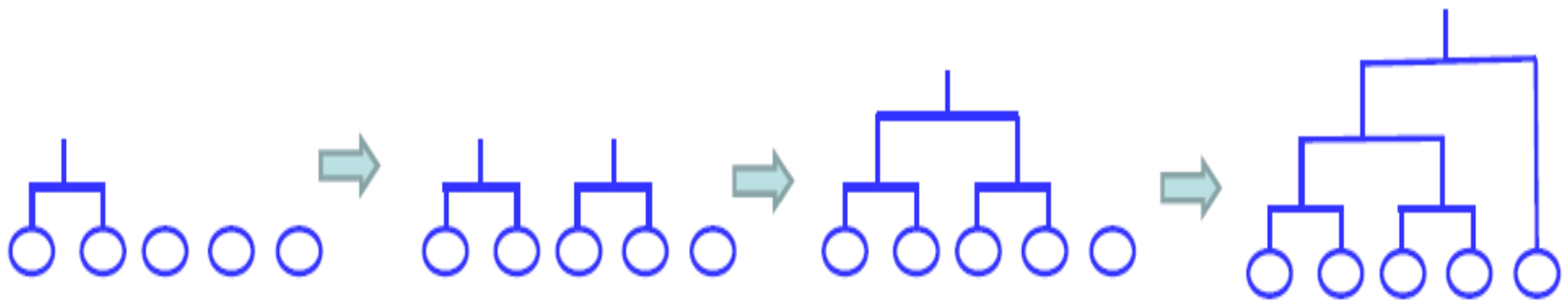


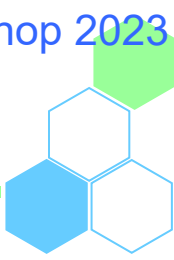


Hierarchical Clustering

Hierarchical clustering (cf. Gordon (1999))

- It is an **agglomerative** approach
- Each unit starts in its own cluster,
 - The method repeatedly **combines the two closest clusters** by some metric for distance among units
 - At the end, all units form **one large cluster**.





Modified Jensen–Shannon divergence, related distance, and impurity measures



Distance/Similarity Measures between distributions

- ◆ S.H. Cha (2007) listed a total of 42 measures in 7 types:
1) L_p Minkowski type, 2) L_1 family, 3) Intersection family, 4) Inner Product family, 5) Fidelity family or Squared-chord family, 6) Squared L_2 family or χ^2 family, 7) Shannon's entropy family
- ◆ **The modified Jensen-Shannon Divergence** (Dhillon et al. 2003)
- ◆ Clustering of histograms using **Wasserstein metric** (Ispiro and Lechevallier, 2006, Ispiro et al. 2014)
- ◆ k-Means using **Mixed α -Divergences** (Nielsen et al. 2014)
- ◆ Fuzzy clustering using **L_1 measure** (Phamtoan et al., 2022, Nguyen-Trang et al. 2023)
- ◆ **Earth Mover's distance** (Henderson et al., 2015)



The Modified Jensen–Shannon divergence

Modified Jensen–Shannon divergence (distance)

For distributions f_1 and f_2 with confidences m_1 and m_2 (> 0), respectively, let $\bar{f}_{\{1,2\}}$ be their weighted average distribution

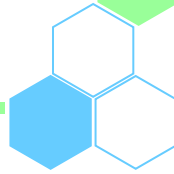
$$\bar{f}_{\{1,2\}} = \frac{1}{m_1 + m_2} (m_1 f_1(x) + m_2 f_2(x)).$$

Then, the modified Jensen–Shannon divergence is defined as

$$D_{\text{MJS}}((f_1, m_1), (f_2, m_2)) = m_1 \text{KL}(f_1 | \bar{f}_{\{1,2\}}) + m_2 \text{KL}(f_2 | \bar{f}_{\{1,2\}})$$

where
$$\text{KL}(f|g) = \int_{\Omega} f(x) \log \frac{f(x)}{g(x)} dx$$

(cf. Dhillon et al., 2003)



The Modified Jensen–Shannon divergence

- ◆ It is symmetric

$$D_{\text{MJS}}((f_1, m_1), (f_2, m_2)) = D_{\text{MJS}}((f_2, m_2), (f_1, m_1))$$

- ◆ No support problem arises since

$$\{x: \bar{f}_{\{1,2\}}(x) > 0\} = \{x: f_1(x) > 0\} \cup \{x: f_2(x) > 0\}$$

- ◆ It can be expressed with the information entropy

$$\begin{aligned} D_{\text{MJS}}((f_1, m_1), (f_2, m_2)) \\ = (m_1 + m_2) H(\bar{f}_{\{1,2\}}) - m_1 H(f_1) - m_2 H(f_2) \end{aligned}$$

where $H(\cdot)$ is the information entropy,

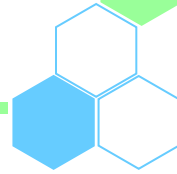
$$H(f) = - \int_{\Omega} f(x) \log f(x) dx$$

- ◆ In the case of multinomial distributions (and histograms),

$D_{\text{MJS}}((\hat{f}_1, m_1), (\hat{f}_2, m_2))$ is the log-likelihood ratio.

Impurity of a group of distributions

(Lennert–Cody et. Al, 2010, 2013)



At each step of **CART**, the binary split that produces the largest decrease in impurity is chosen.

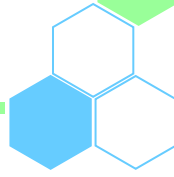
Impurity of a group of distributions

For a group of distributions $\mathcal{G} = \{(f_i, m_i), i \in G\}$, let m_G and \bar{f}_G be its confidence and weighted average distribution, respectively,

$$m_G = \sum_{i \in G} m_i, \bar{f}_G = \frac{\sum_{i \in G} m_i f_i}{m_G}.$$

We define the KL-impurity of $\mathcal{G} = \{(f_i, m_i), i \in G\}$ as

$$\text{Imp}_{\text{KL}}(\mathcal{G}) = \sum_{i \in G} m_i \text{KL}(f_i | \bar{f}_G)$$



An Expression with the information entropy

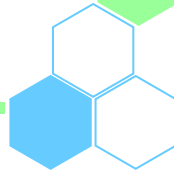
KL-impurity can be expressed with the information entropy

$$\text{Imp}_{\text{KL}}(\mathcal{G}) = m_{\mathcal{G}} H(\bar{f}_{\mathcal{G}}) - \sum_{i \in \mathcal{G}} m_i H(f_i)$$

because

$$\begin{aligned} \text{Imp}_{\text{KL}}(\mathcal{G}) &= \sum_{i \in \mathcal{G}} m_i \text{KL}(f_i | \bar{f}_{\mathcal{G}}) \\ &= \sum_{i \in \mathcal{G}} m_i \int_{\Omega} f_i(x) \log \frac{f_i(x)}{\bar{f}_{\mathcal{G}}(x)} dx \\ &= - \int_{\Omega} \sum_{i \in \mathcal{G}} (x) m_i f_i(x) \log \bar{f}_{\mathcal{G}}(x) dx \\ &\quad + \sum_{i \in \mathcal{G}} \int_{\Omega} m_i f_i(x) \log f_i(x) dx \\ &= m_{\mathcal{G}} H(\bar{f}_{\mathcal{G}}) - \sum_{i \in \mathcal{G}} m_i H(f_i) \end{aligned}$$

Distance between two groups of distributions



Hierarchical clustering repeatedly combines two closest clusters as measured by some measure of “distance”

Distance between two groups of distributions

$$D_{\text{MJS}}(\mathcal{G}_1, \mathcal{G}_2) = D_{\text{MJS}}\left(\left(\bar{f}(\mathcal{G}_1), \sum_{i \in \mathcal{G}_1} m_i\right), \left(\bar{f}(\mathcal{G}_2), \sum_{i \in \mathcal{G}_2} m_i\right)\right)$$

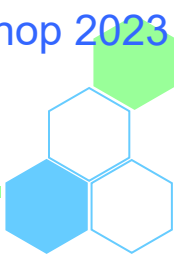
where

$$\bar{f}_{\mathcal{G}_j} = \frac{\sum_{i \in \mathcal{G}_j} m_i f_i}{m_{\mathcal{G}_j}}, j = 1, 2$$

It can be shown that

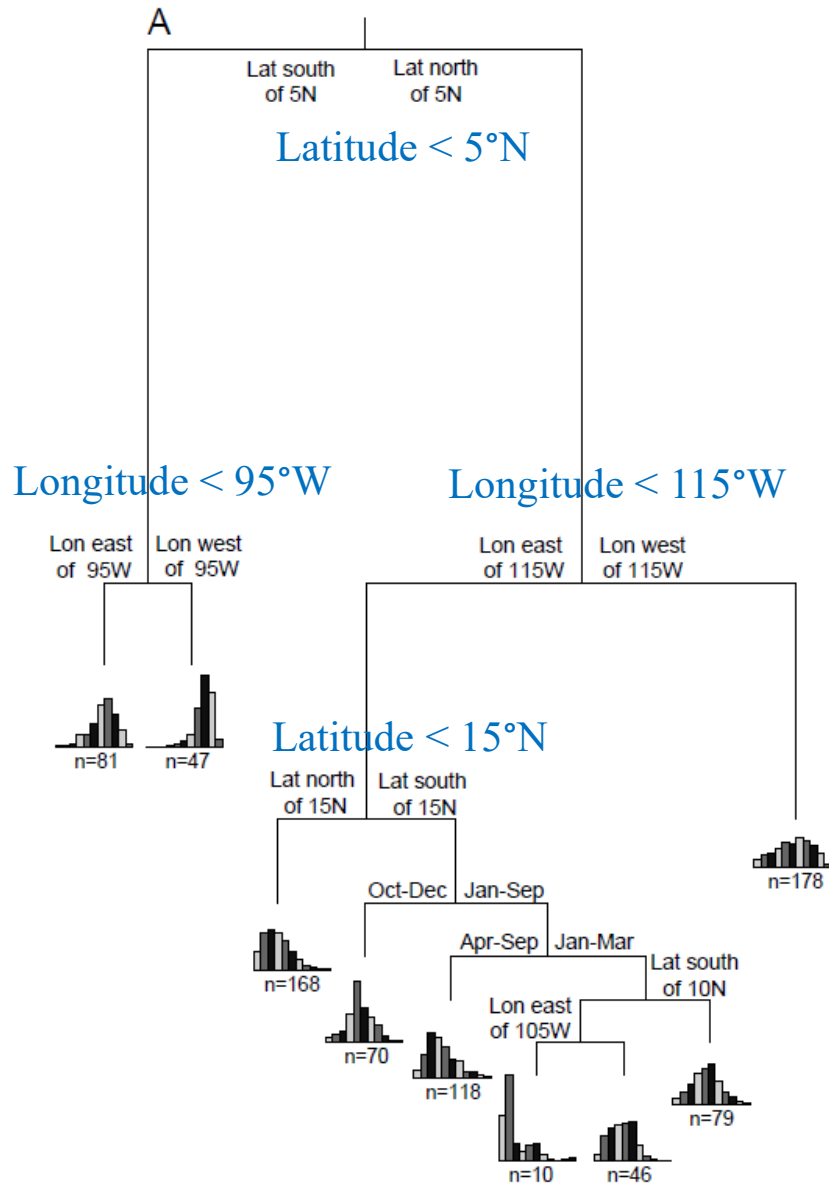
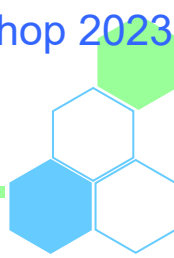
$$\begin{aligned} D_{\text{MJS}}(\mathcal{G}_1, \mathcal{G}_2) &= m_{\mathcal{G}_1 \cup \mathcal{G}_2} H(\mathcal{G}_1 \cup \mathcal{G}_2) - m_{\mathcal{G}_1} H(\mathcal{G}_1) - m_{\mathcal{G}_2} H(\mathcal{G}_2) \\ &= \text{Imp}_{\text{KL}}(\mathcal{G}_1 \cup \mathcal{G}_2) - \text{Imp}_{\text{KL}}(\mathcal{G}_1) - \text{Imp}_{\text{KL}}(\mathcal{G}_2) \end{aligned}$$

Increase of impurity by the merge / Decrease of impurity by partition



Regression Tree for Histograms of Yellowfin tuna fork length

Regression tree for histograms of tuna body length

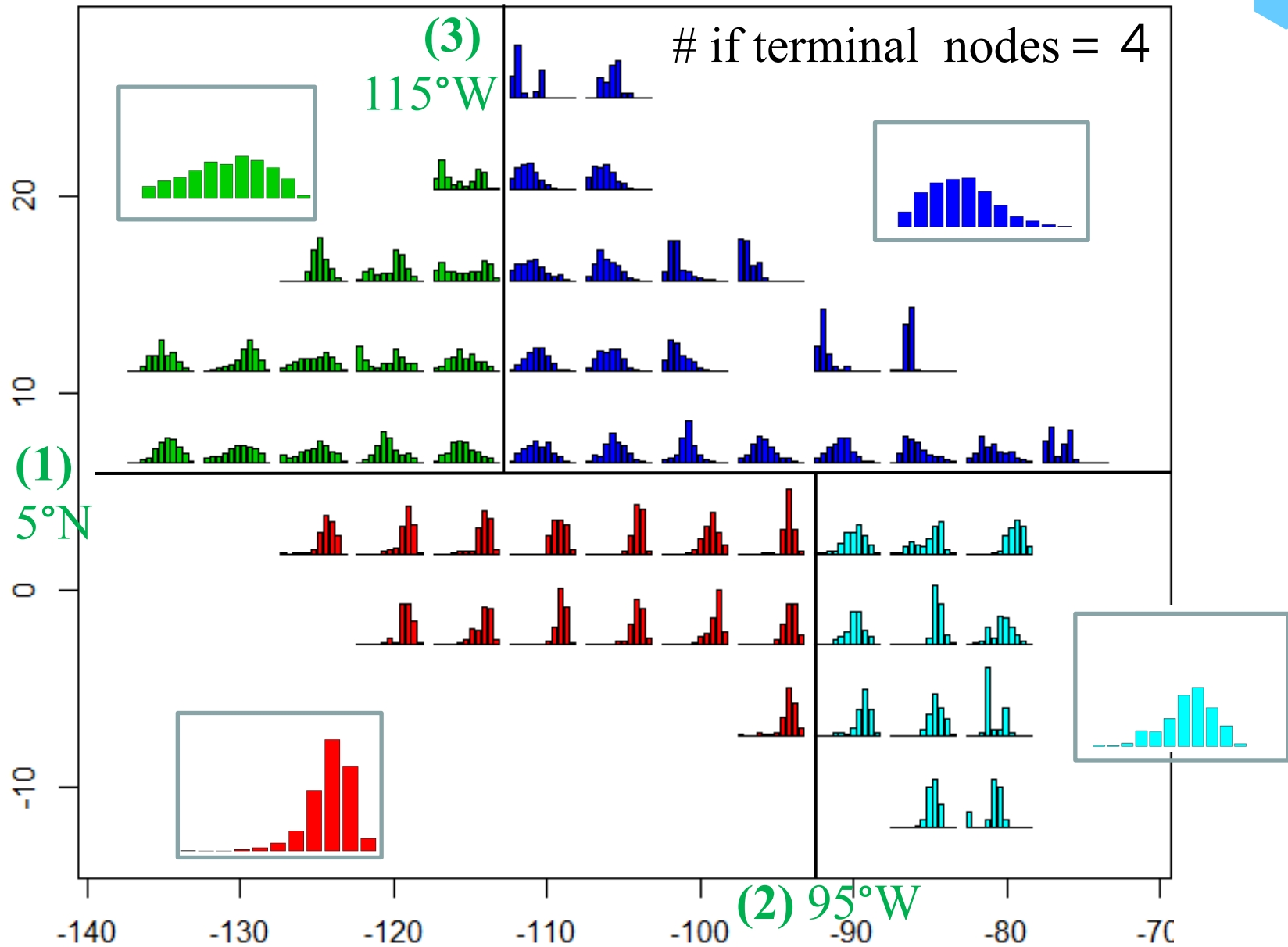
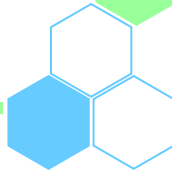


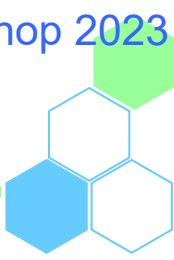
Responses :
distributions of body length

Explanatory variables :
season, latitude, longitude

Cleridy E. Lennert-Cody, Mihoko Minami, Patrick K. Tomlinson, Mark N. Maunder, Fisheries Research (2010)

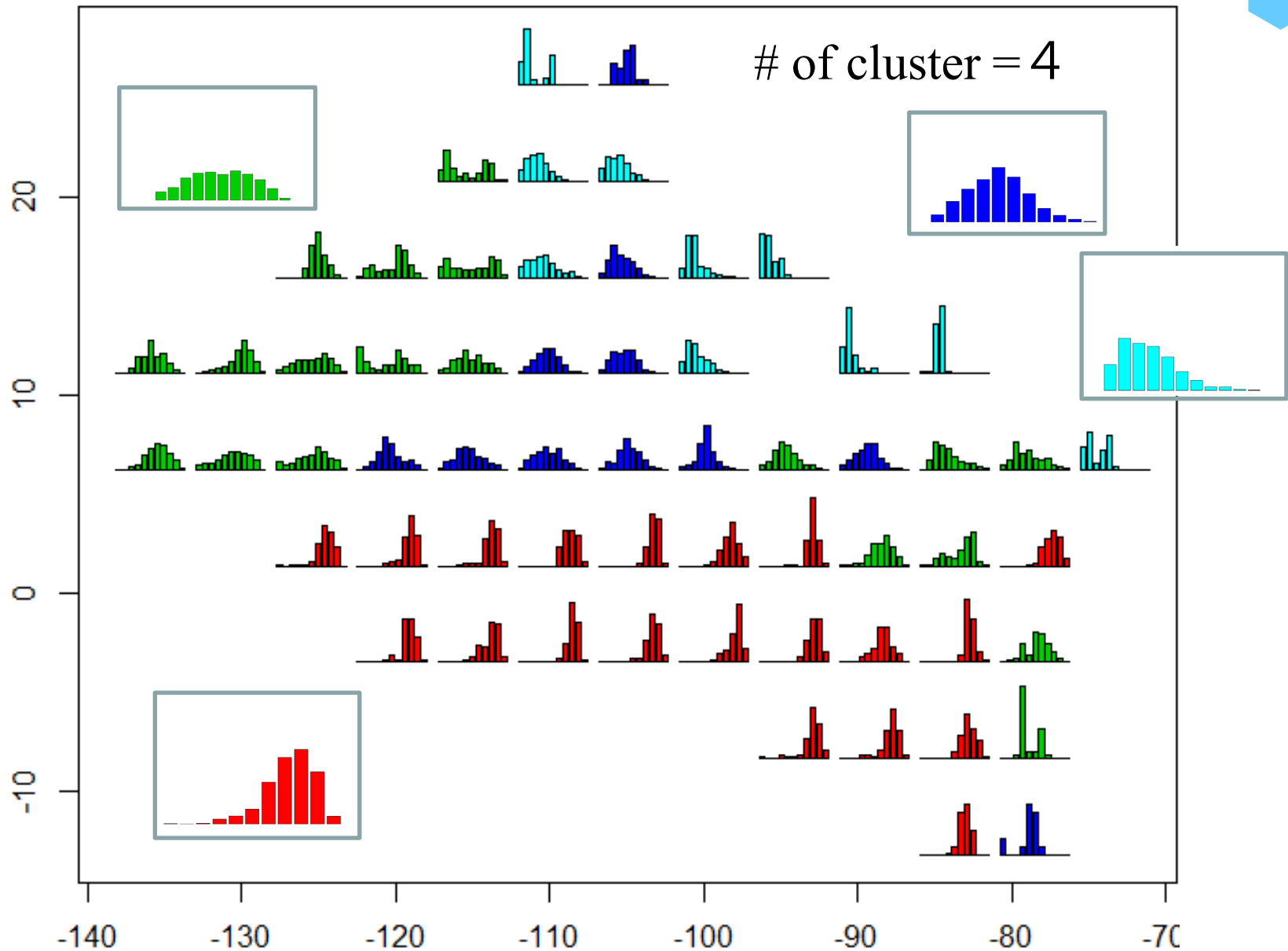
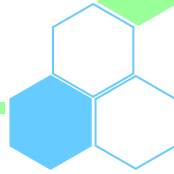
Partition of the Eastern Pacific Ocean by Regression Tree





Hierarchical Clustering for Histograms of Yellowfin tuna fork length

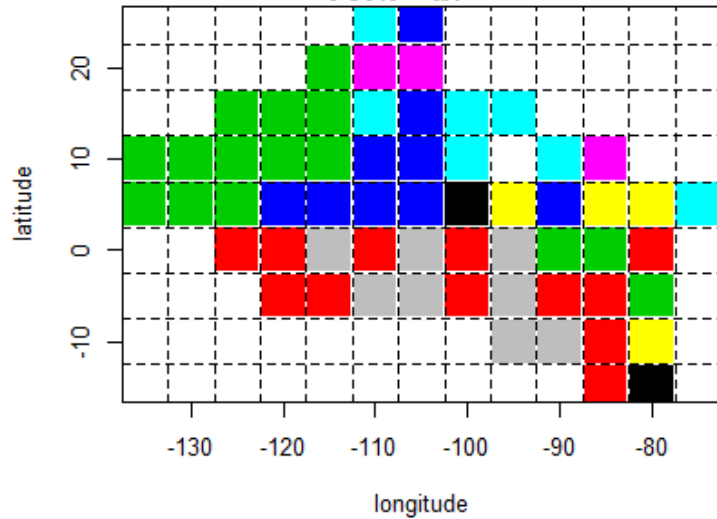
Partition of the Eastern Pacific Ocean by Regression Tree



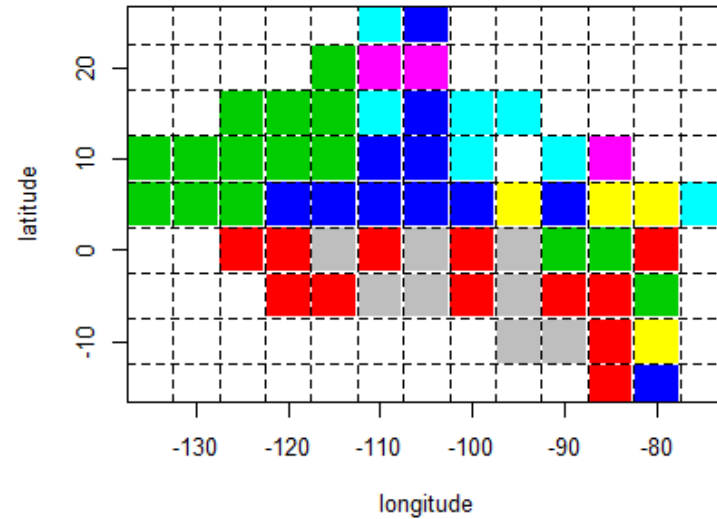
Partitions with different numbers of clusters



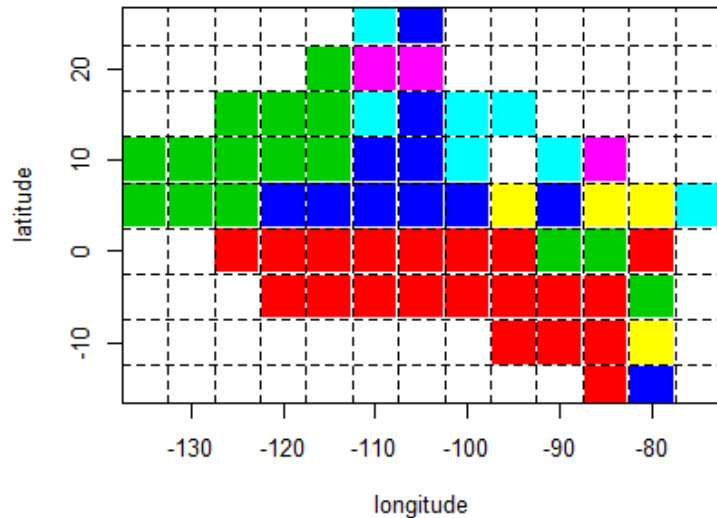
of clusters = 8



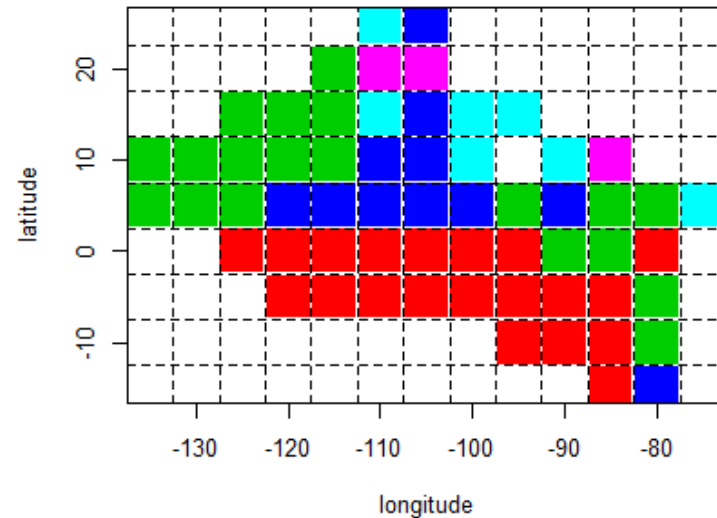
of clusters = 7

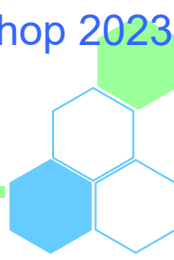


of clusters = 6



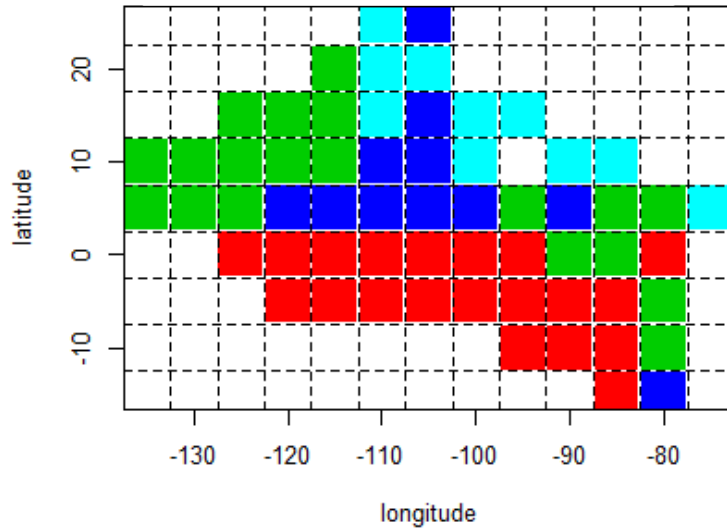
of clusters = 5



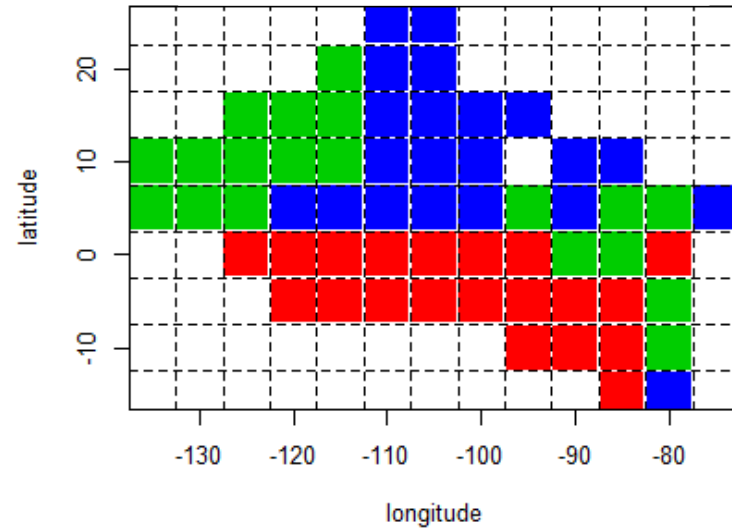


Partitions with different numbers of clusters

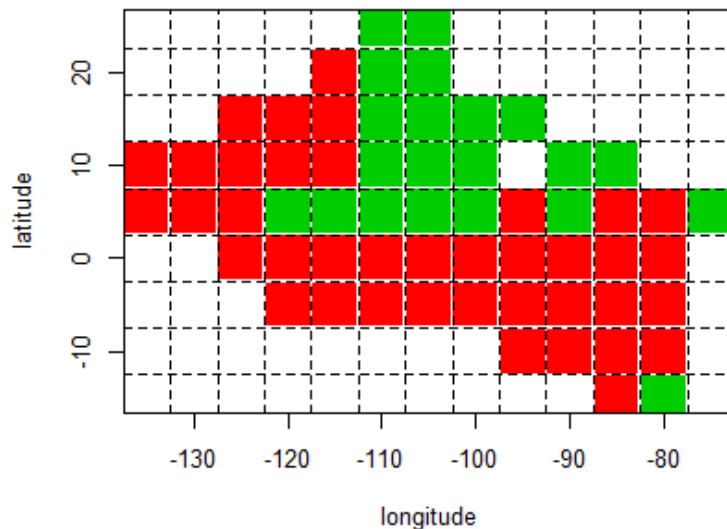
of clusters = 4



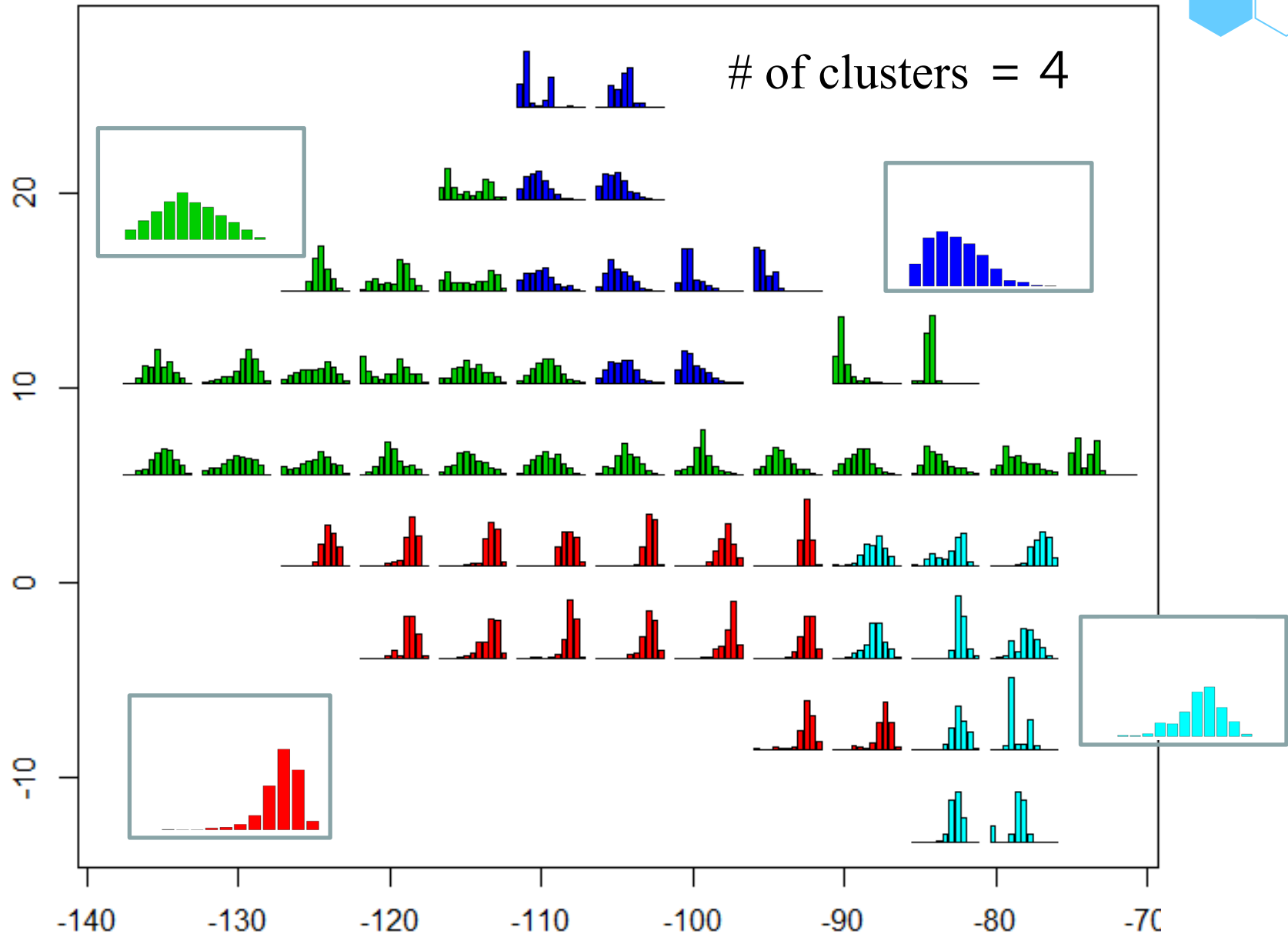
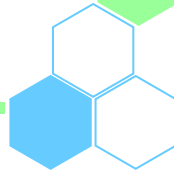
of clusters = 3

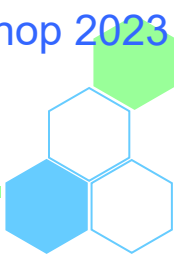


of clusters = 2



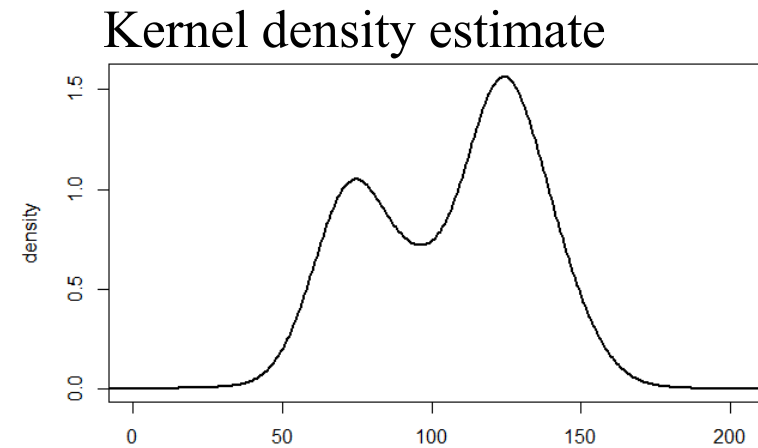
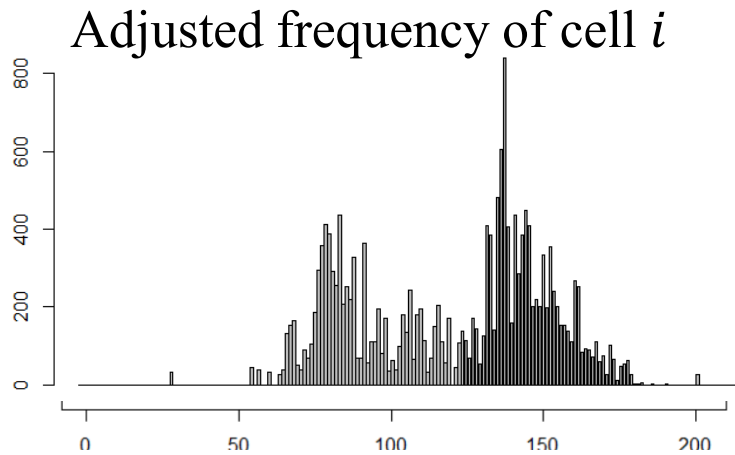
Clustering under connectivity restriction





Regression Tree for Density Estimates of Yellowfin tuna fork length

Kernel density estimation for tuna fork length



Kernel density estimation of tuna body length for cell i

$$\hat{f}_i(x) = \frac{1}{h m_i} \sum_{j=1}^M b_{ij} K\left(\frac{x - x_j}{h}\right)$$

where

$$K(x) = \frac{1}{\sqrt{2\pi}} e^{-x^2/2} \quad \text{Gauss kernel}$$

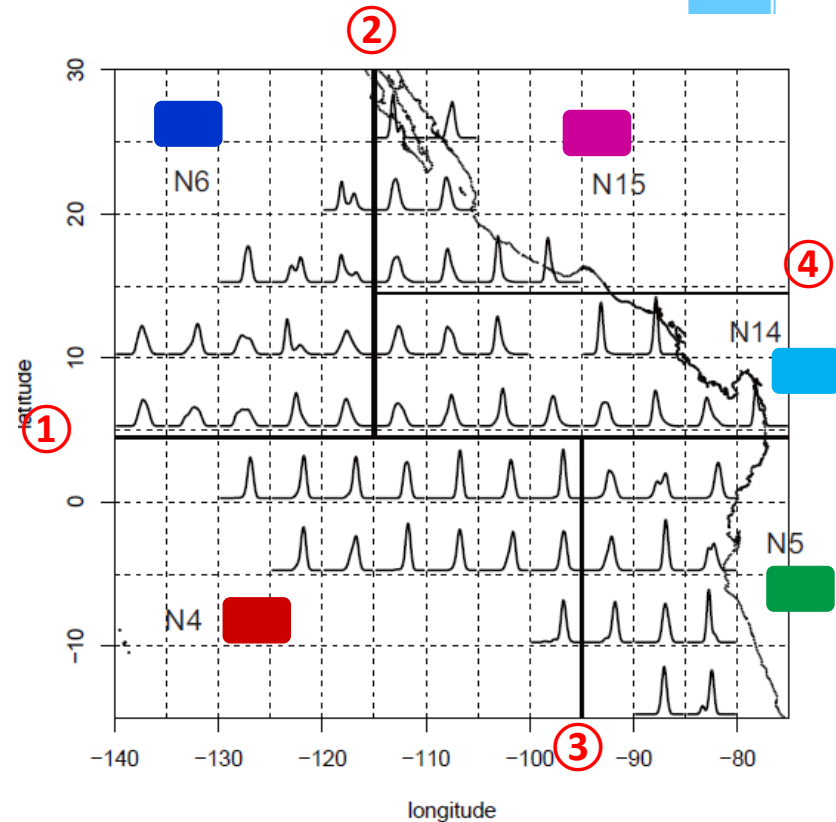
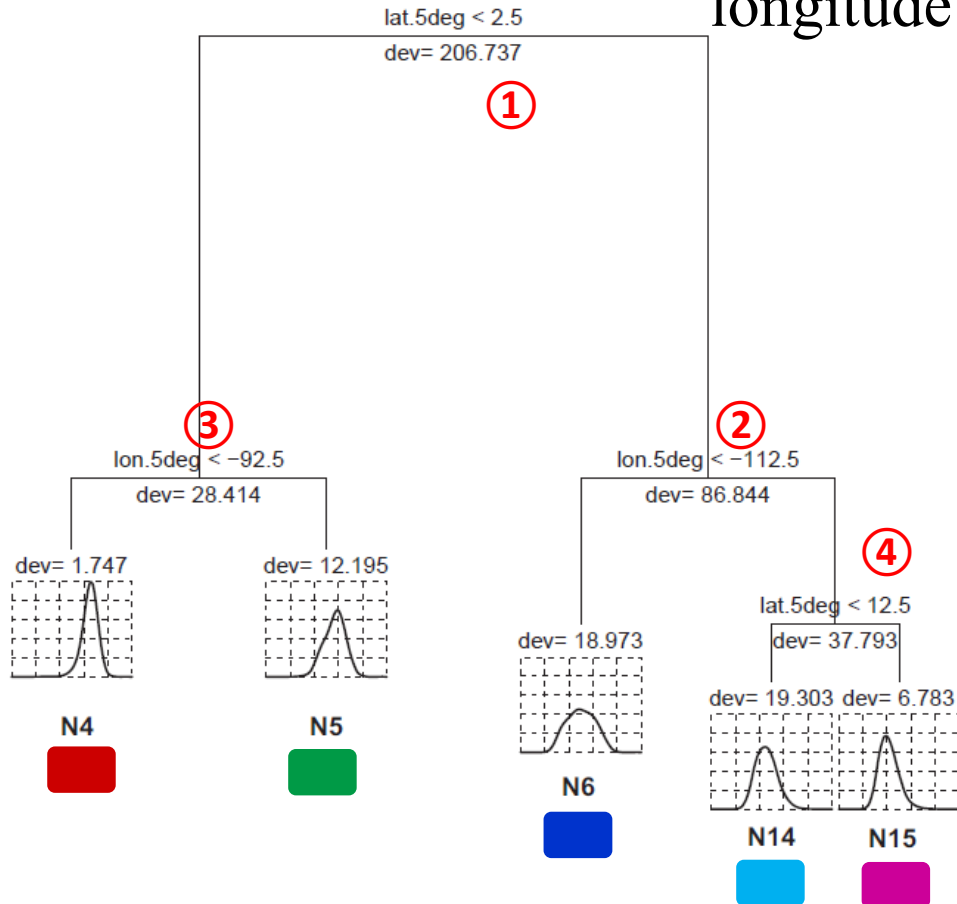
x_j body length of size j ($j = 1, 2, \dots, M$)

b_{ij} adjusted frequency of size j in cell i

$$m_i = \sum_{j=1}^M b_{ij} \quad \text{Total frequency of cell } i \quad \text{sample size} = \text{“confidence”}$$

Regression tree and partitions of the Ocean

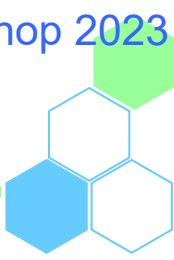
Response variable: density estimates
 Explanatory variables : latitude,
 longitude



(c.f. Lennert-Cody et. al, 2010, 2013)

Stopping rule:

- (1) leaf impurity < $0.1 \times$ root impurity
- (2) confidence < 20



Clustering for Density Estimates of Yellowfin tuna fork length

Distance measures between clusters of distributions



For comparison, we performed clustering with three distances:

◆ **Modified Jensen–Shannon divergence (MJS)**

$$D_{\text{MJS}}(\mathcal{G}_1, \mathcal{G}_2) = m_1 \text{KL}(\bar{f}_1 | \bar{f}_{\{1,2\}}) + m_2 \text{KL}(\bar{f}_2 | \bar{f}_{\{1,2\}})$$

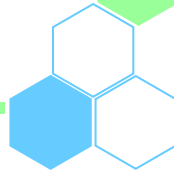
◆ **Earth Mover ‘s distance (EMD)** (Henderson et al. 2015)

$$D_{\text{EMD}}(\mathcal{G}_1, \mathcal{G}_2) \equiv \int_0^1 |\bar{F}_1^{-1}(y) - \bar{F}_2^{-1}(y)| dy = \int_{-\infty}^{\infty} |\bar{F}_1(x) - \bar{F}_2(x)| dx$$

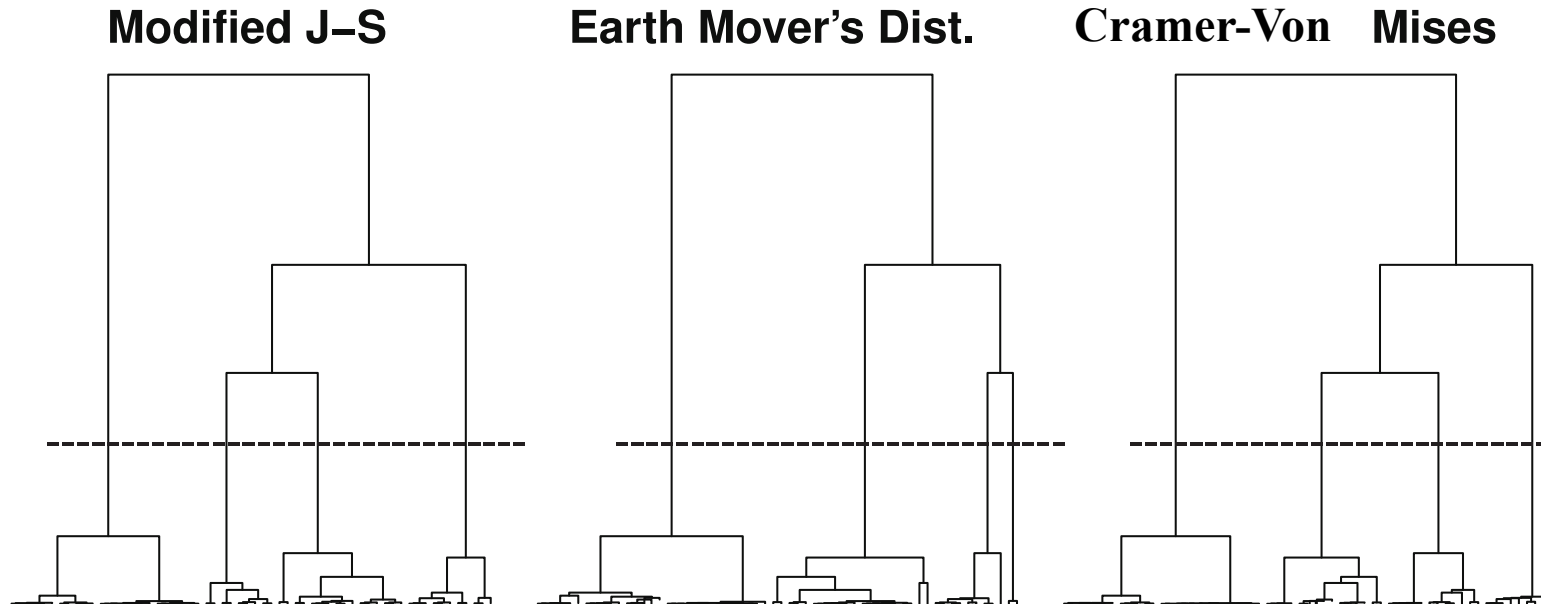
◆ **Cramér–Von Mises type distance**

$$D_{\text{CVM}}(\mathcal{G}_1, \mathcal{G}_2) \equiv \frac{m_1 \cdot m_2}{(m_1 + m_2)} \int_{-\infty}^{\infty} (\bar{F}_1(x) - \bar{F}_2(x))^2 dF^b(x)$$

where F^b is the overall average distribution function as F^b

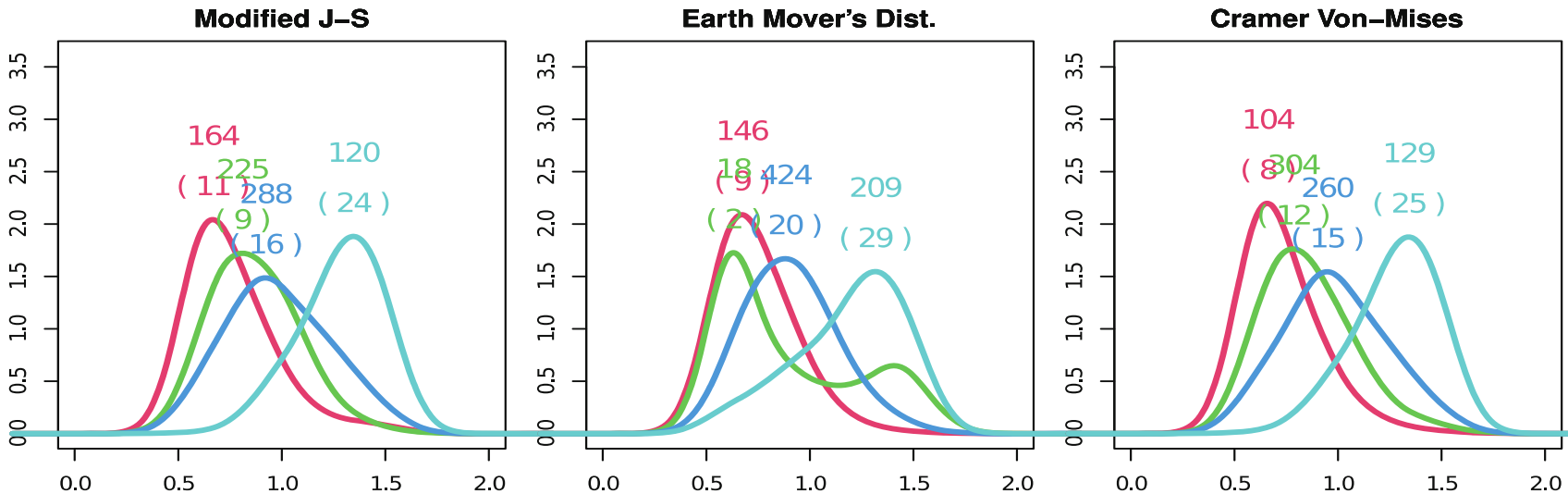
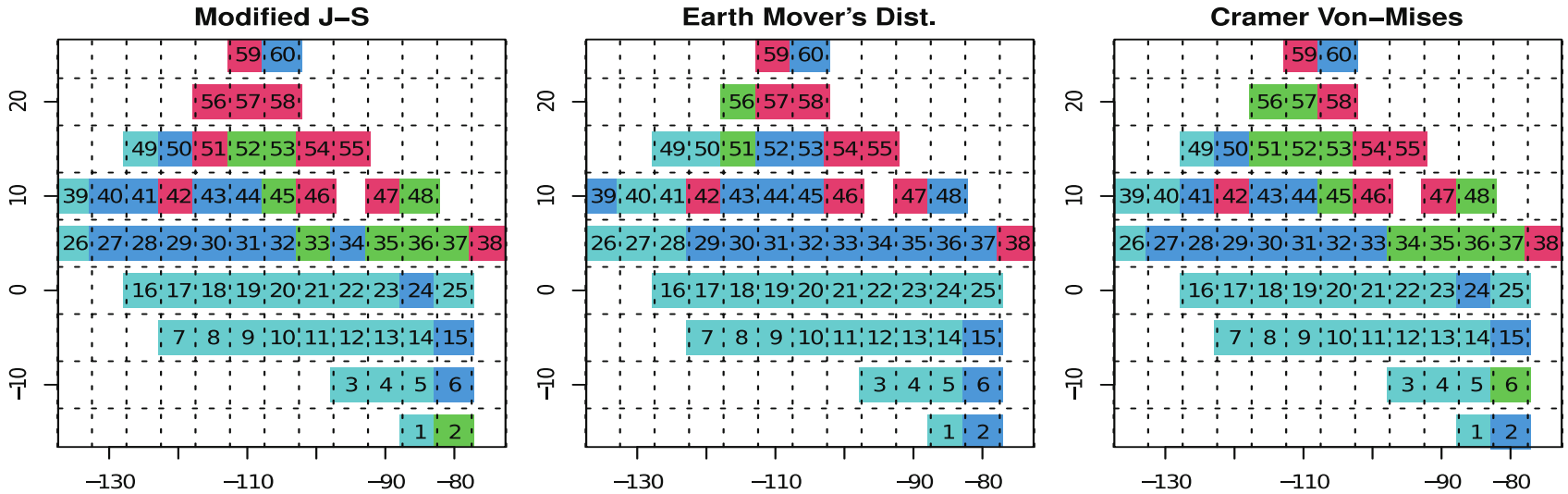
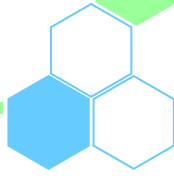


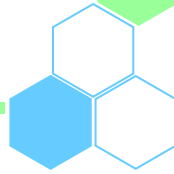
Dendrograms of $5^\circ \times 5^\circ$ cells



- ◆ clusters by modified J-S divergence and Cramer Von-Mises distance are somehow similar compared to the clusters by Earth Mover's distance
- ◆ Earth Mover's distance produced a cluster with a small confidence (sample size)

Clusters and distributions of $5^\circ \times 5^\circ$ cells





Comparison of results by three distances

◆ Modified Jensen–Shannon divergence (MJS)

$$D_{\text{MJS}}(\mathcal{G}_1, \mathcal{G}_2) = m_1 \text{KL}(\bar{f}_1 | \bar{f}_{\{1,2\}}) + m_2 \text{KL}(\bar{f}_2 | \bar{f}_{\{1,2\}})$$

◆ Earth Mover ‘s distance (EMD) (Henderson et al. 2015)

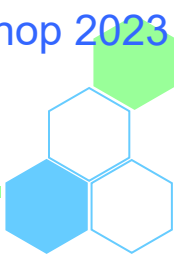
$$D_{\text{EMD}}(\mathcal{G}_1, \mathcal{G}_2) \equiv \int_0^1 |\bar{F}_1^{-1}(y) - \bar{F}_2^{-1}(y)| dy = \int_{-\infty}^{\infty} |\bar{F}_1(x) - \bar{F}_2(x)| dx$$

c.f. $L_1(\mathcal{G}_1, \mathcal{G}_2) \equiv \int_{-\infty}^{\infty} |\bar{f}_1(x) - \bar{f}_2(x)| dx$ (T. Nguyen-Trang et al. 2022)

◆ Cramér–Von Mises type distance

$$D_{\text{CVM}}(\mathcal{G}_1, \mathcal{G}_2) \equiv \frac{m_1 \cdot m_2}{(m_1 + m_2)} \int_{-\infty}^{\infty} (\bar{F}_1(x) - \bar{F}_2(x))^2 dF^b(x)$$

Measure	Modified J-S			Earth Mover’s Dist.			Cramér Von-Mises		
	# of cells	# of sets	mean length(m)	# of cells	# of sets	mean body length(m)	# of cells	# of sets	mean body length(m)
cluster 1	11	164	0.771	9	146	0.754	8	104	0.741
cluster 2	9	225	0.881	2	18	0.913	12	304	0.866
cluster 3	16	288	0.994	20	424	0.919	15	260	0.999
cluster 4	24	120	1.294	29	209	1.197	25	129	1.289



Testing procedures for homogeneity and the minimal homogeneous tree structure

Testing homogeneity of kernel density estimates

Cao and Keilegom (2006) considered the problem to test

$$H_0: F_1 = F_2$$

using the kernel density estimates \hat{f}_1 and \hat{f}_2 obtained from two independent random samples,

$$X_1, \dots, X_n \sim^{i.i.d.} F_1, \quad Y_1, \dots, Y_m \sim^{i.i.d.} F_2.$$

They defined a test statistics, showed its asymptotic distribution and proposed a testing procedure for $H_0: F_1 = F_2$

However, their testing procedure cannot be applied to test the homogeneity of a cluster because member distributions or samples are results of previous merges and are not independent.

Randomization sample for testing homogeneity



Suppose $\mathcal{G} = \{(\hat{f}_i, m_i), i \in G\}$, $G = G_L \cup G_R$ and $D_{\text{MJS}}(\mathcal{G}_L, \mathcal{G}_R) = d$.

H_0 : \mathcal{G} is homogeneous (i.e., all samples used to estimate \hat{f}_i are from the same distribution)

◆ **Generation of randomization sample $\mathbf{t} = (t_1, t_2, \dots, t_K)$**

For $k = 1, 2, \dots, K$, repeat the following steps

1. Generate a sample \mathbf{x}_i^k of size m_i from \bar{f}_G for all $i \in G$
2. Compute density estimate \hat{g}_i^k with \mathbf{x}_i^k for all $i \in G$
3. Perform clustering with $\mathcal{G}^k = \{(\hat{g}_i^k, m_i), i \in G\}$.
4. Let $t_k = D_{\text{MJS}}(\mathcal{G}_L^k, \mathcal{G}_R^k)$ where \mathcal{G}_L^k and \mathcal{G}_R^k are two clusters combined at the last step, that is, $G = G_L^k \cup G_R^k$.

Testing procedure for homogeneity of a cluster

- ◆ Generate a randomization sample \mathbf{t} of small size K_1 (e.g. 100)
- ◆ If $\#\{t_k | t_k > d, k = 1, 2, \dots, K_1\} \geq N_{d1}$ (e.g. 10)
then, **[p value for d]** is greater than $\frac{N_{d1}}{K_1}$ (e.g. 0.1)
- ◆ If $\#\{t_k | t_k > d, k = 1, 2, \dots, K_1\} < N_{d1}$ (e.g. 10)
then, compute the Chebyshev's upper bound U of $P(X \geq d)$
under H_0 computed with sample mean and variance of sample \mathbf{t}
 - ◆ If the upper bound $U \leq \epsilon$ (e.g. 0.001)
then, **[p value for d]** is less than U ($U = \frac{svar(\mathbf{t})}{(d - \bar{t})^2}$)
 - ◆ If the upper bound $U > \epsilon$ (e.g. 0.001)
then, generate a randomization sample of large size
 K_2 (e.g. 1000) and let $N_{d2} = \#\{t_k | t_k > d, k = 1, 2, \dots, K_2\}$
[p value for d] = $\inf\{p | P(W \leq N_{d2}) \leq 0.05, W \sim \text{Bin}(K_2, p)\}$

Hierarchical testing procedure for homogeneity



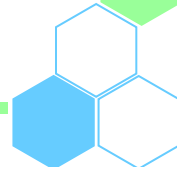
Let $H_0^{(j)}$ be “Cluster \mathcal{G}_j is homogeneous”.

Because the null hypotheses have the Hierarchical structure,

$$\mathcal{G}_i \supset \mathcal{G}_j \text{ implies } H_0^{(i)} \Rightarrow H_0^{(j)},$$

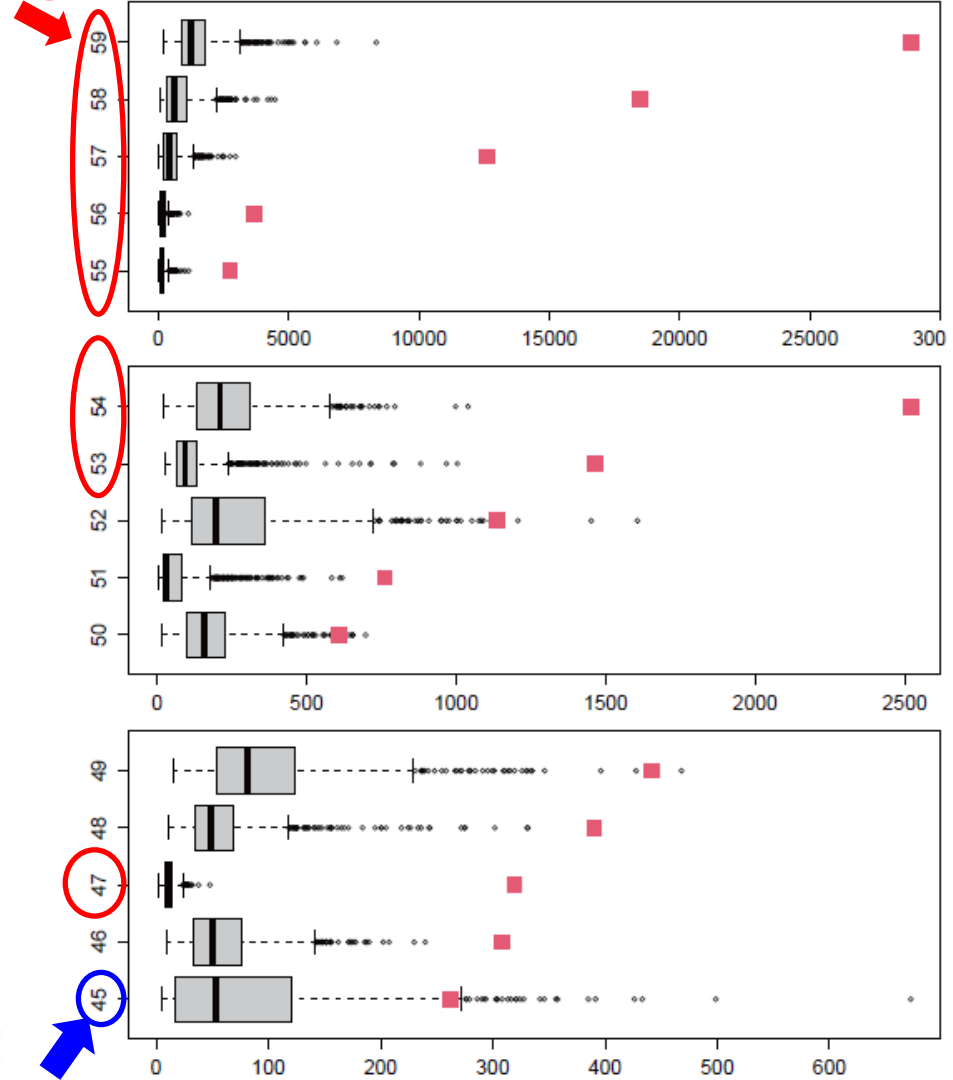
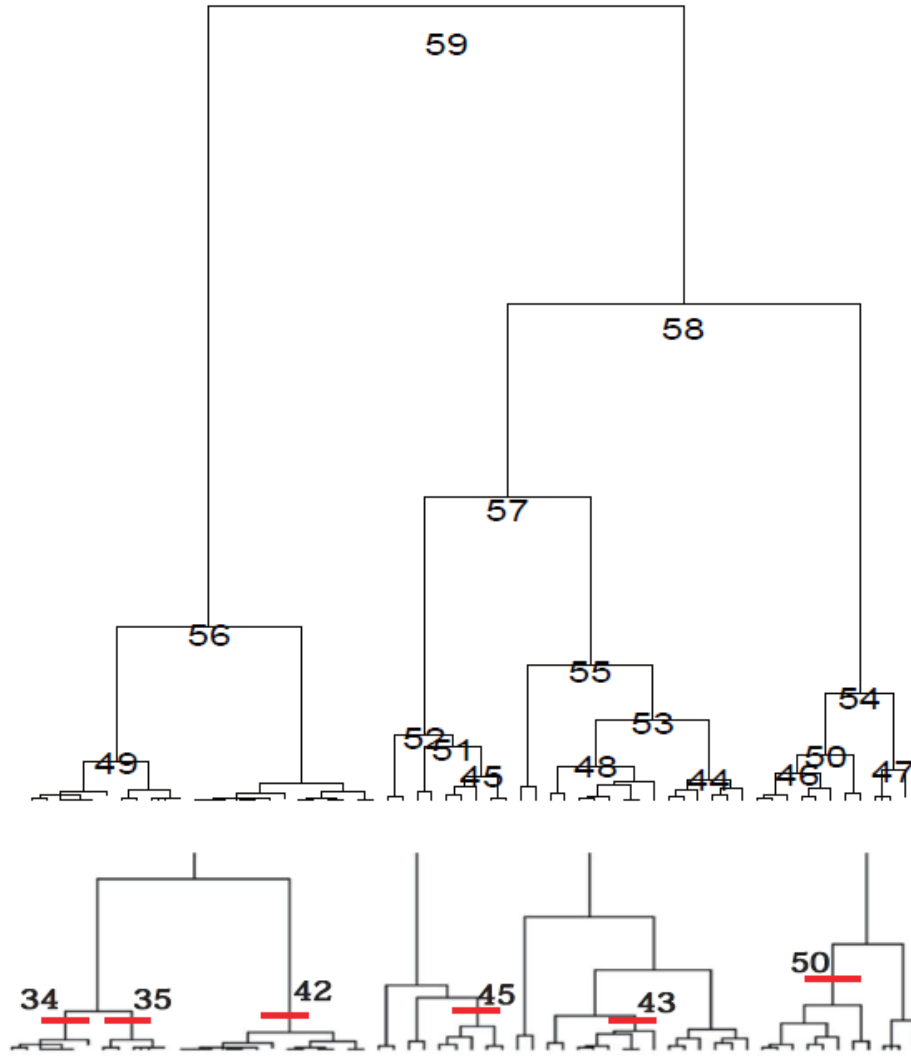
the family-wise error rate is controlled at the significance level α with the following hierarchical testing procedure:

- ◆ Test starts from the top, the cluster of all distributions
- ◆ If the hypothesis is significant at α ,
then, the hypothesis “the cluster is homogeneous” is rejected
The test proceeds to child clusters.
- ◆ If the hypothesis is not significant at α ,
then, the hypothesis “the cluster is homogeneous” is accepted
The child clusters are not tested.



Boxplot of randomization samples with size K_1

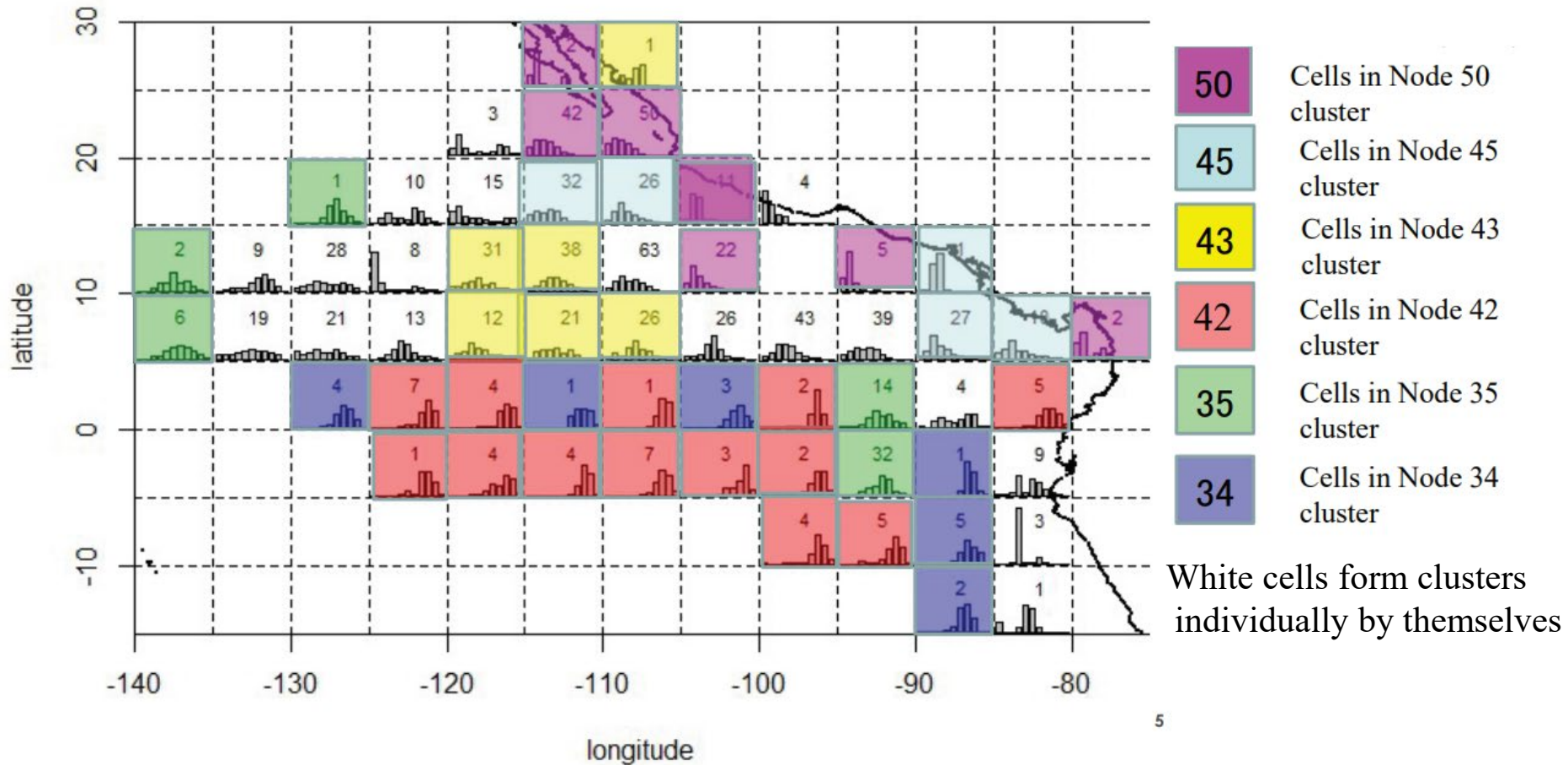
Not homogeneous at $\alpha = 0.01$

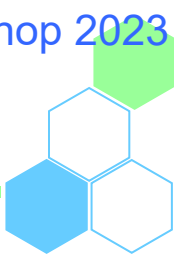


Cannot be said “Not homogeneous”

The minimal Homogeneous tree structure

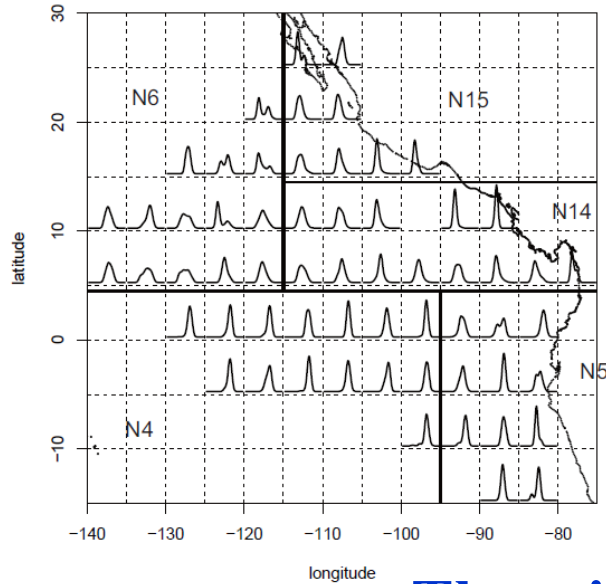
with $\alpha = 0.01$



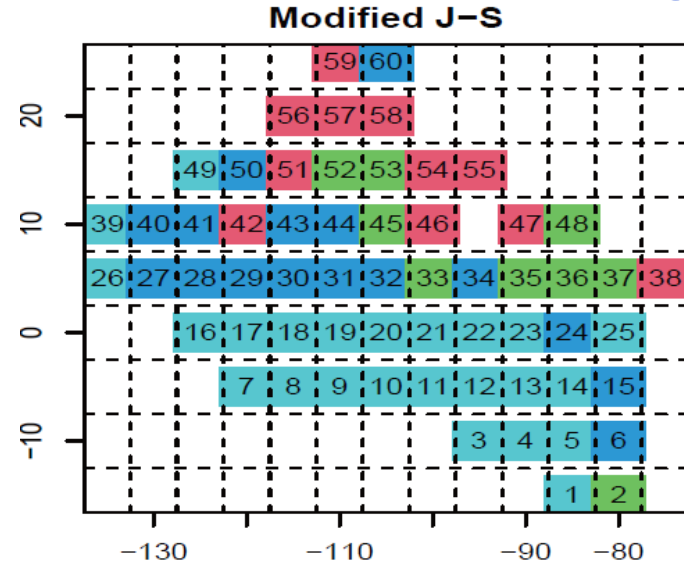


Comparison of the results

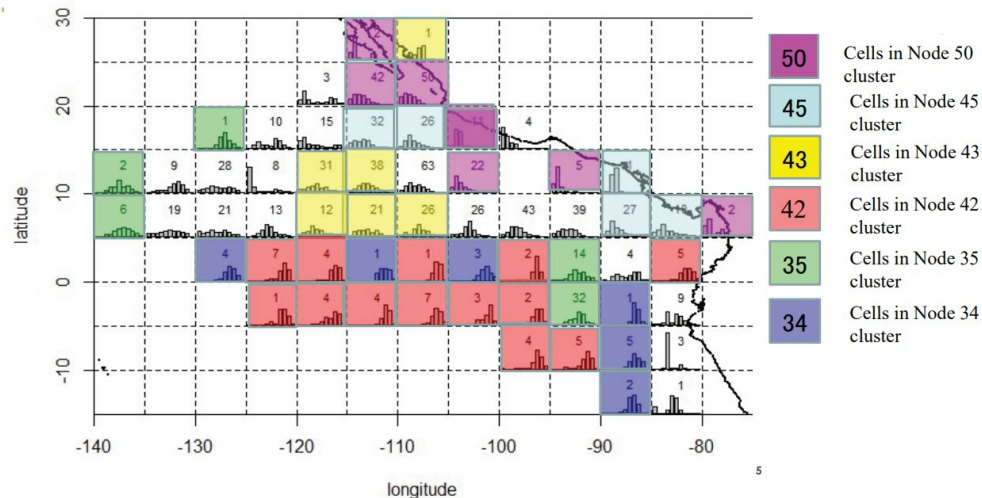
Regression tree

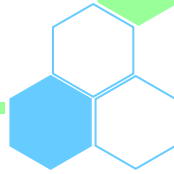


Hierarchical clustering



The minimal homogeneous structure





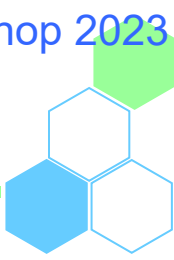
Near-homogeneous tree structure

Homogeneity might be too demanding for defining a cluster.

As a relaxed concept, we regard a cluster whose Chebyshev's upper bound U obtained with random sample is less than θ , is near-homogeneous with cut-off value θ .

Cut-off value interval	# of nodes	Terminal nodes in the minimal near-homogeneous tree
[31.32, ∞]	1	59
[26.81, 31.32)	3	58, 49, 42
[15.08, 26.81)	5	55, 54, 52, 49, 42
[14.20, 15.08)	7	55, 52, 50, 49, 42, (42), 21
[9.94, 14.20)	10	52, 50, 49, 48, 44, 42, (42), (6), (34), 21
[9.41, 9.94)	11	52, 50, 49, 44, 43, 42, (42), (6), (34), 21, 30
[6.50, 9.41)	12	52, 50, 49, 44, 43, 42, (42), (6), (34), 21, (15), (29)
[9.41, 9.94)	16	52, 50, 49, 43, 42, (42), (6), (34), 21, (15), (29), 23, (40), (28), (27), (41)

Numbers without parenthesis are node(merge) numbers and numbers in parenthesis are original cell numbers

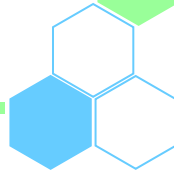


Summary and future work

- We consider regression tree and hierarchical clustering methods for distributions based on the modified Jensen–Shannon divergence.
- We presented a testing procedure for homogeneity of a cluster and a hierarchical testing procedure to find the minimal homogeneous/near-homogeneous tree structure of distributions.
- These methods and procedures are applied to the yellowfin tuna fork length data

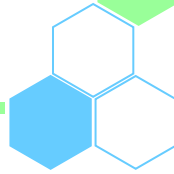
Future work

- ◆ We would like to extend the method to Bayesian clustering for distributions with prior for partitions.



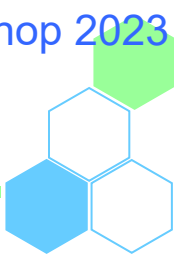
Reference

- Breiman, L. and Friedman, J. H. and Olshen, R. A. and Stone, C. J. (1984). *Classification and Regression Trees*, Wadsworth and Brooks, Monterey, CA.
- Cao, R. and van Keilegom I. (2006) Empirical Likelihood Tests for Two-Sample Problems via Nonparametric Density Estimation, *Can J Stat* 34, 61–77
- Cha, S.-H. (2007). Comprehensive survey on distance/similarity measures between probability density functions. *Int. J. Math. Model. Meth. Appl. Sci.*(1), 300–307
- Dhillon, I. S., S. Mallela, and R. Kumar (2003). A divisive information-theoretic feature clustering algorithm for text classification. *Journal of Machine Learning Research*(3), 1265–1287.
- Henderson, K., Gallagher, B., Eliassi-Rad, T. (2015). EP-MEANS: An Efficient Nonparametric Clustering of Empirical Probability Distributions. *SAC '15* Proceedings of the 30th Annual ACM Symposium on Applied Computing, 893–900.
- Gordon, A.D. (1999). *Classification*, 2nd Edition, CRC Press,
- Irpino, A. and Lechevallier, Y. 2006. Dynamic clustering of histograms using Wasserstein metric. In *COMPSTAT 2006 – Advances in Computational Statistics* pp 869-876.
- Irpino, A., Verde, R., de Carvalho, F. de A.T. 2014. Dynamic clustering of histogram data base don adaptive squared Wasserstein distances. *Expert Systems with Applications* 41: 3351-3366. DOI | 10.1016/j.eswa.2013.12.001

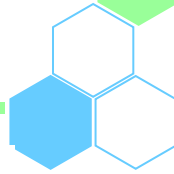


Reference

- Jiang, B., J. Pei, Y. Tao and X. Lin(2013) Clustering uncertain data based on probability distribution similarity. *IEEE Transactions on Knowledge and Data Engineering*, 25(4), 751--763.
- Lennert-Cody, C. E., Minami, M. Patrick K. Tomlinson, Mark N. Maunder (2010), Exploratory analysis of spatial-temporal patterns in length-frequency data: An example of distributional regression trees, *Fisheries Research*, 102, 323-326
- Lennert-Cody, C. E., Maunder, M. N., Aires-da-Silva, A., & Minami, M. (2013). Defining population spatial units: Simultaneous analysis of frequency distributions and time series. *Fisheries Research*, 139, 85-92.
- Minami, M. and Lennert-Cody, C.E. (2023). Regression Tree and Clustering for Distributions, and Homogeneous Structure of Population Characteristics, submitted
- Nielsen, F., Nock, R. and Amari, S-I. 2014. On clustering histograms with k-Means by using Mixed alpha-Divergences. *Entropy* 16 (6): 3273-3301.
- Thao Nguyen-Trang, Trung Nguyen-Thoi1, Tai Vo-Van (2023) Globally automatic fuzzy clustering for probability density functions and its application for image data, *Applied Intelligence*,
- Dinh Phamtoan, Tai Vovan (2022), Automatic fuzzy clustering for probability density functions using the genetic algorithm, *Neural Computing and Applications* (2022) 34:14609-14625.



Thank you for your attention



Standardized mJS of Nodes

Node ID	Std. mJS	Upper bound	out of 100	out of 1000	p-value	parent node	child nodes	Significance
59	31.32	0.0010	0		< 0.0010	—	56, 58	**
58	26.81	0.0014	0		< 0.0014	59	57, 54	**
57	28.99	0.0012	0		< 0.0007	58	52, 55	**
56	36.61	0.0007	0		< 0.0007	59	49, 42	**
55	14.20	0.0050	0		< 0.0050	57	40, 53	**
54	15.08	0.0044	0		< 0.0044	58	50, 47	**
53	18.14	0.0030	0		< 0.0030	55	48, 44	**
52	4.28	0.0547	1	3	0.0077	57	20, 51	**
51	9.35	0.0114	0	0	0.0030	52	32, 45	**
50 [†]	4.47	0.0502	1	6	0.0118	54	46, 31	*
49	6.03	0.0275	0	1	0.0047	56	34, 35	**
48	9.94	0.0101	0	0	0.0030	53	30, 43	**
47	59.27	0.0003	0		< 0.0003	54	-42 [†] , 21	**
46	9.01	0.0123	0	0	0.0030	50	25, 38	*
45 [†]	2.38	0.1759	5	31	0.0416	51	39, 19	*
44	6.50	0.0237	0	0	0.0030	53	36, 37	**
43 [†]	2.45	0.1662	6	37	0.0484	48	41, -43	*
42 [†]	2.59	0.1490	1	19	0.0278	56	29, 33	*
41	9.64	0.0108	0	0	0.0030	43	22, 9	*