



सिद्धिपूर्व प्रवन्धनम्  
भा. प्र. सं. इन्दौर  
IIM INDORE

# Hierarchical Clustering

## An Application on Wholesale Market Data

---

Aarjav Sethi  
Apurv Chaudhary  
Harshvardhan  
Pradeep Charan  
Tamanna Gupta

IPM 2016-21 Batch

# TABLE OF CONTENTS

1. Introduction
2. Problem
3. Hierarchical Clustering
4. Application

# Introduction

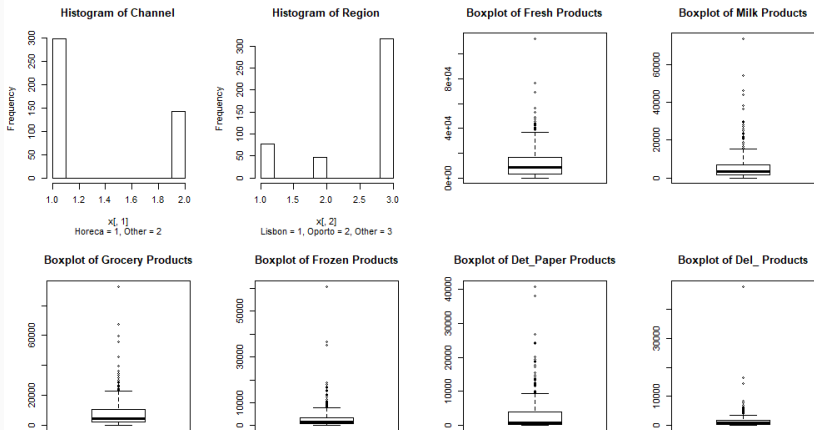
---

## DESCRIPTION OF DATA

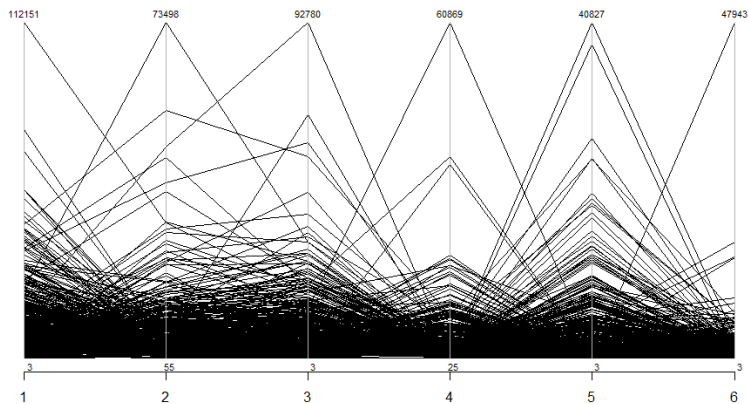
Data Source: University of California, Irvine's Machine Learning Repository

Attribute	Annual Spending/Description	Type
FRESH	Fresh Products	Continuous
MILK	Milk Products	Continuous
GROCERY	Grocery Products	Continuous
FROZEN	Frozen Products	Continuous
DETERGENTS	Detergents and Paper Products	Continuous
DELICATESSEN	Delicatessen Products	Continuous
CHANNEL	Horeca or Retail Channel	Nominal
REGION	Lisbon, Oporto or Other	Nominal

## Histograms and Boxplots for descriptive knowhow of data



# PARALLEL COORDINATE PLOT



# Problem

---

PROBLEM: *Group* wholesale customers on the basis of their buying patterns

METHOD: Use *Hierarchical Clustering* to group based on “similar” buying patterns



Using *Unsupervised Learning Methods* to group data without any previously known attributes of such groups

Assume these values are output of some internal function unknown to us

## K-means VS Hierarchical

- Complexity:  $O(n)$  and  $O(n^2)$
- Consistency: *K-means* renders different results with every run
- No previously known  $k$ !

# Hierarchical Clustering

---

Builds *clusters*, i.e. groups that have maximum similarity with each other in the same cluster and maximum dissimilarity with other clusters

Results in *hierarchy of clusters*

## Agglomerative

“bottom up” approach; each starts in its own cluster and then merged as we move up the hierarchy

## Divisive

“top down” approach; all observations start in one cluster, and splits as one moves down the hierarchy

*How to measure similarity between two clusters?*

- Euclidean Distance

$$d_0^2 = (\mathbf{x}_1 - \mathbf{x}_2)^T (\mathbf{x}_1 - \mathbf{x}_2)$$

- Correlation Distance

$$d_0 = \frac{\text{Cov}(X_1, X_2)}{\sqrt{\text{Var}(X_1) \cdot \text{Var}(X_2)}} \text{ (Pearson's Correlation Coefficient)}$$

## Examples

- **Similarity based on Sizes:** grouping based on buying sizes; grouping which fruits are larger in size and which are smaller
- **Similarity based on Attributes:** grouping based on buying patterns; grouping which fruits are edible and which aren't

*How to measure distance between a cluster and a point of another cluster?*

- **Single** - minimum of all possible distances; chaining effects
- **Complete** - maximum of all possible distances; no chaining; affected by outliers
- **Average** - unweighted average; compromise b/w single and complete
- **Centroid** - average of points in the cluster then distance is calculated

## An Example

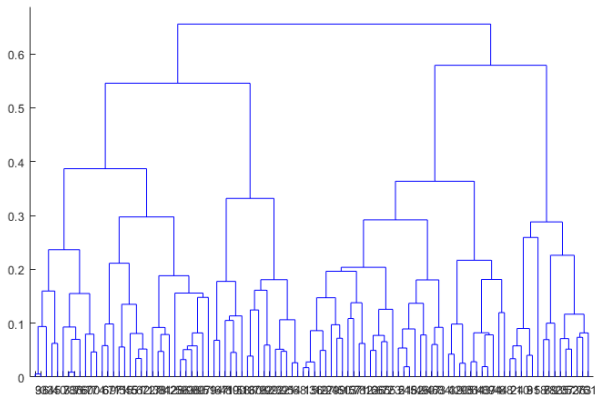
Consider clusters -  $\{1, 4\}$  and  $\{2\}$

Linkage	Calculation	Distance
Single	$\min\{(2 - 1), (4 - 2)\}$	1
Complete	$\max\{(2 - 1), (4 - 2)\}$	2
Average	$\frac{(2-1)+(4-2)}{2}$	1.5
Centroid	$\frac{1+4}{2} - 2$	0.5

# DENDROGRAMS

A *tree based* diagram to represent *hierarchies* of the hierarchical clustering

## Example





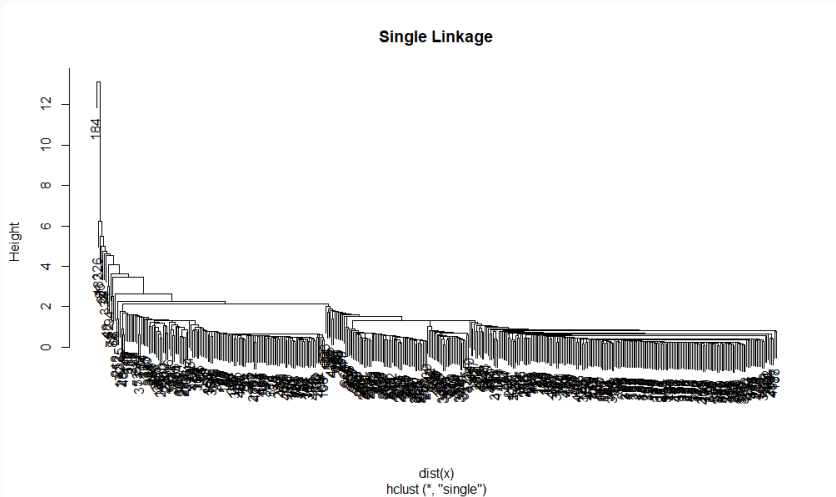
# Application

---

# Dendrograms

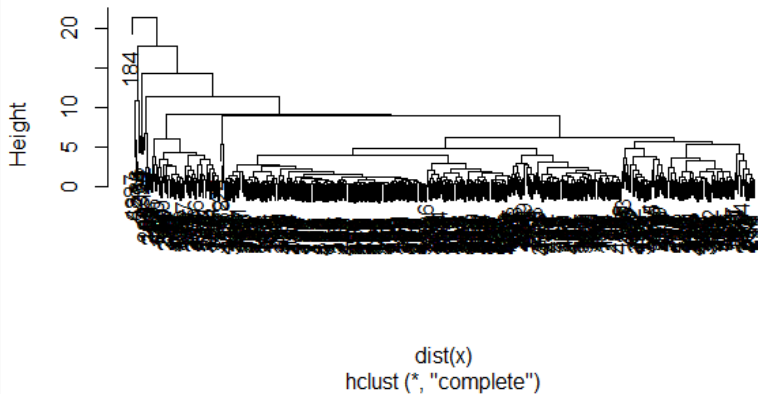
---

# DENDROGRAM FOR SINGLE-EUCLIDEAN

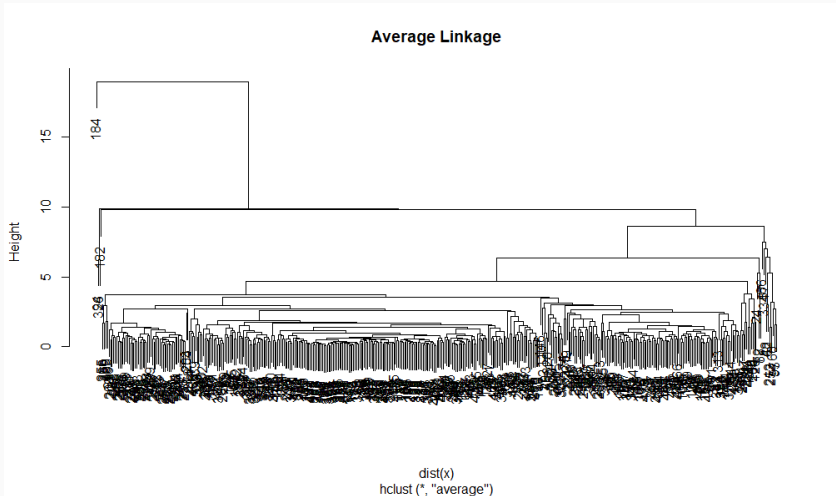


# DENDROGRAM FOR COMPLETE-EUCLIDEAN

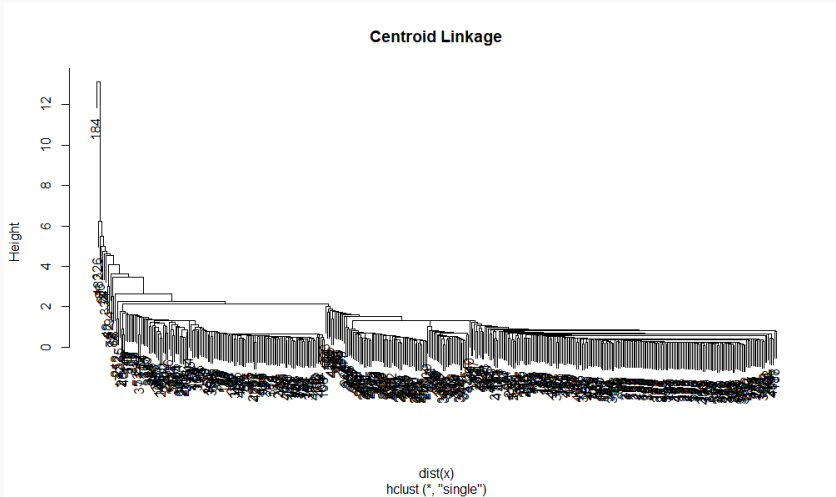
## Complete Linkage



# DENDROGRAM FOR AVERAGE-EUCLIDEAN

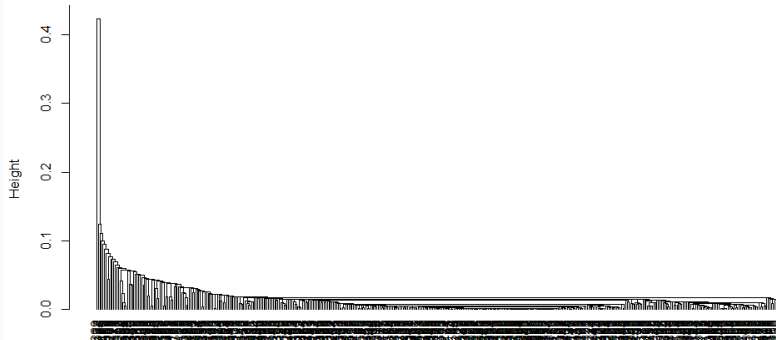


# DENDROGRAM FOR CENTROID-EUCLIDEAN

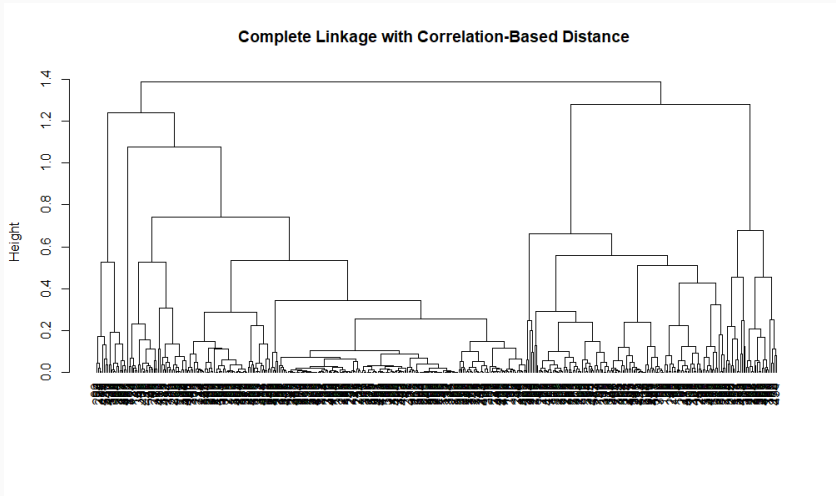


# DENDROGRAM FOR SINGLE-CORRELATION

Single Linkage with Correlation-Based Distance

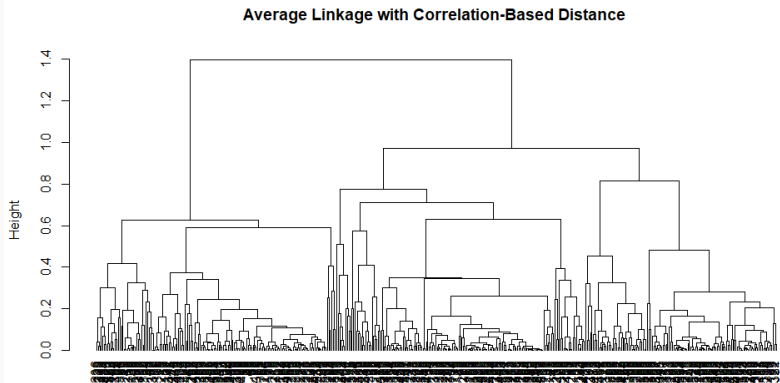


# DENDROGRAM FOR COMPLETE-CORRELATION

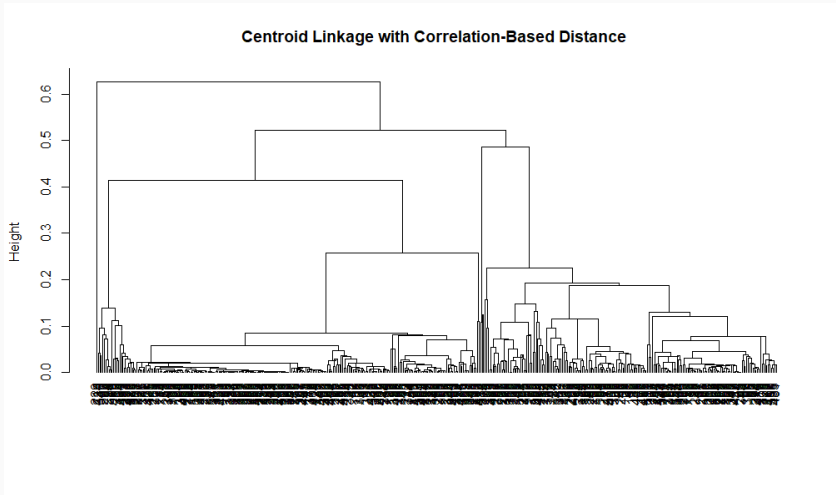




## DENDROGRAM FOR AVERAGE-CORRELATION



# DENDROGRAM FOR CENTROID-CORRELATION



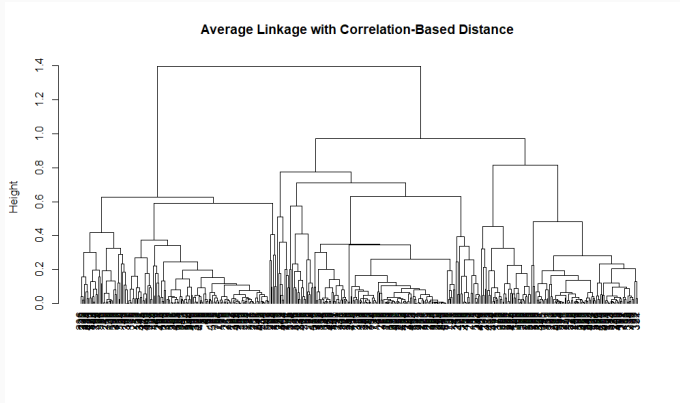
# CHOOSING DISTANCE AND SIMILARITY MEASURE FOR POST ANALYSIS

**Distance Measure** *average* over single, complete and centroid

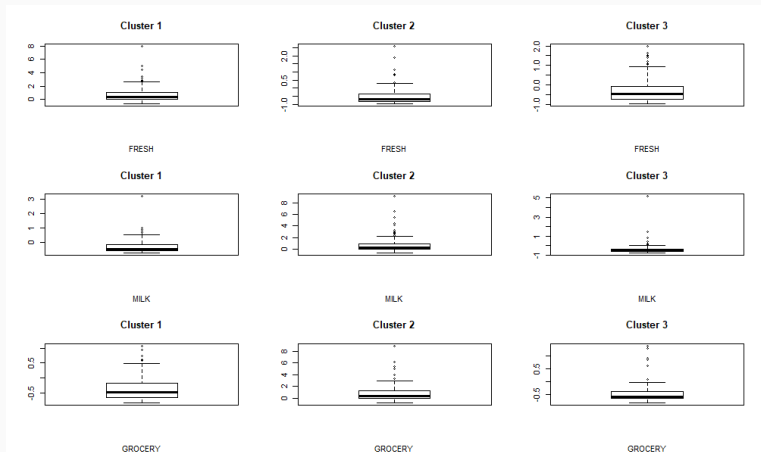
**Similarity Measure** *correlation* as we wanted to find customers with similar buying pattern rather than size of customers

# NUMBER OF CLUSTERS

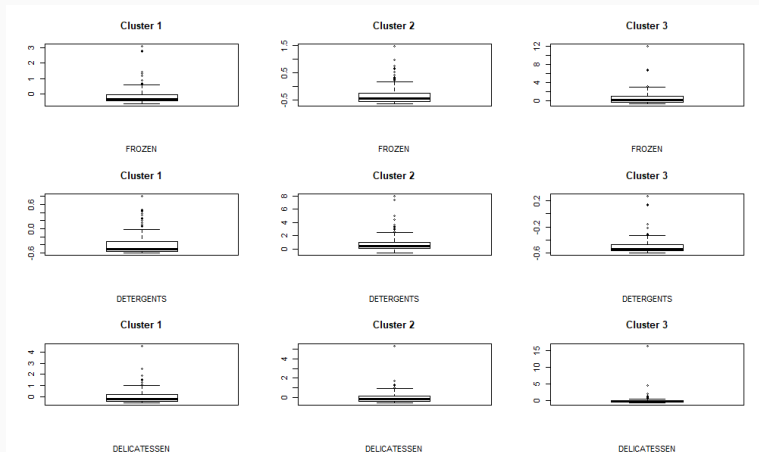
Choose *three* clusters for further analysis.



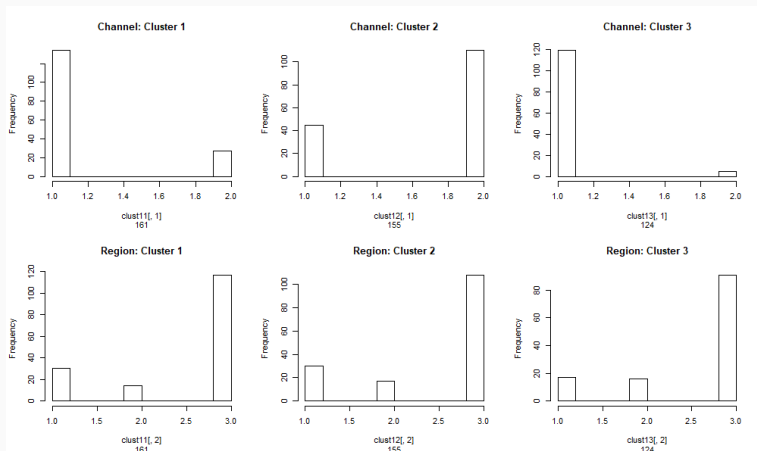
# INTERPRETATING CLUSTERS...



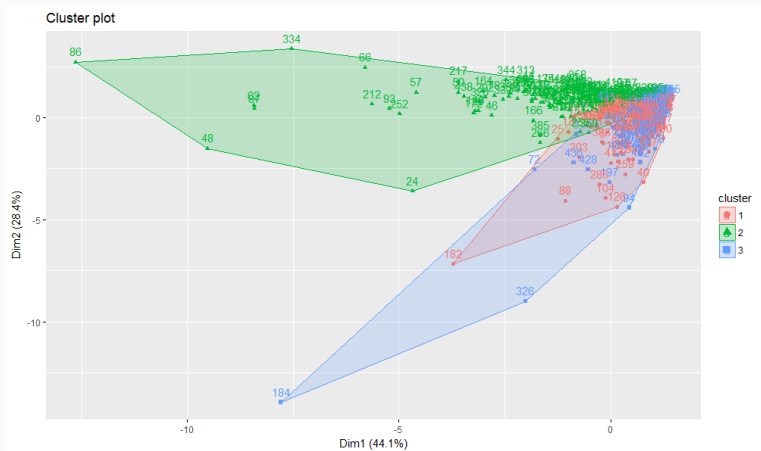
# INTERPRETATING CLUSTERS...



# INTERPRETATING CLUSTERS...

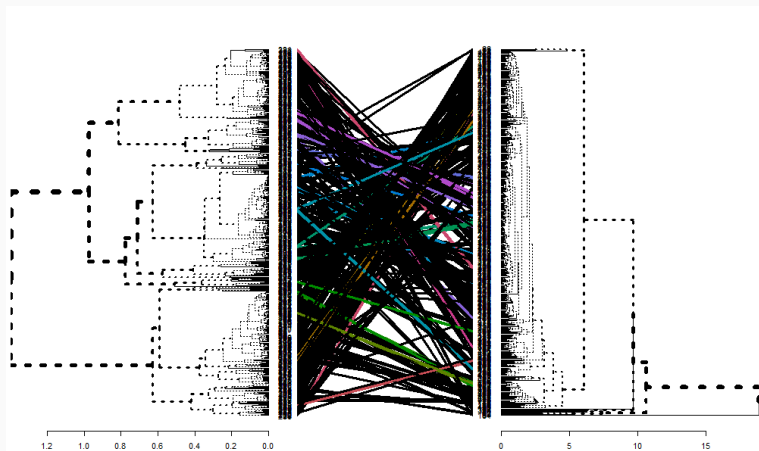


# INTERPRETATING CLUSTERS...





# TANGLEGRAM



# Conclusion

---

# THANK YOU!

Thank you all for your presence and for this opportunity!

Questions?