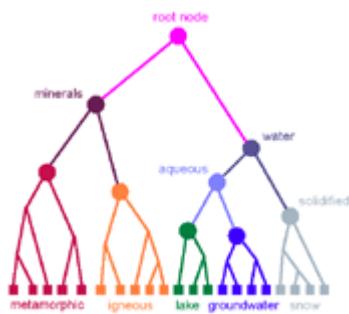
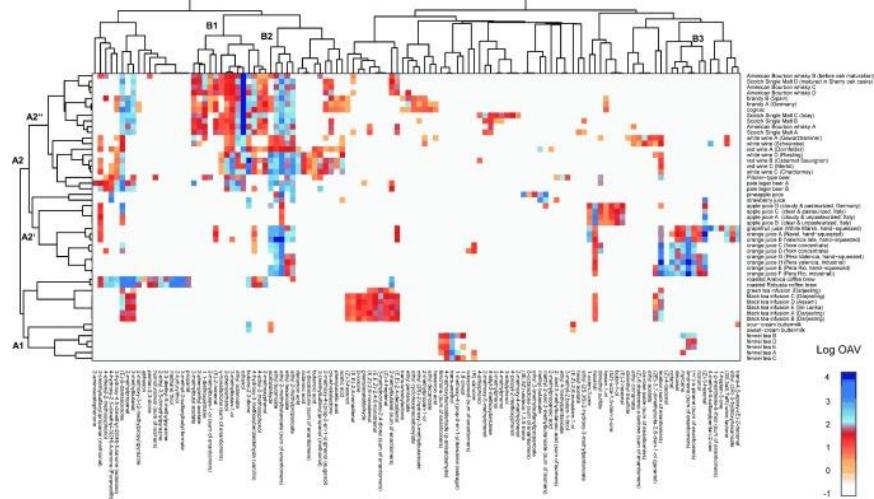
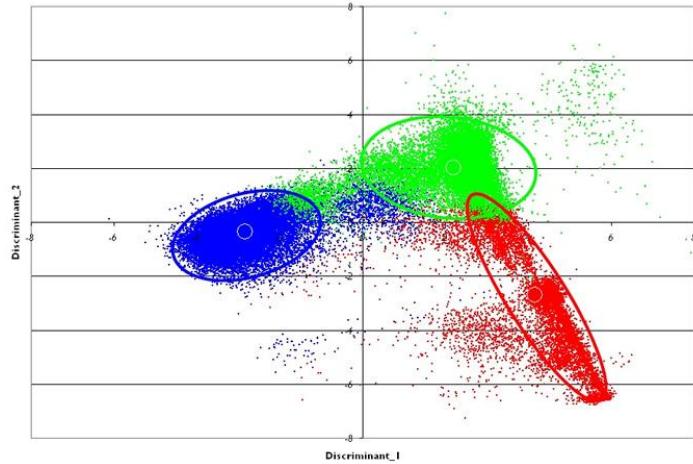


DISCOVERING PATTERNS

Module 7 : Cluster Analysis



Dr Rita Chakravarti
Institute of Systems Science
National University of Singapore
Email: rita@nus.edu.sg



© 2018 NUS. The contents contained in this document may not be reproduced in any form or by any means, without the written permission of ISS, NUS, other than for the purpose for which it has been supplied.

CLUSTER ANALYSIS

Agenda

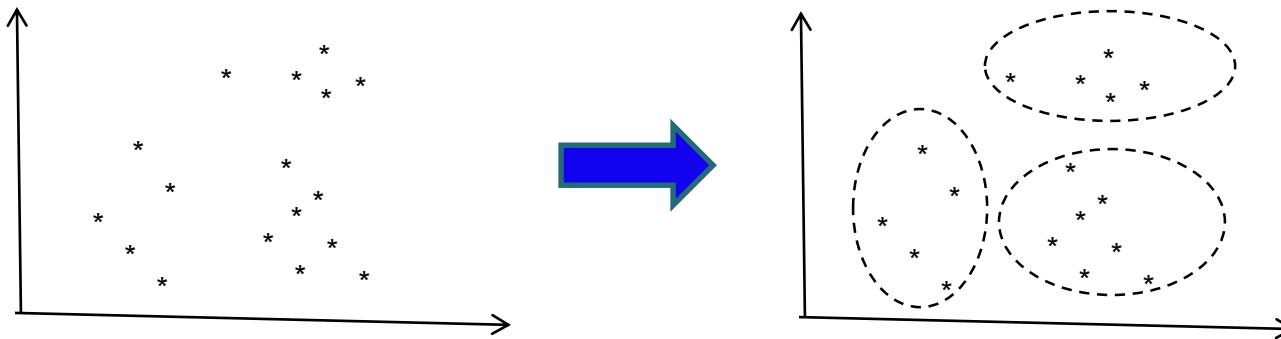
At the end of this module, you will be able to:

- Understand what is Cluster Analysis
- Application of Cluster Analysis
- Types of Clustering Methods
- How to Profile Clusters
- How to Validate the Created Clusters
- Understand how to perform these techniques using JMP, R & SPSS
- Limitations of Cluster Analysis

What is Cluster Analysis?

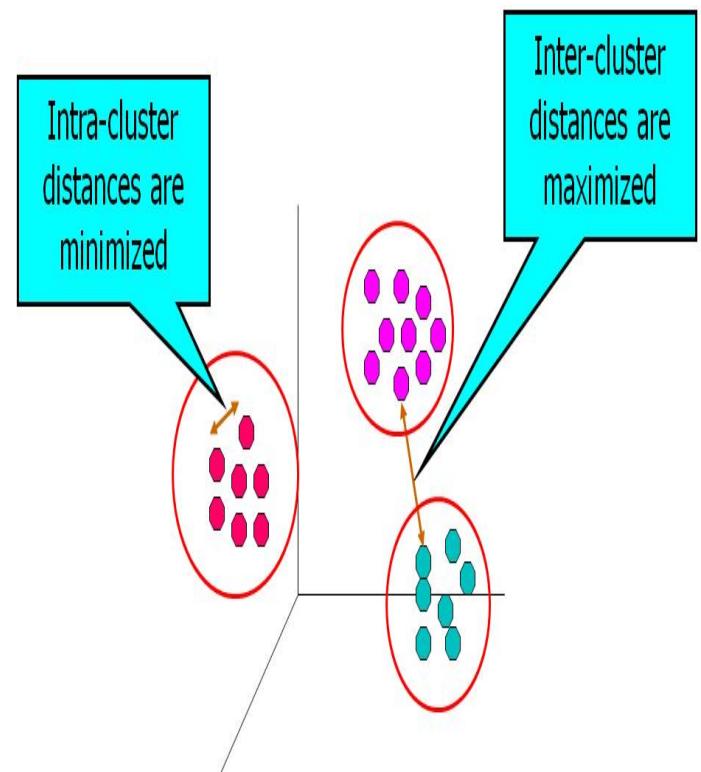
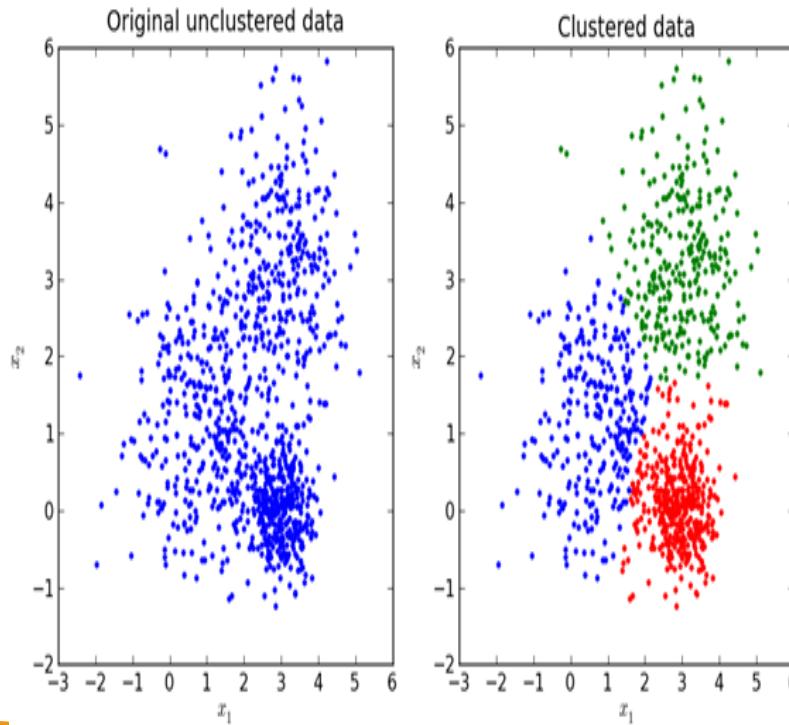
Goal of Clustering

- To segment the data into a set of **homogeneous** clusters (i.e.: members within same clusters are similar enough) of observations for the purpose of generating insight to **formulate cluster strategy**
- i.e. each object in a cluster should have more similarity within the same group such as characteristic, lifestyle, purchase preferences, product holdings etc. **but**
- Less similarity between groups



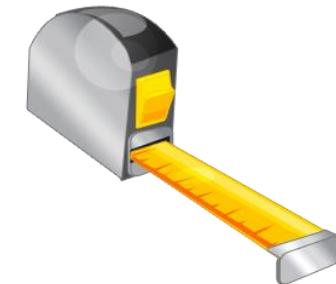
Characteristics of Cluster Analysis

**Cluster Analysis is an unsupervised learning -
there are no pre-defined classes**



Technicalities of Cluster Analysis

- **Need to begin by determining a distance measure** which will represent how far apart two data points are
- If each data point is represented by straightforward **numerical variables** then it is relatively easy;
- Measures of Distance
 - ***Euclidean distance.***
 - Squared (or absolute) Euclidean distance.
 - City-block (Manhattan) distance.
 - Chebychev distance.
 - Mahalanobis distance (D2)
 -
- In most cases we use Euclidean distance
- However it is important that the measures of distance are ***scaled or normalized*** so that one component of distance does not dominate
 - One measure of scale is the Z score
 - Transform the data so that the distances have mean 0, and variance of 1

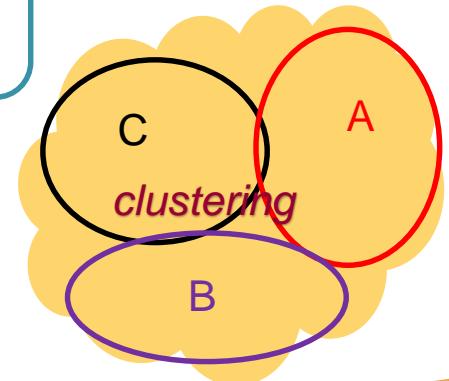


Utility of Cluster Analysis

Decide action

- Specific action on particular group, e.g.:
 - Targeting on group B (customers with potential interest in some new service)
 - Increasing influence to C (customers with possibility in developing some interest in new service)
 - Identify similar cities for test marketing
 - Identify group of firms that are prime target for takeover

Description of each group is achieved through **profiling**



Cluster Analysis : Be Mindful About Data Type

With **categorical*** variables it is more difficult to define distances

- If the categories can be related to levels of some measurement then we can define distance/association measures based on this
- Example: Your measurement scale may be (**Likert Scale**)
 - Strongly like
 - Like
 - Neutral
 - Dislike
 - Strongly dislike
- If one object scores “Like” and another scores “Neutral” then the distance could be 1
- If one object scores “Like” and another scores “Dislike” then the distance could be 2

	Strongly Disagree	Disagree	Undecided	Agree	Strongly Agree
Scale Week is a worthwhile feature on The Research Bunker Blog.	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>
I would like to read more posts about survey rating scales.	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input checked="" type="radio"/>
Vance Marriner is, without a doubt, the most insightful contributor to The Research Bunker Blog.	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>

*<https://pdfs.semanticscholar.org/aa3c/6df25998ed83e005bb5ea13acdcf5d58b500.pdf>

Distance Measure Used For Clustering

Euclidean Distance is the most popular distance measure

- Given two cases i , and j (in the p -dimensional space) the distance is defined by

$$d_{ij} = \sqrt{(x_{i1} - x_{j1})^2 + (x_{i2} - x_{j2})^2 + \dots + (x_{ip} - x_{jp})^2}$$

- Euclidean distance is **highly scale dependent**.
- Normalization on continuous measurements is recommended before computing the Euclidean distance, to convert all measurements to the same scale

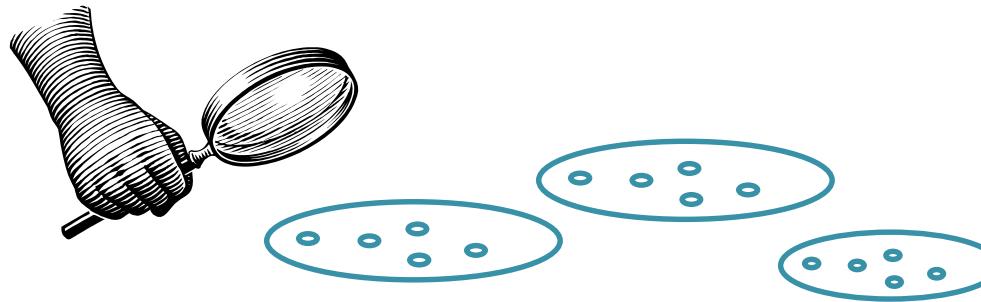
Other measures

- Manhattan distance for numerical data
- There are also measures for categorical* data and mixed data

*Reference: <http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.140.8831&rep=rep1&type=pdf>

Cluster Interpretation For Strategy Formulation

Interpreting Clusters



Can the clustering be explained in practical terms?

Examine the value of each variable of the cluster centroid, (the cluster profile) of each cluster.

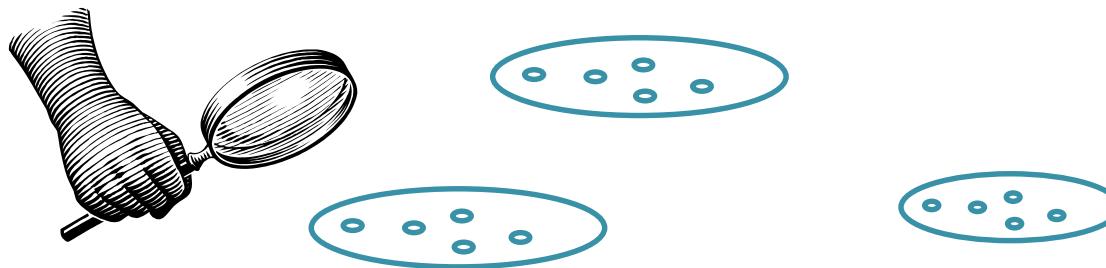
Look at the distinguishing characteristics of each cluster's profile and identifying substantial differences between clusters

Cluster solutions failing to show substantial variation between clusters indicate that this may be a spurious clustering

Other cluster solutions should be examined.

Interpretation May Require Research!!

Interpreting Clusters



The clusters should also be assessed against the analyst's



If no practical interpretation can be offered, more original research may be needed into the background of the data



May use case studies, specific measures, or pre-existing labels

Prior expectations based on theory

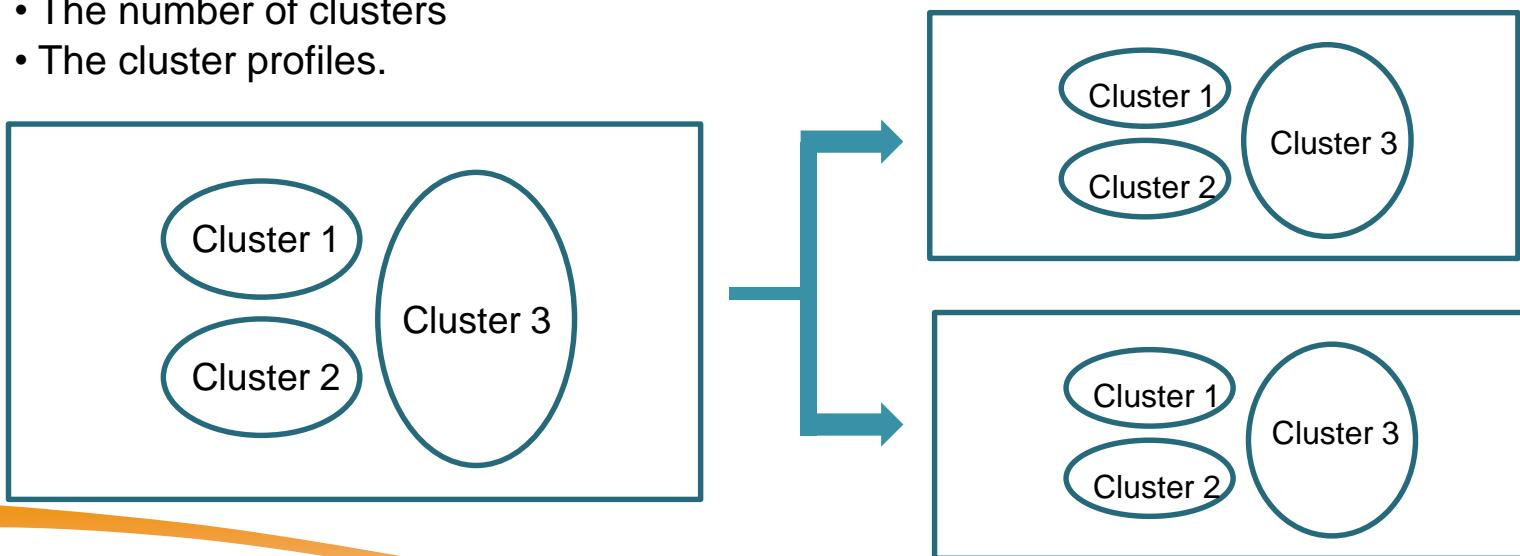
Or practical experience in previous work

Validation of Clusters

Validation is essential in cluster analysis since the clusters are descriptive of structure and require additional support for their relevance:

One method is Cross-validation

- Cross-validation empirically validates a cluster solution by
- Creating two sub groups from the original group of data points by randomly splitting the group
- Perform cluster analysis on **each** group
- Comparing the cluster solutions from each subgroup for consistency with respect to the
 - The number of clusters
 - The cluster profiles.

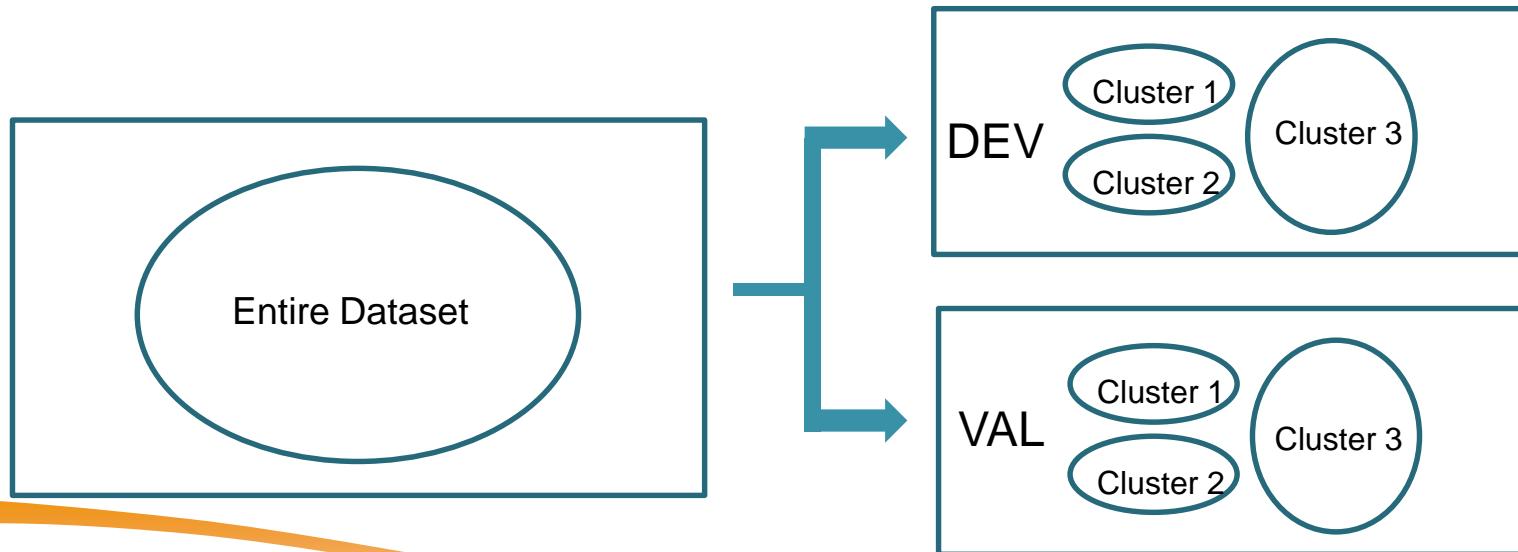


Validation of Clusters

Second method is to split randomly the clustering data set into two parts :

Development and Validation (50:50)

- Perform cluster analysis on development sample **only**
 - Create cluster profiles
 - Validate the results on the validation sample using development sample clustering algorithm (which becomes the implementation code)



Profiling of Clusters

CLUSTER PROFILING

Cluster 1	Cluster 2	Cluster 3	Cluster 4
<ul style="list-style-type: none">• Discount Hungry• Buys large amount• Very consistent• On an average pays less	<ul style="list-style-type: none">• Buys less• Somewhat consistent• On an average pays less	<ul style="list-style-type: none">• Not so discount hungry• Buys less• Very Consistent• On an average spends good	<ul style="list-style-type: none">• Not so discount hungry• Buys somewhat good amount• Somewhat consistent• On an average spends good

Profiling of Clusters should be done with those variables where their pattern or distribution is different within each cluster

e.g. If age is a clustering variable and the age distribution is same in the parent (development) data set and within all clusters then age is not a discriminating variable and **should not** be used for profiling

Efficient Method of Cluster Analysis

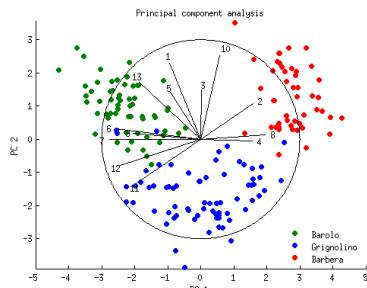
Data Exploration

Dimension Reduction

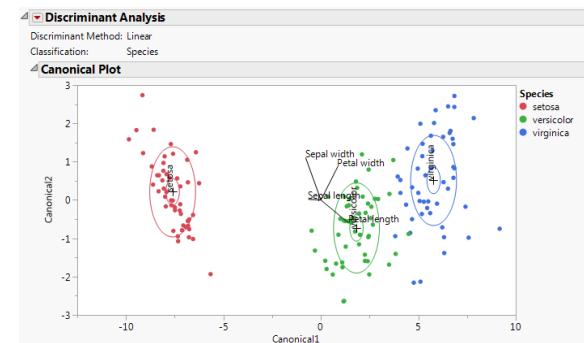
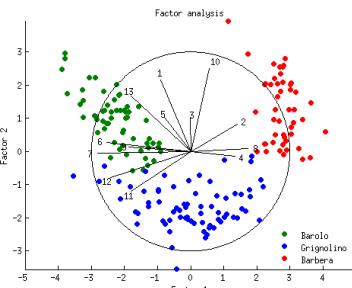
Clustering

Profiling

- **Factor analysis** reduces the dimensions (thus the number of variables), and minimizes multicollinearity effects.
- **Cluster analysis** identifies the grouping (thus structures within the data)
- **Discriminant analysis** checks the goodness of fit of the model and profiles the clusters.



Factor analysis



Discriminant analysis

How does Cluster Analysis help us?

Determine if different clusters exist and if there is a “natural” grouping for members within the data set.

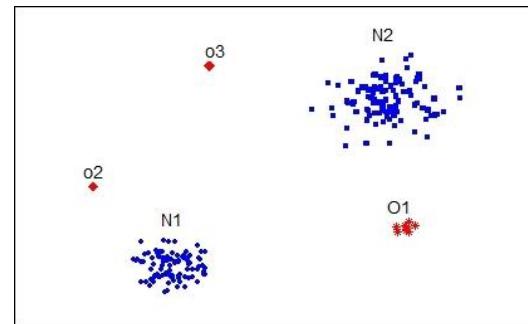
Identify those variables (or ranges of variables) that most strongly define the clusters.

Determine the rationale underlying the clusters. The results of the clusters should make sense from a business point of view & should generate business opportunities.

Identify anomalies or outliers; i.e. members who do not fit into any cluster

Suggests a classifying scheme for new data members

Dealing with Outliers



Outliers are data points that do not naturally fit into any cluster

- They only join groups at the end of the clustering process; i.e. when everything fits into one cluster

They may distort your clustering by changing your similarity measures

Usually they have some unusually high or low variable value

It may be necessary to remove these from your analysis to get a more well-defined clustering

But it may be good to investigate where the outlier came from

- Was it faulty data collection or transcription?
- Or was the outlier a representative of some different, new population?

Application of Cluster Analysis

Cluster analysis has been largely used in research such as:

- **Healthcare/Medical Research**
 - e.g. classifying patients on the basis of clinical and/or laboratory type observations.
- **Urban Planning**
 - e.g. Identifying groups of houses according to their house type, value, and geographical location

Application of Cluster Analysis

Marketing

- Clustering of neighborhoods using U.S. postal zip codes has been used successfully to group neighborhoods by lifestyle
- Group neighborhoods into 40 clusters using various measures of consumer expenditure and demographics.
- Knowledge of lifestyles can be used to estimate the potential demand for products (e.g. sports utility vehicles) and services (e.g. pleasure cruises)

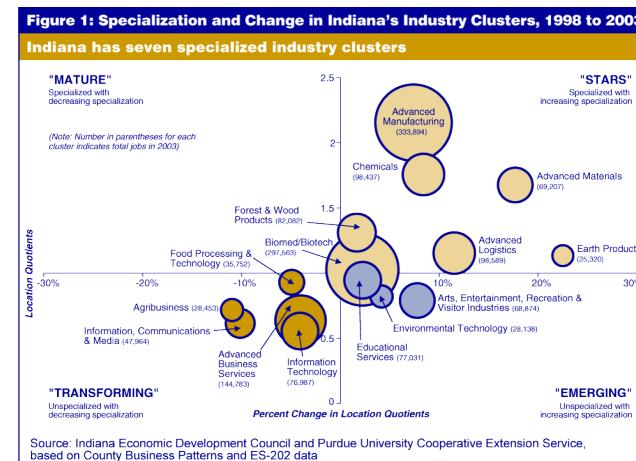


<https://ocw.mit.edu/courses/sloan-school-of-management/15-062-data-mining-spring-2003/lecture-notes/lec11.pdf>

Application of Cluster Analysis

Industry analysis

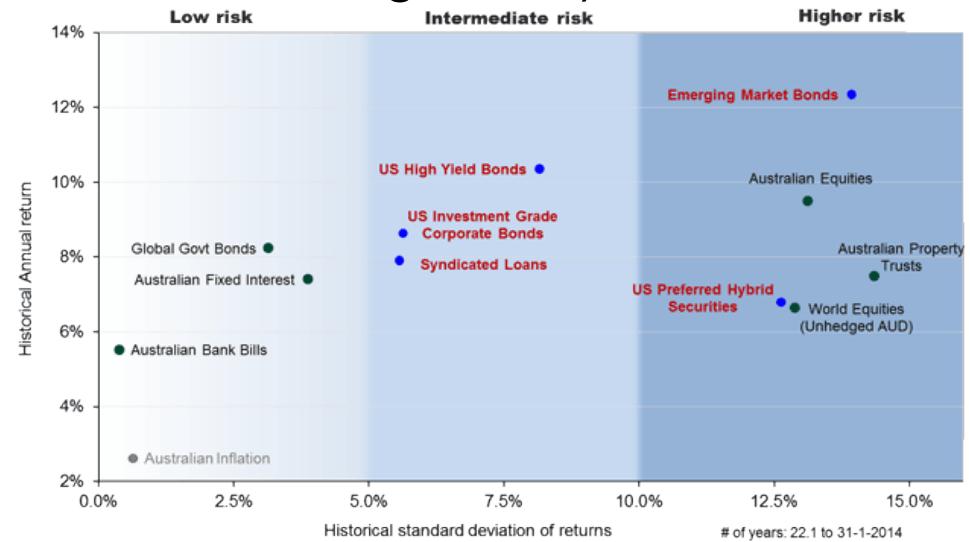
- For a given industry, we are interested in finding groups of similar firms based on measures such as growth rate, profitability, market size, product range, and presence in various international markets.
- These groups can then be analyzed in order to understand industry structure and to determine, for instance, who is a competitor, and how can we capture new opportunity



Application of Cluster Analysis

Finance: Balancing Portfolios

- Given data on a variety of investment opportunities (e.g., stocks), one may find clusters based on performance variables such as return, volatility, and other characteristics, such as industry and market capitalization.
- Selecting securities from different clusters can help create a balanced portfolio (for a better risk management)



Application of Cluster Analysis

Apparel Industry: Uniforms for women in U.S. Army

- The design of a new set of sizes for army uniforms for women
- The cluster analysis study came up with a new clothing size system with only 20 sizes to fit different body types
- 20 sizes are combinations of five measures: chest, neck, shoulder circumference, sleeve outseam, and neck-to-buttock length

*A completely new
insightful view can be
gained by examining
clusters of records*

(Berry & Lincoff 1997)



Data Sources for Clustering

Big Data

- Model systems may store and process very big data (e.g. weblogs) Example, Google's MapReduce Framework use Map Functions to distribute the computation across different machines.

Multimedia Data

- Clustering Video / Audio Clips
- Temporal ordering of records represents its meaning (Flickr, YouTube)

Biological Data

- Clustering Images (Position of Pixel represent its context)
- Gene/Protein Sequences

Social Network Data

- Clustering different typed nodes / links together (e.g. NetClus)

Customer Behavioural Data (Bank, Insurance, Telco....)

- Clustering different types of customers

Manufacturing Data

- Clustering different types of machines based on performance

TYPES OF CLUSTERING METHODS

Types of Cluster Analysis

Hierarchical Clustering

- Agglomerate Method
- Divisive Method

K-Means Method

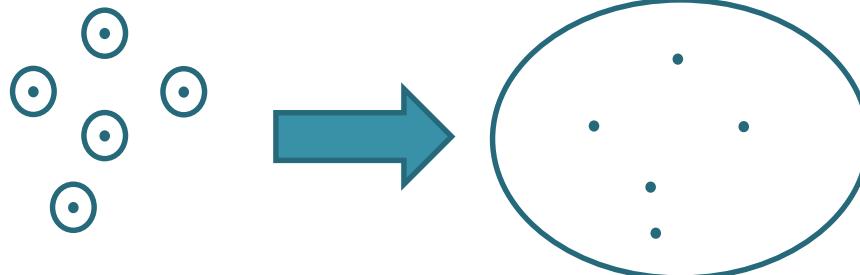
Kohonen's Clustering Method (Self Organizing Maps SOM)

K Nearest Neighbour (KNN)

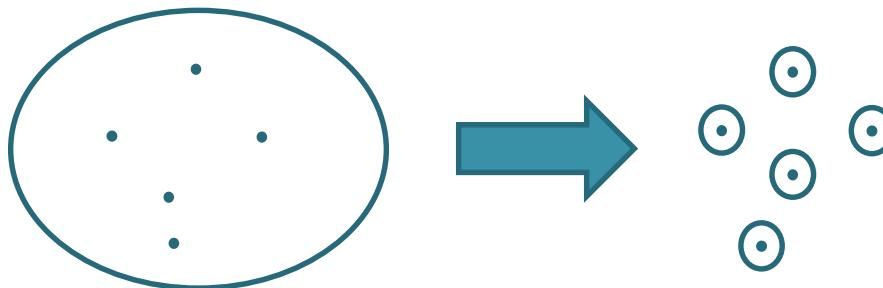
Two Step Clustering

Source: <http://www.stat.columbia.edu/~madigan/W2025/notes/clustering.pdf>

Two Approaches to Hierarchical Clustering



- **“Agglomerative (bottom up)” Methods**
 - We build up clusters from a group of disparate individuals

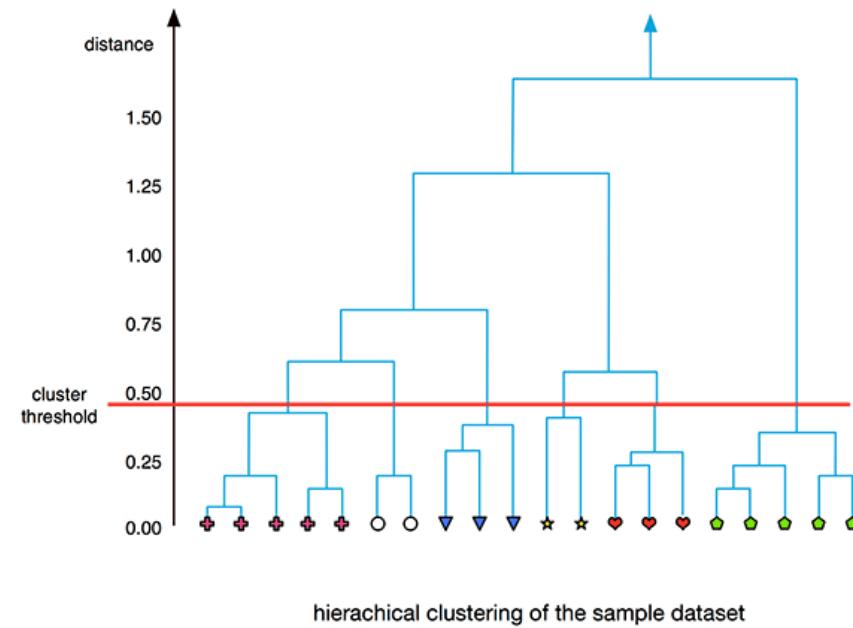
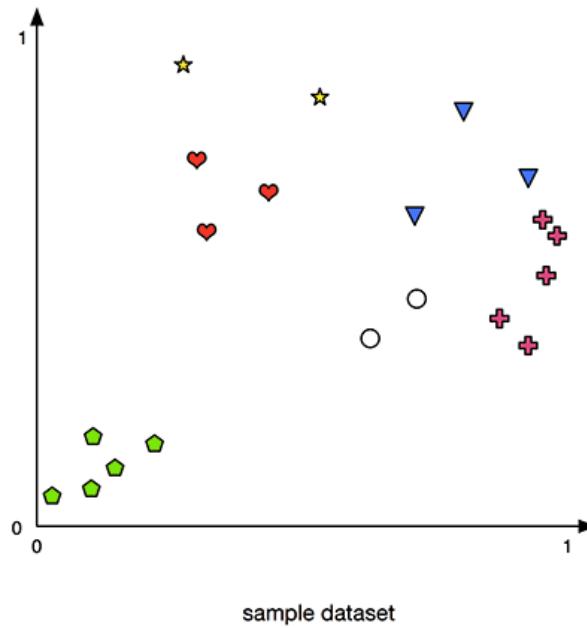


- **“Divisive(top down)” Methods (Partitional)**
 - Breakdown a single cluster (of all individuals) into several different clusters

Agglomerate or Hierarchical Clustering

Agglomerate or Hierarchical Clustering

- Start with each cluster comprising exactly one record and then progressively agglomerating (combining) the two **nearest** clusters until there is just one cluster left at the end, which consists of all the records.



Agglomerate or Hierarchical Clustering

Quite outdated nowadays

It is called hierarchical or agglomerate because it starts with a solution where each record comprises a cluster & gradually groups records up to the point where all of them fall into one supercluster.

In each step it calculates the distances between all pairs of records & groups the most similar ones

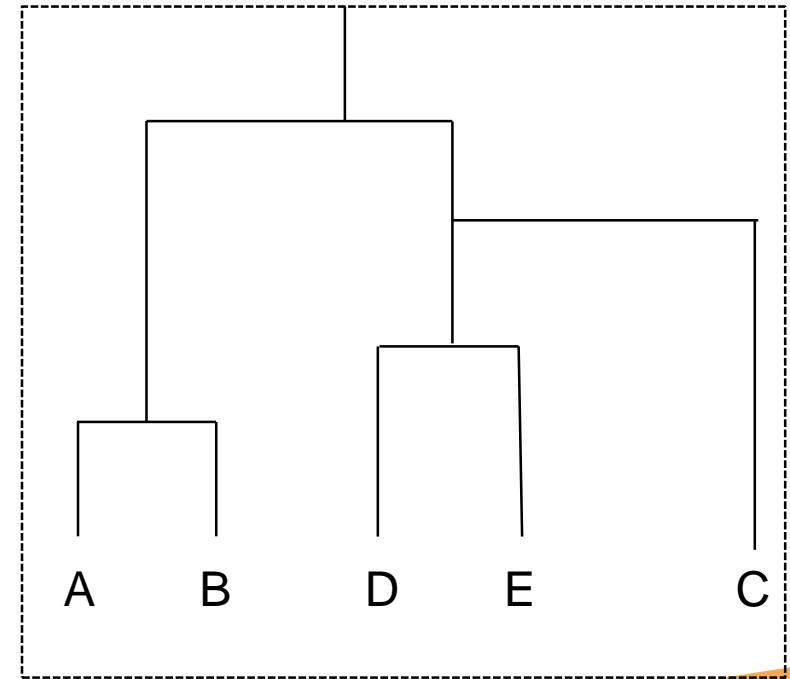
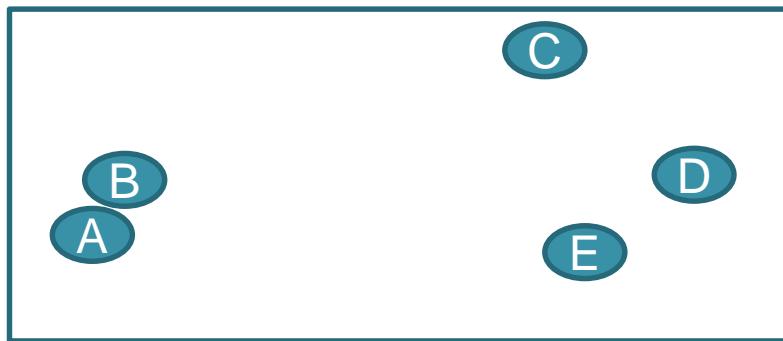
The analyst should consult this information, & identify the point where the algorithm starts to group disjoint cases, & then decide on the number of clusters to retain.

This algorithm cannot effectively handle more than a few thousand cases. Thus, it cannot be directly applied in most business clustering tasks.

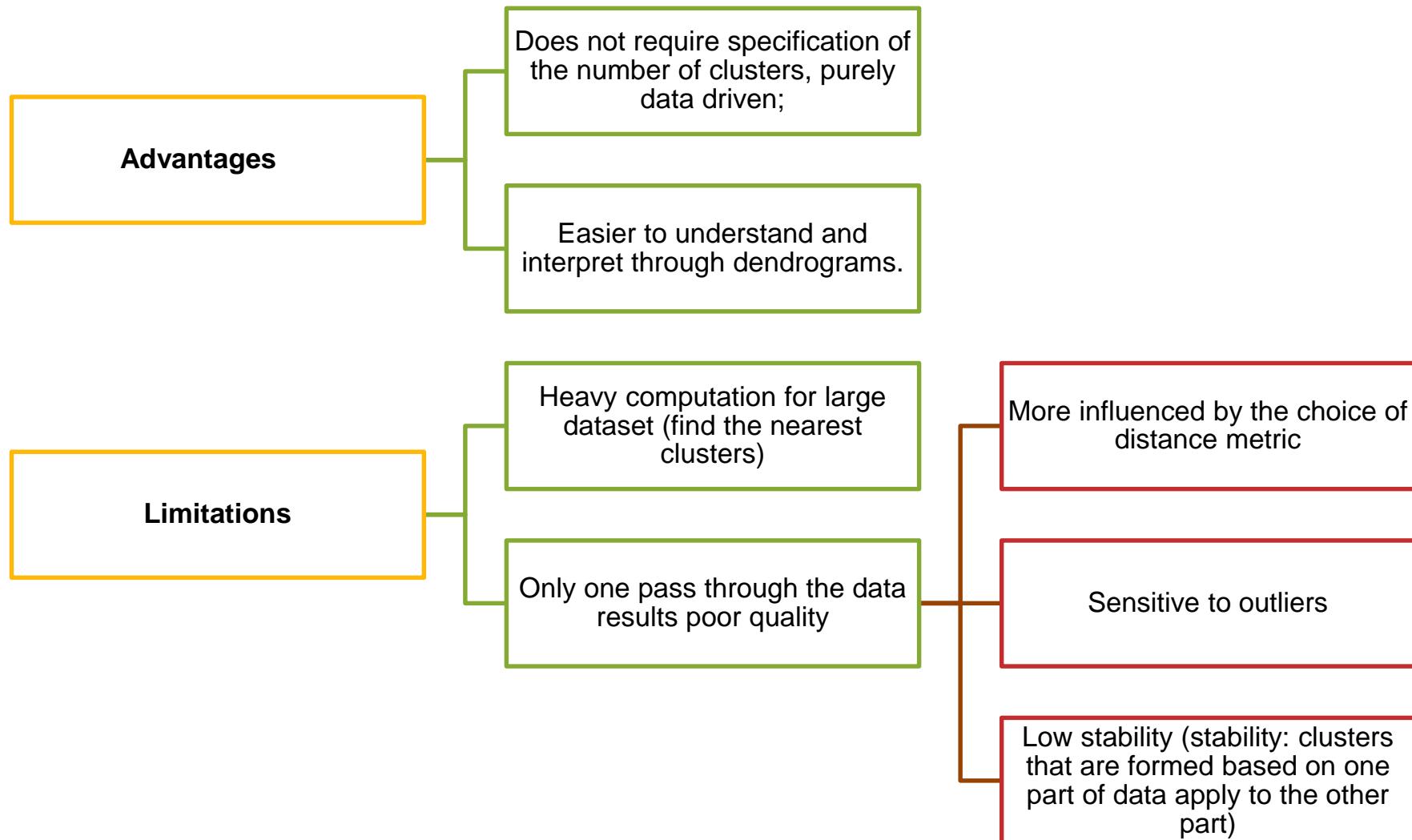
Agglomerate or Hierarchical Clustering

A dendrogram is a tree diagram used to illustrate the arrangement of the clusters produced by hierarchical clustering.

Example, suppose the below data is to be clustered using Euclidean distance as the distance metric.



Agglomerate or Hierarchical Clustering



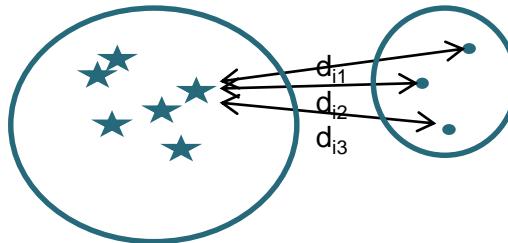
Agglomerate or Hierarchical Clustering

1. Start with all the data points

- Determine the distances between every data point
- Work out which two data points are closest together
- Combine these two data points together to form your first cluster

2. Then recalculate the distances

- Between every data point
- Also between every data point and every cluster
- Also between every cluster



3. Identify the two clusters with the smallest average distance between them

4. Combine these two clusters together to form a cluster

5. Repeat steps 2 to 4 until all data points are in a single cluster

Divisive Clustering Methods

A cluster hierarchy can also be generated top-down

This variant of hierarchical clustering is known as divisive clustering

- One starts with all objects in one cluster
- Then the cluster is split using a flat* cluster algorithm (typically used in document clustering)
- The procedure is applied recursively until each object is in its own singleton clusters
- Top-down clustering is conceptually more complex than bottom-up clustering since a second, flat clustering algorithm as a ``subroutine'' is required
- It has the advantage of being more efficient if one does not generate a complete hierarchy all the way down to individual objects
- k-means is perhaps the most widely used flat clustering algorithm due to its simplicity and efficiency

*<http://software.ucv.ro/~cmihaescu/ro/teaching/AIR/docs/Lab6-ClusteringWithKMeans.pdf>

Nonhierarchical Clustering

***k*-Means** is a typical nonhierarchical clustering algorithm

- The goal is to divide the sample into *predetermined* number k of non-overlapping clusters so that clusters are as homogeneous as possible with respect to the measurements used

Kohonen clustering (sometimes called *k*-Means with a twist)

- Using Kohonen network, which is an important unsupervised learning neural network performing *Self Organizing Map*
- Is more data driven, without the requirement of prespecified number of clusters

***k*-Nearest Neighbors algorithm (KNN):**

- It is a non-parametric method used for classification and regression.
- In both cases, the input consists of the k closest training examples in the feature space.
- The output depends on whether k -NN is used for classification or regression

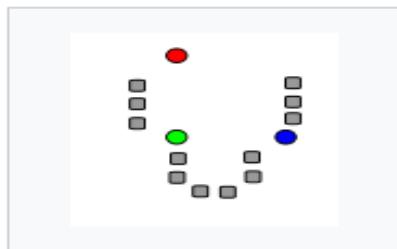
K – Means Clustering

This is an efficient & perhaps the fastest clustering algorithm that can handle both long (many records) and wide datasets (many input fields)

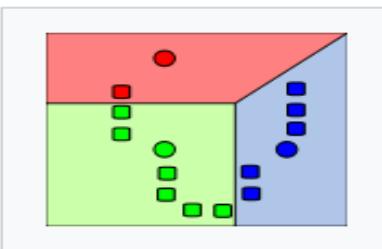
It is distance based and unlike the hierarchical algorithms, it does not need to calculate the distance between all pairs of records

The number of clusters to be formed is predetermined and specified in advance

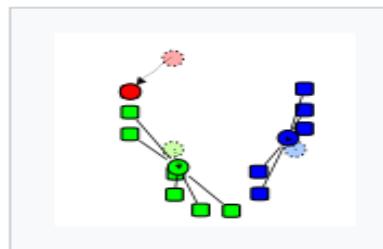
Demonstration of the standard algorithm



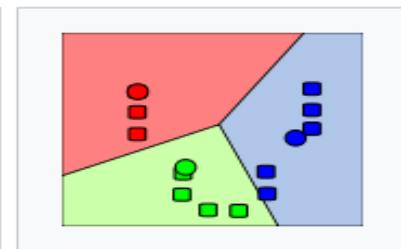
1. k initial "means" (in this case $k=3$) are randomly generated within the data domain (shown in color).



2. k clusters are created by associating every observation with the nearest mean. The partitions here represent the Voronoi diagram generated by the means.



3. The centroid of each of the k clusters becomes the new mean.

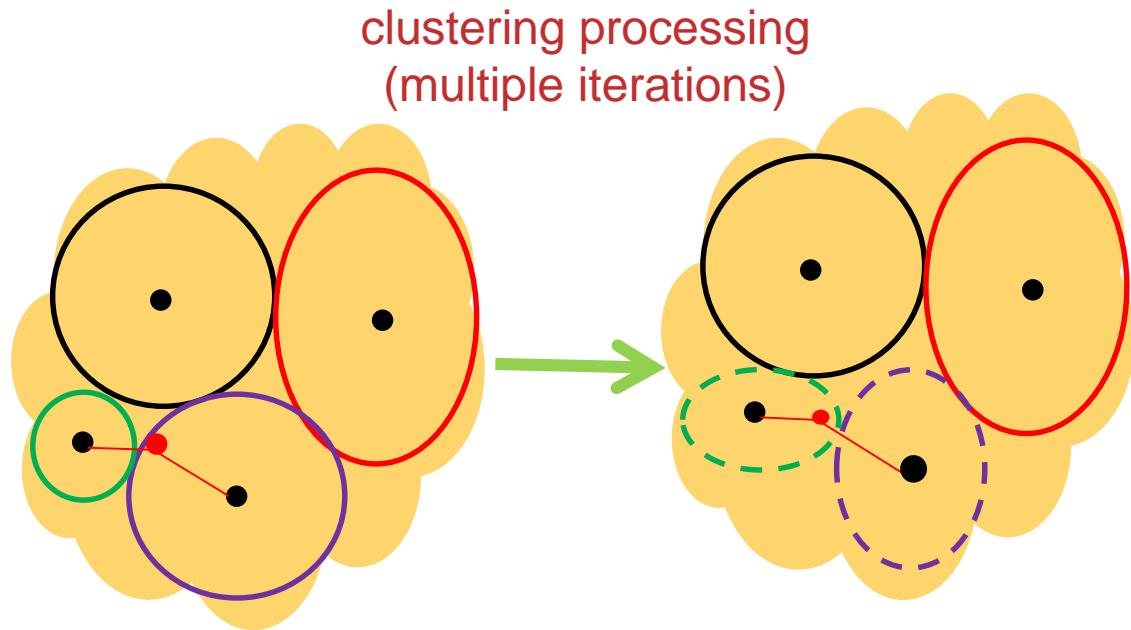


4. Steps 2 and 3 are repeated until convergence has been reached.

Reference : Applied Multivariate Statistical Analysis : Johnson & Wichern

K – Means Clustering

- Usually a number of different solutions should be tried & evaluated before approving the most appropriate.
- K-Means Clustering is best for handling **continuous** clustering fields



K – Means Clustering



Step 1 : Start with k initial clusters (user choose k)

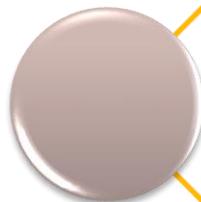


Step 2 : At every step, each data record is reassigned to the cluster with “closest” centroid



Step 3 : Re-compute the centroids of clusters that lost or gained a record, and

- repeat the previous step as the next iteration



Step 4 : Stop when moving any more records between clusters increases cluster dispersion

K – Means Clustering

The choice of the number of clusters can

- Either be driven by external consideration (e.g.: previous knowledge, practical constraints, etc)
- Or we can try a few different values for k and compare the resulting clusters

Initial partition into k clusters

- Any available external information suggesting a certain partitioning, or the centroids of the k clusters, should be used
- In case no information is available for the initial partition, the algorithm can rerun with different randomly generated starting partitions to reduce the chance of resulting poor solution

Kohonen Network (Self Organizing Maps)

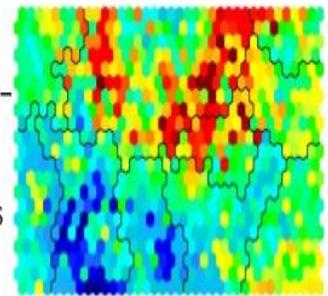
Kohonen Networks are based on Neural Networks & typically a two-dimensional grid or map of the clusters, hence the name self organising maps.

Kohonen Networks usually take a longer time to train than k-Means & Two Step algorithms, & provide a different view on clustering that is worth trying.

Self-Organising Maps

A Self-Organising Map (SOM) is a form of unsupervised neural network that produces a low (typically two) dimensional representation of the input space of the set of training samples.

- First described by Teuvo Kohonen (1982) ("Kohonen Map")
- Over 10k citations referencing SOMs – most cited Finnish scientist.
- Multi-dimensional input data is represented by a 2-D "map" of nodes
- Topological properties of the input space are maintained in map



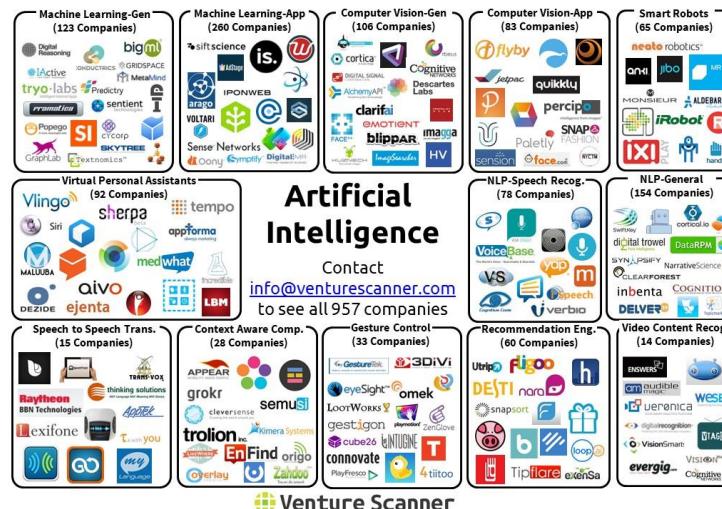
Dublin R / SOMs / Shane Lynn

5

Segmentation – A Strategic Tool – Application of Clustering Technique

One size does not fit all

Segmentation is the practice of dividing a customer base/market/softwares/tools/machines into groups that are similar in specific ways so that the **description/profile** of the segment can be used for strategic purposes



Bank Case Study

Behaviour Segmentation For Salaried Customers

Behavior Segmentation - Approach

“Segmentation”



Aim: To identify actionable customer clusters based on their past behavior

Step 1: Group the customers into clusters based on a set of criteria

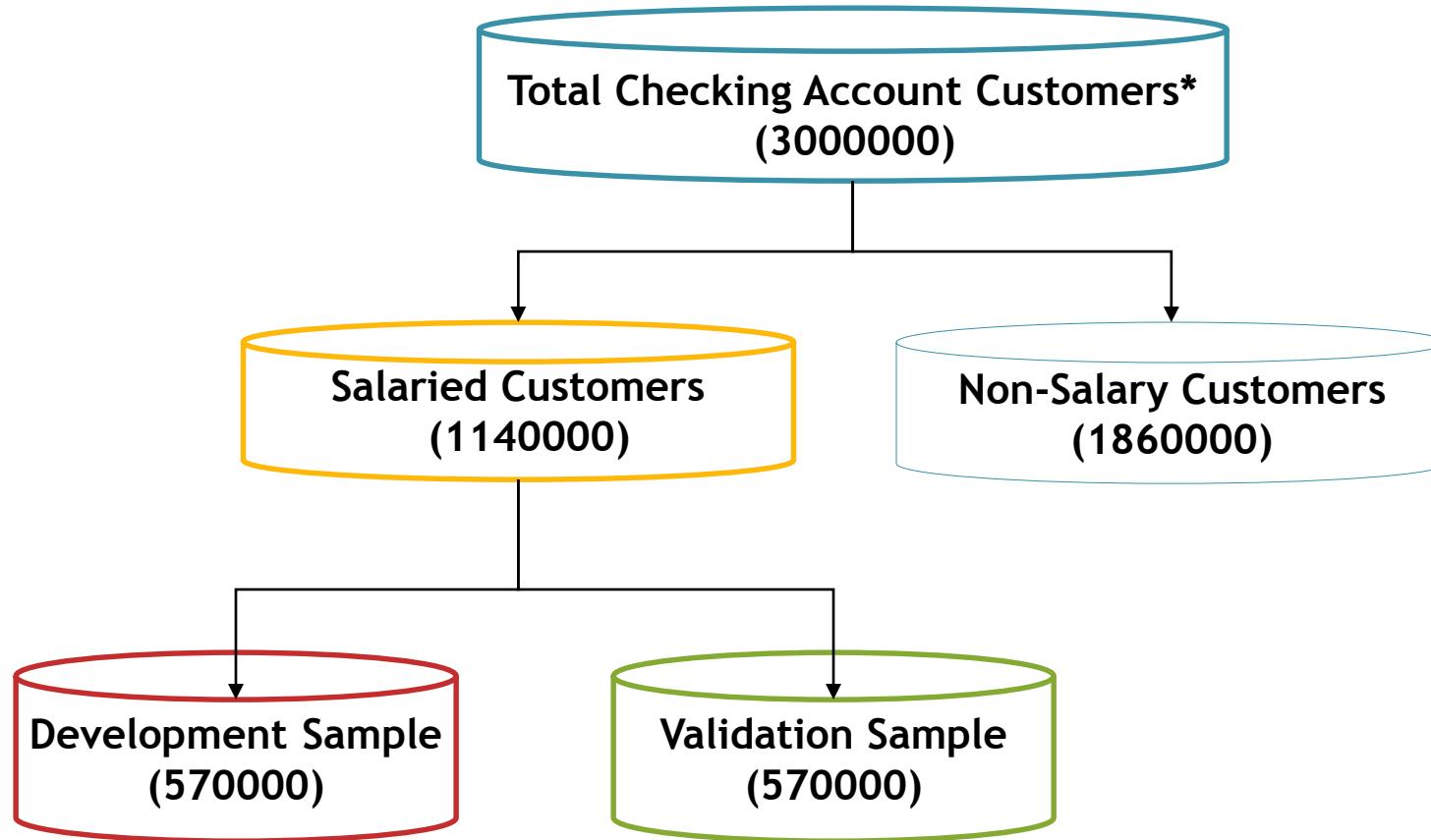
Step 2: Compare & identify the clusters

Step 3: Profile each cluster to know our customer better

Step 4: Develop strategy for each cluster

Step 5: Develop suitable actionable insights for the clusters

Segmentation was performed on a sample of 570K customers

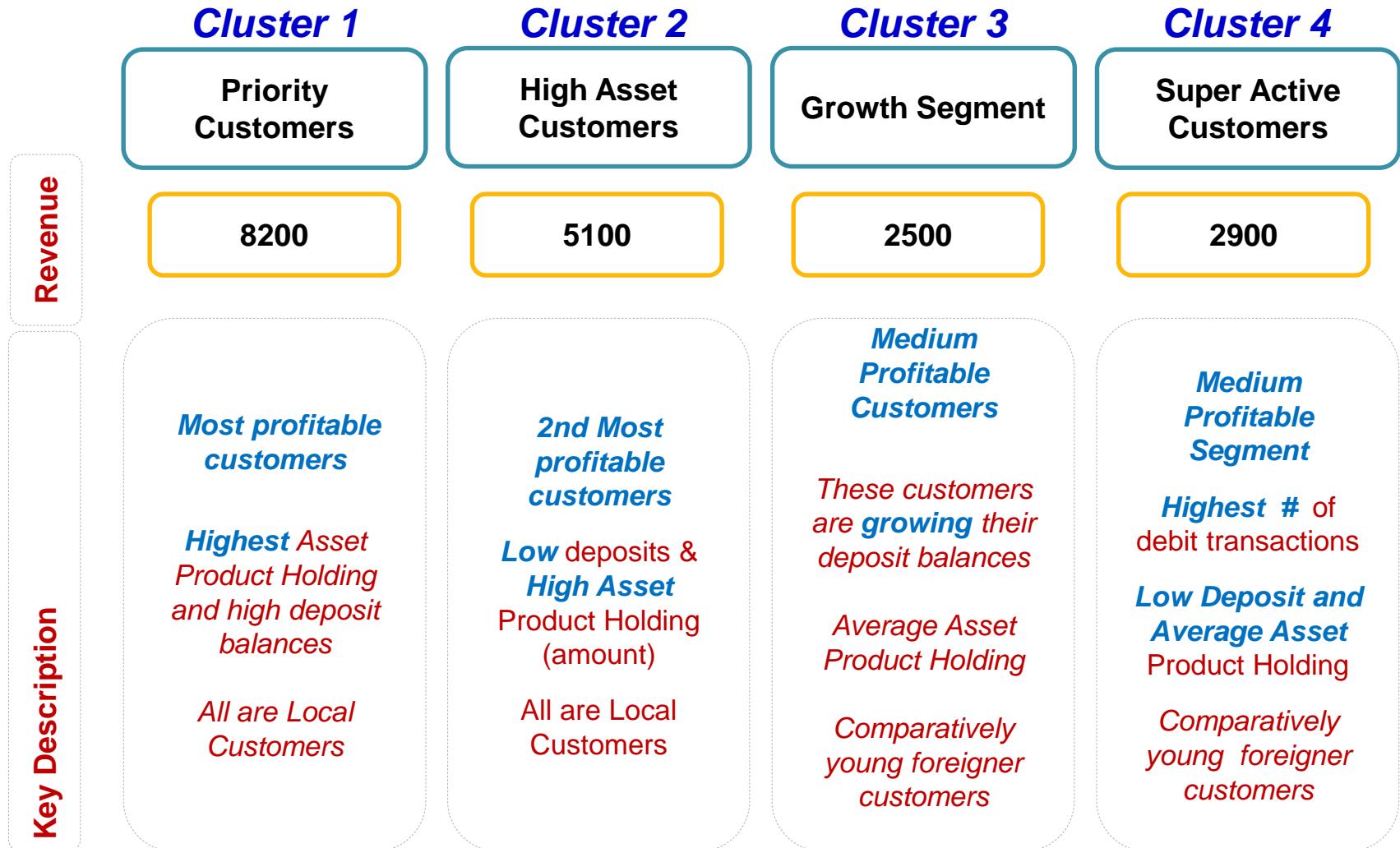


*As of October 04

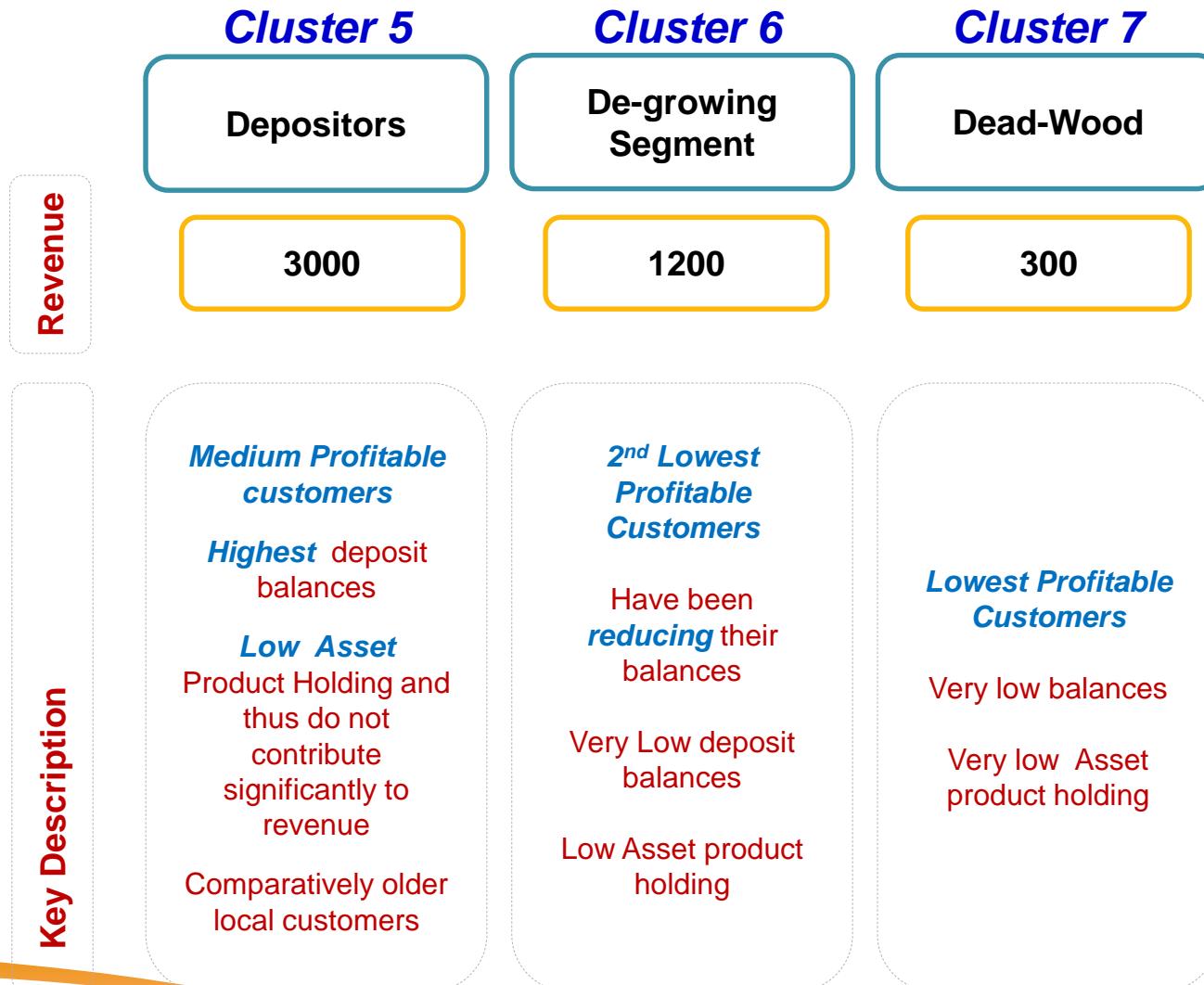
7 Key Segments were identified in the Salaried Base



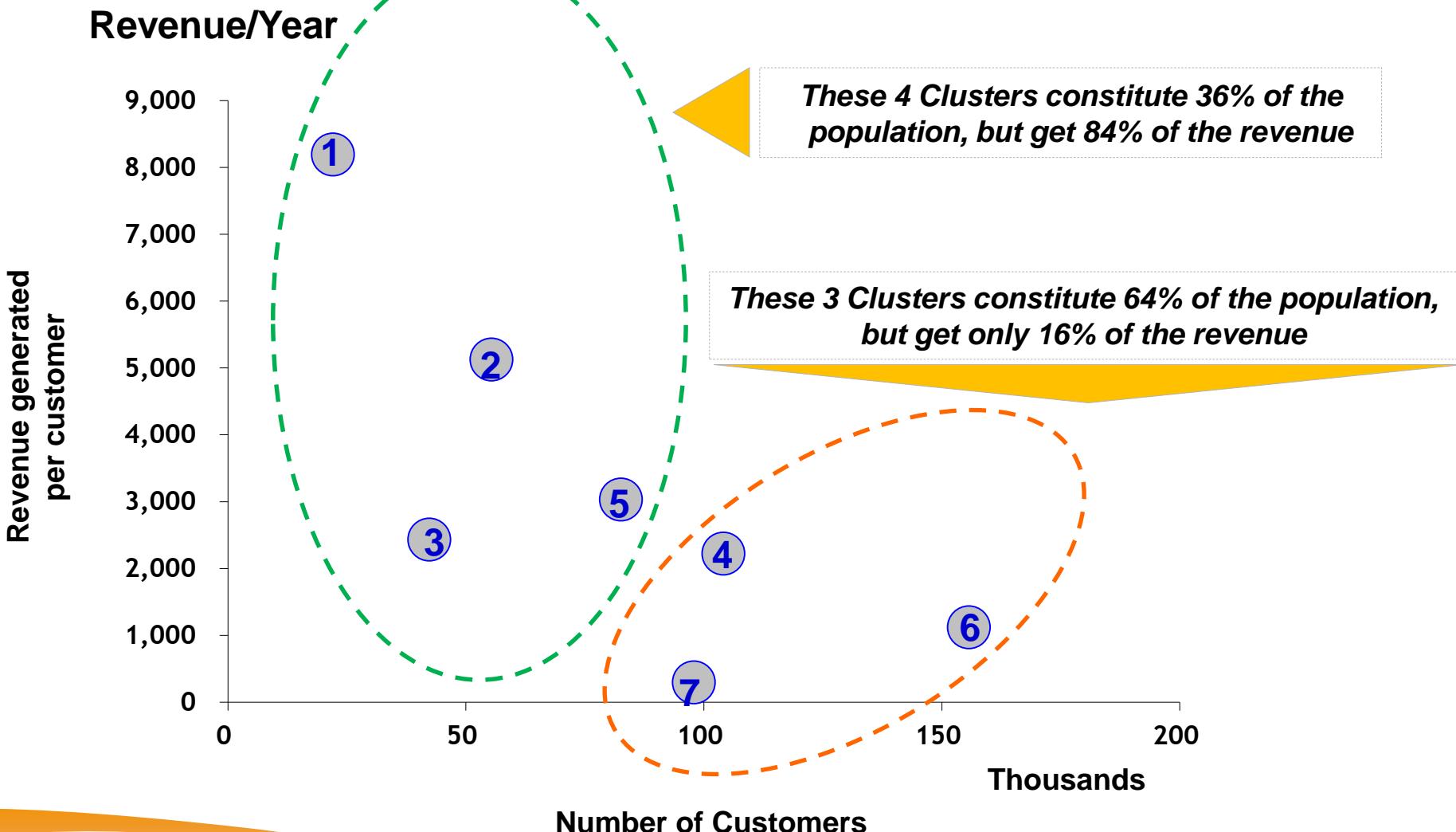
Description of each cluster



Description of each cluster



Top 4 clusters share as high as 84% of revenue with 36% of the population



Profile of 7 clusters based on the behavior across products...

Cluster 1 Cluster 2 Cluster 3 Cluster 4

	Priority Customers	High Asset Customers	Growth Segment	Super Active Customers	Population Average
Size	4%	10%	8%	19%	100%
Avg. Revenue per customer	8200	5100	2500	2900	2223
Deposit Avg. Month-End balance	15382	3951	3259	6879	9261
Deposit Avg. bal growth	0.5	0.5	0.8	0.4	0.45
Liabilities - # of Debit Transactions/ Month	271	167	136	284	144
Tenure – Checking Account (months)	145	116	87	95	96
Loan holding	100%	100%	93%	73%	52%
Overdraft holding	95%	96%	89%	69%	49%
Card Holding	45%	36%	24%	21%	17%
Total Loan – Disbursed Amount	290833	195599	89579	72422	66680
Overdraft – Disbursed Amount	264132	176187	82650	66869	60728
Other Asset products – Disbursed Amt.	26701	19412	6928	5554	5952
% Local	100%	100%	98%	98%	90%
Age (yrs)	39	35	31	31	35

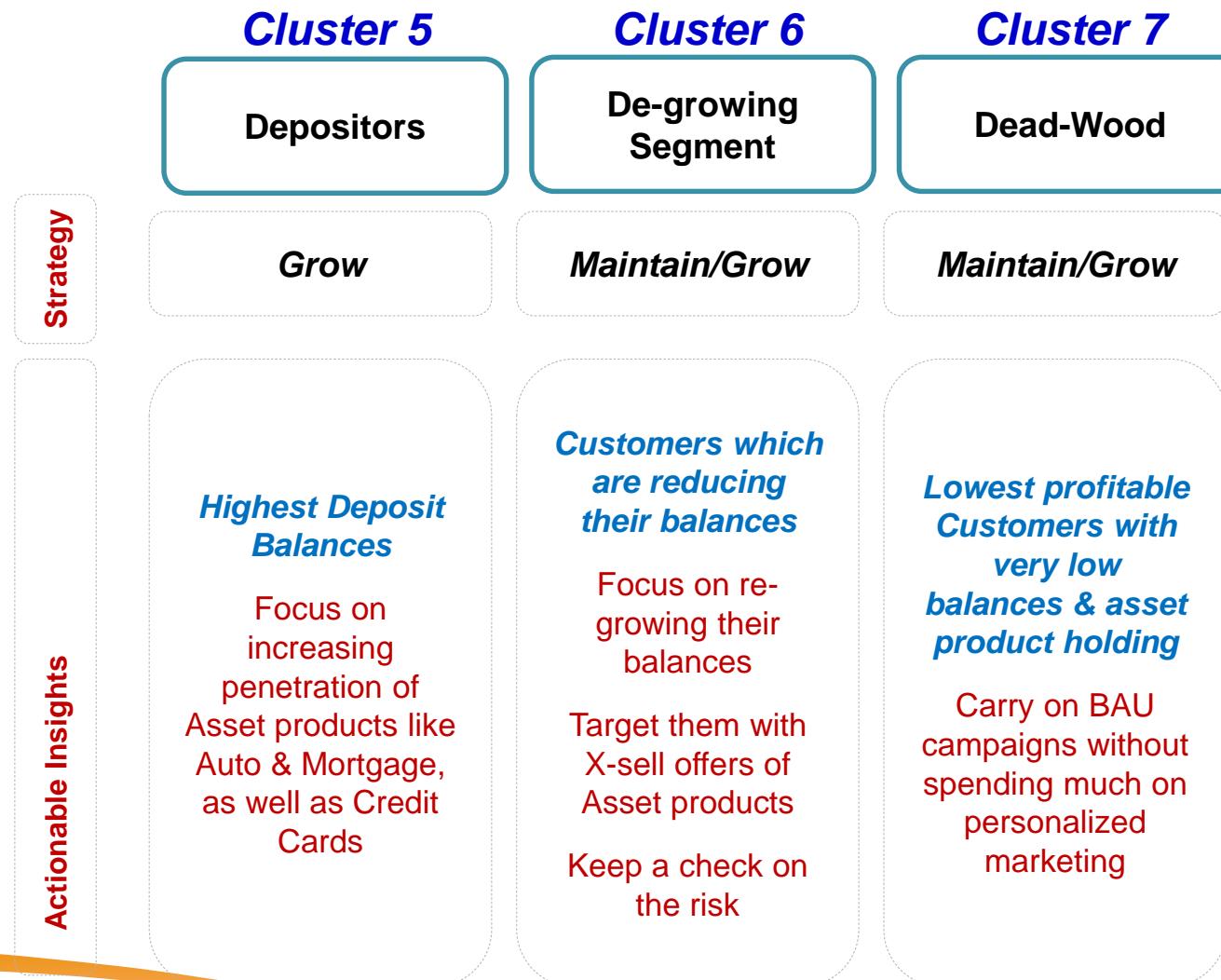
Profile of 7 clusters based on the behavior across products...

	<i>Cluster 5</i>	<i>Cluster 6</i>	<i>Cluster 7</i>	
Size	Depositors	De-growing Segment	Dead-Wood	Population Average
Avg. Revenue per customer	15%	28%	17%	100%
Deposit Avg. Month-End balance	3000	1200	300	2223
Deposit Avg. bal growth	39290	1794	2519	9261
Liabilities - # of Debit Transactions/ Month	0.5	0.2	0.7	0.45
Tenure – Checking Account (months)	123	90	66	144
Loan holding	118	86	80	96
Overdraft holding	23%	45%	11%	52%
Card Holding	21%	42%	10%	49%
Loan – Disbursed Amount	11%	12%	4%	17%
Overdraft – Disbursed Amount	26990	40576	5920	66680
Other Asset products – Disbursed Amt.	9,176	36881	5363	60728
% Local	2379	3696	557	5952
Age (yrs)	89%	87%	77%	90%
	39	35	37	35

Recommendations for each Segment

	Cluster 1	Cluster 2	Cluster 3	Cluster 4
Strategy	Priority Customers	High Asset Customers	Growth Segment	Super Active Customers
Actionable Insights	<p>Most Profitable customers</p> <p>Engage them with personal care & keep them happy</p> <p>Further enhance the relationship by offering investment & insurance products</p>	<p>2nd Most profitable customers</p> <p>Need to engage with these customers to keep them happy</p> <p>Focus on growing their deposit balances with the bank</p>	<p>A Growing Segment</p> <p>As they are growing their balances, we can proactively target them with Up-sell/Cross-sell offers to increase the penetration of asset products</p>	<p>Highly active in terms of debit transactions</p> <p>Can encourage them to maintain their balances</p> <p>Increase penetration of Asset products</p> <p>Can be good target for Cards</p>

Recommendations for each Segment



Cluster Analysis for Customer Segmentation

Advantages

- Clustering of customers according to their attribute preferences
- Cluster customer with similar behaviours/characteristics
- Clusters of similar brands/products can help identifying competitors / market opportunities

Disadvantages

- Choice of cluster-forming variables can be random
- Determine the right number of cluster often time-consuming
- Highly dependent on the analyst's interpretation

HANDS ON EXERCISE

Statistical Tools- JMP, R, SPSS Modeler

Cluster Analysis Using JMP Pro 13

Data Set: European Countries

K Means Cluster by JMP

K MEANS CLUSTER

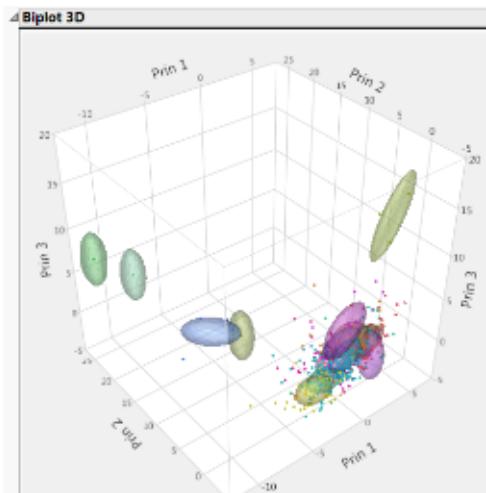
[PREVIOUS](#) | [NEXT](#)

Group Observations Using Distances

Use the K Means Cluster platform to group observations that share similar values across a number of variables. Use the *k-means* method with larger data tables, ranging from approximately 200 to 100,000 observations.

The K Means Cluster platform constructs a specified number of clusters using an iterative algorithm that partitions the observations. The method, called *k-means*, partitions observations into clusters so as to minimize distances to cluster centroids. You must specify the number of clusters, *k*, in advance. However, you can compare the results of different values of *k* to select an optimal number of clusters for your data.

3D Biplot

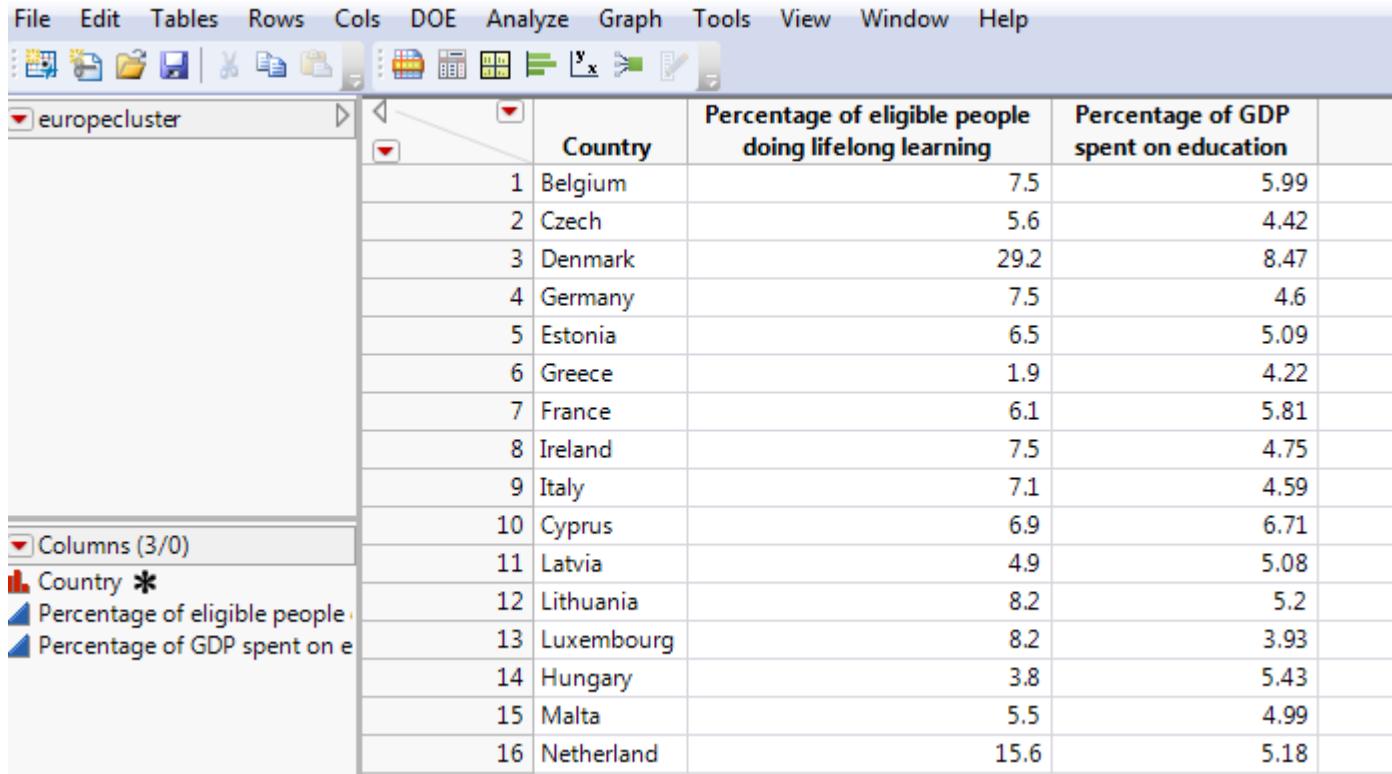


Summary of Clustering Methods

Method	Data Type or Modeling Type	Data Table Size	Specify Number of Clusters
Hierarchical Cluster	Any	With Fast Ward, up to 200,000 rows	No
		With other methods, up to 5,000 rows	
K Means Cluster	Numeric	Up to millions of rows	Yes
Normal Mixtures	Numeric	Any size	Yes
Latent Class Analysis	Nominal or Ordinal	Any size	Yes

<http://www.jmp.com/support/help/13/>

Select Appropriate Data Set



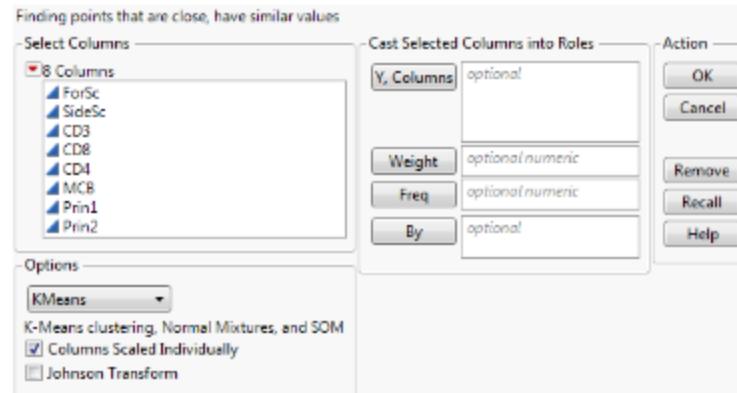
The screenshot shows a software interface for data analysis, likely SPSS or a similar tool. The menu bar includes File, Edit, Tables, Rows, Cols, DOE, Analyze, Graph, Tools, View, Window, and Help. The toolbar has various icons for file operations and data manipulation. A left sidebar shows a tree view with 'europecluster' selected, and a list of columns: Country, Percentage of eligible people doing lifelong learning, and Percentage of GDP spent on education. Below this is a list of variables: Country *, Percentage of eligible people doing lifelong learning, and Percentage of GDP spent on education. The main area displays a table with 16 rows, each representing a country and its two corresponding metrics. The table is as follows:

	Country	Percentage of eligible people doing lifelong learning	Percentage of GDP spent on education
1	Belgium	7.5	5.99
2	Czech	5.6	4.42
3	Denmark	29.2	8.47
4	Germany	7.5	4.6
5	Estonia	6.5	5.09
6	Greece	1.9	4.22
7	France	6.1	5.81
8	Ireland	7.5	4.75
9	Italy	7.1	4.59
10	Cyprus	6.9	6.71
11	Latvia	4.9	5.08
12	Lithuania	8.2	5.2
13	Luxembourg	8.2	3.93
14	Hungary	3.8	5.43
15	Malta	5.5	4.99
16	Netherland	15.6	5.18

Data File: European-Countries

K Means Cluster - Launch

K Means Cluster Launch Window



Y, Columns

The variables used for clustering observations.

Note: K-Means clustering supports only numeric columns.

Weight

A column whose numeric values assign a weight to each row in the analysis.

Freq

A column whose numeric values assign a frequency to each row in the analysis.

By

A column whose levels define separate analyses. For each level of the specified column, the corresponding rows are analyzed. The results are presented in separate reports. If more than one By variable is assigned, a separate analysis is produced for each possible combination of the levels of the By variables.

Select Clustering (K Means) Method

The screenshot shows the SPSS interface with the Analyze menu open. The Clustering option is selected, and the K Means Cluster option is highlighted. A tooltip provides the following description: "Clusters rows based on numeric variables in data tables with up to millions of rows. You must specify the number of clusters in advance."

File Edit Tables Rows Cols DOE Analyze Graph Tools View Window Help

europocluster

Columns (3/0)

Country *
Percentage of eligible people earning
Percentage of GDP spent on education

eligible people earning	Percentage of GDP spent on education
7.5	5.99
5.6	4.42
29.2	8.47
7.5	4.6
6.5	5.09
1.9	4.22
6.1	5.81
7.5	4.75
7.1	4.59

16 | Netherland

Distribution

Fit Y by X

Tabulate

Text Explorer

Fit Model

Predictive Modeling

Specialized Modeling

Screening

Multivariate Methods

Clustering

Hierarchical Cluster

K Means Cluster

Normal Mixtures

Latent Class Analysis

Cluster Variables

Clusters rows based on numeric variables in data tables with up to millions of rows. You must specify the number of clusters in advance.

Select Variables (Numeric) To be Used For Clustering

Percentage of eligible people doing lifelong learning	Percentage of GDP spent on education
7.5	5.99
5.6	4.42
29.2	8.47

Clustering - JMP Pro

Finding points that are close, have similar values

Select Columns

3 Columns

Country
Percentage of eligible people doing lifelong learning
Percentage of GDP spent on education

Cast Selected Columns into Roles

Y, Columns Percentage ...ng learning
Percentage ...n education
optional

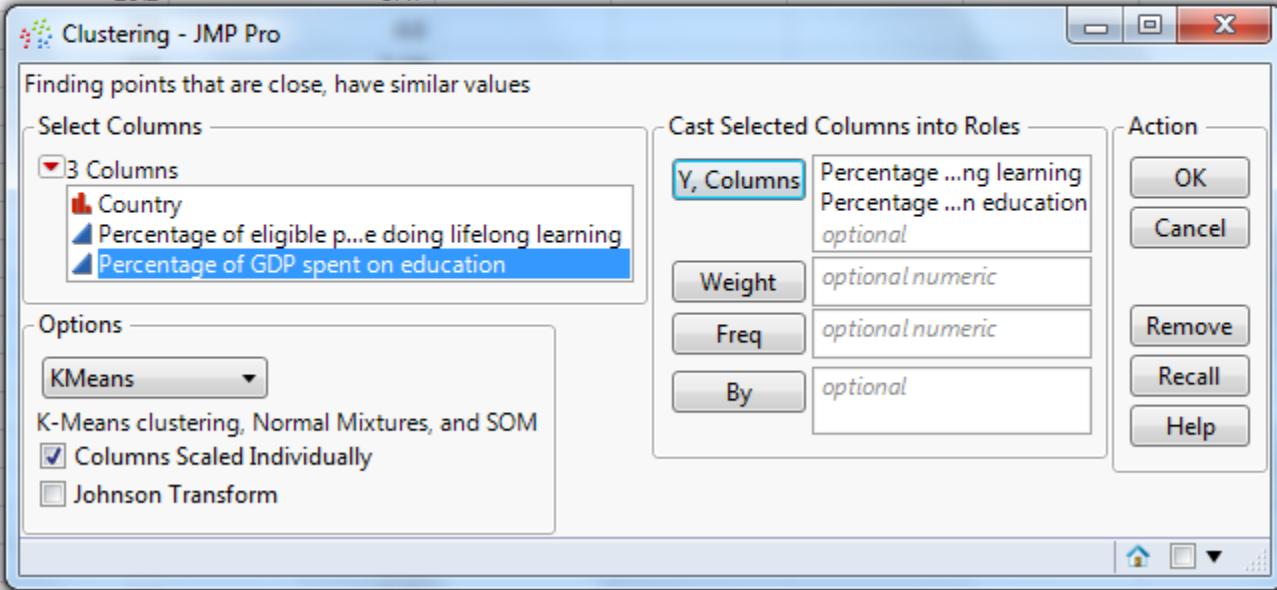
Weight optional numeric

Freq optional numeric

By optional

Action

OK
Cancel
Remove
Recall
Help



Select Number of Clusters (Individual or Range)

	Country	Percentage of eligible people doing lifelong learning	Percentage of GDP spent on education
1	Belgium	7.5	5.99
2	Czech	5.6	4.42
3	Denmark	29.2	8.47
4	Germany		
5	Estonia		
6	Greece		
7	France		
8	Ireland		
9	Italy		
10	Cyprus		
11	Latvia		
12	Lithuania		
13	Luxembourg		
14	Hungary		
15	Malta		
16	Netherlands		
17	Austria		
18	Poland		
19	Portugal		
20	UK		

europecluster - K Means Cluste...

Iterative Clustering
Columns Scaled Individually

Control Panel

Outlier cleanup: Declutter

Method: K-Means Clustering

Number of Clusters: Range of Clusters (Optional)
3 6

Go Help

Single Step
 Use within-cluster std deviations
 Shift distances using sampling rates

Cluster Output

Iterative Clustering

Cluster Comparison

Method	NCluster	CCC	Best
K-Means Clustering	3	-1.6313	
K-Means Clustering	4	-0.8029	
K-Means Clustering	5	-0.4991	
K-Means Clustering	6	0.55953	Optimal CCC

Columns Scaled Individually

Control Panel

K Means NCluster=3

Columns Scaled Individually

Cluster Summary

Cluster	Count
1	17
2	2
3	1

Step

Criterion

Cluster Means

Cluster	Percentage of eligible people doing lifelong learning	Percentage of GDP spent on education
1	64	5.11647059
2	21.1	5.235
3	29.2	8.47

Cluster Standard Deviations

K Means NCluster=4

Columns Scaled Individually

Cluster Summary

Cluster	Count
1	4
2	13
3	1

Step

Criterion

Cluster Means

Cluster	Percentage of eligible people doing lifelong learning	Percentage of GDP spent on education
1	6.0777778	4.63
2	5.46666667	5.77666667
3	29.2	8.47
4	12.3	5.27666667
5	26.6	5.29

K Means NCluster=4

Cluster Summary

Cluster	Count
4	2

Cluster Means

Cluster	Percentage of eligible people doing lifelong learning	Percentage of GDP spent on education
1	8.4	5.99
2	5.78461538	4.84769231
3	29.2	8.47
4	21.1	5.235

Cluster Standard Deviations

K Means NCluster=5

Columns Scaled Individually

Cluster Summary

Cluster	Count
1	9
2	6
3	1
4	3
5	1

Step

Criterion

Cluster Means

Cluster	Percentage of eligible people doing lifelong learning	Percentage of GDP spent on education
1	14.35	5.315
2	6.3	4.41833333
3	29.2	8.47
4	6.9	6.71
5	5.66666667	5.36777778
6	26.6	5.29

Cluster Standard Deviations

K Means NCluster=6

Columns Scaled Individually

Cluster Summary

Cluster	Count
1	2
2	6
3	1
4	1
5	9
6	1

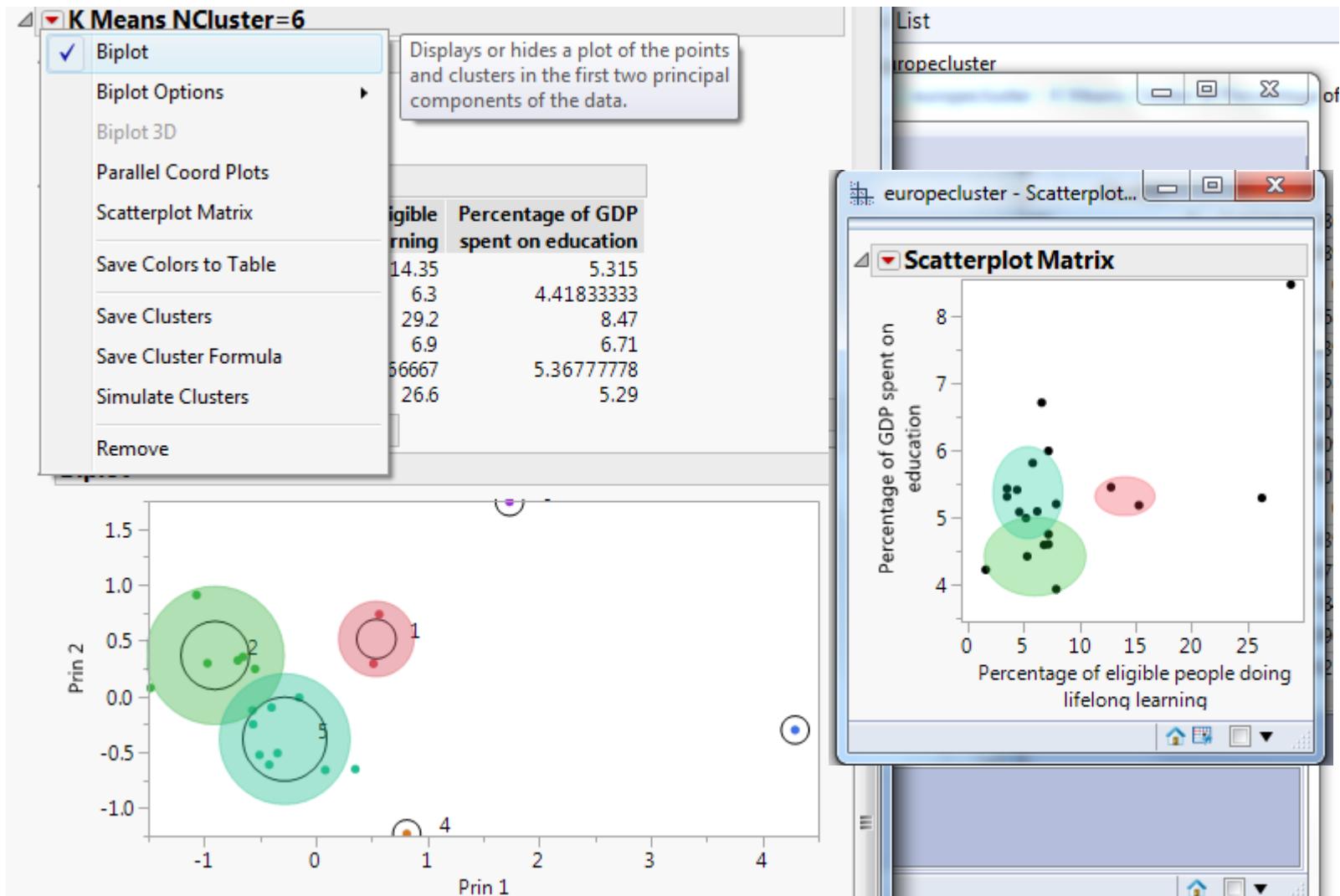
Assign Cluster

The screenshot shows the SPSS interface with several windows open:

- Cluster Means**: A table showing cluster means for two variables:

Cluster	Percentage of eligible people doing lifelong learning	Percentage of GDP spent on education
1	6.0777778	4.63
2	5.46666667	5.77666667
3	29.2	8.47
4	12.3	5.27666667
5	26.6	5.29
- Cluster Standard Deviations**: A table showing cluster standard deviations.
- K Means NCluster=6**: A window showing various options:
 - Biplot
 - Biplot Options
 - Biplot 3D
 - Parallel Coord Plots
 - Scatterplot Matrix
 - Save Colors to Table
 - Save Clusters** (highlighted)
 - Save Cluster Formula
 - Simulate Clusters
 - RemoveA tooltip for "Save Clusters" indicates: "Creates a new column in the data table containing the assigned cluster number."
- Data View**: A table showing individual data points with their assigned cluster numbers (Cluster column) and distances from the centroid (Distance column).

Cluster Visualization



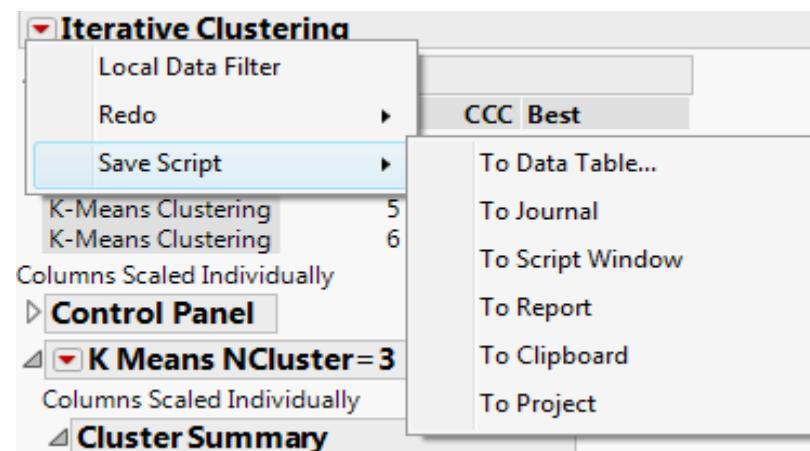
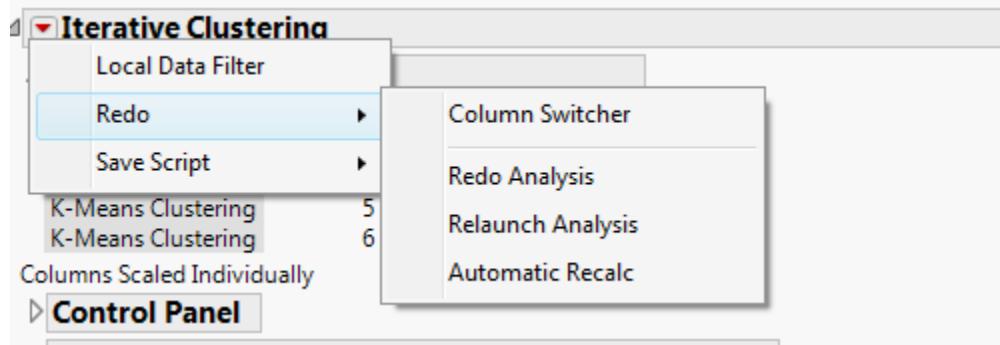
Cluster Script

The screenshot shows the JMP Pro interface with a script window open. The title bar says "Script Window - JMP Pro". The menu bar includes File, Edit, Tables, DOE, Analyze, Graph, Tools, View, Window, and Help. The toolbar has various icons for data manipulation. The main window displays a JSL script for K-Means clustering:

```
K Means Cluster(
  Y(
    :Percentage of eligible people doing lifelong learning,
    :Percentage of GDP spent on education
  ),
  {Number of Clusters( 3 ), Name( "K-Means Clustering" ), Go},
  {Number of Clusters( 4 ), Name( "K-Means Clustering" ), Go},
  {Number of Clusters( 5 ), Name( "K-Means Clustering" ), Go},
  {Number of Clusters( 6 ), Name( "K-Means Clustering" ), Go},
  Go( Parallel Coord Plots, Biplot( 1 ) },
  SendToReport( Dispatch( {}, "Control Panel", OutlineBox, {Close( 1 )} ) )
);
```

A context menu is open over the "Iterative Clustering" section of the script, showing options like Local Data Filter, Redo, Save Script, and CCC Best. The CCC Best value is listed as -1.6313.

Additional Necessary Functionality



Cluster Analysis Using R

Data Set: cluster_children

Launch R Studio

```
R version 3.3.2 (2016-10-31) -- "sincere Pumpkin Patch"
Copyright (C) 2016 The R Foundation for Statistical Computing
Platform: x86_64-w64-mingw32/x64 (64-bit)
```

```
R is free software and comes with ABSOLUTELY NO WARRANTY.
You are welcome to redistribute it under certain conditions.
Type 'license()' or 'licence()' for distribution details.
```

```
R is a collaborative project with many contributors.
Type 'contributors()' for more information and
'citation()' on how to cite R or R packages in publications.
```

```
Type 'demo()' for some demos, 'help()' for on-line help, or
'help.start()' for an HTML browser interface to help.
Type 'q()' to quit R.
```

```
> library(rattle)
Rattle: A free graphical interface for data mining with R.
Version 4.1.0 Copyright (c) 2006-2015 Togaware Pty Ltd.
Type 'rattle()' to shake, rattle, and roll your data.
> |
```

Launch Rattle

The screenshot shows the Rattle graphical user interface. At the top is a menu bar with Project, Tools, Settings, and Help. To the right of the menu is a status bar indicating "Rattle Version 4.1.0 togaware.com". Below the menu is a toolbar with icons for Execute, New, Open, Save, Report, Export, Stop, Quit, and Connect R. A tab bar below the toolbar has tabs for Data, Explore, Test, Transform, Cluster, Associate, Model, Evaluate, and Log, with Data selected. Under the Data tab, there are options for Source (Spreadsheet, ARFF, ODBC, R Dataset, RData File, Library, Corpus, Script), a Filename input field (None), Separator and Decimal settings, and a Header checkbox. There are also fields for Partition (70/15/15), Seed (42), View, and Edit. Below these are Input (green circle), Ignore (red circle), and Weight Calculator fields, along with Target Data Type options (Auto, Categorical, Numeric, Survival). The main window contains a welcome message and several informational paragraphs about Rattle's features and licensing.

Rattle is a free graphical user interface for Data Mining, developed using R. R is a free software environment for statistical computing and graphics. Together they provide a sophisticated environments for data mining, statistical analyses, and data visualisation.

See the Help menu for extensive support in using Rattle. The book Data Mining with Rattle and R is available from Amazon. The Togaware Desktop Data Mining Survival Guide includes Rattle documentation and is available from datamining.togaware.com

Rattle is licensed under the GNU General Public License, Version 2. Rattle comes with ABSOLUTELY NO WARRANTY. See Help -> About for details.

Rattle Version 4.1.0. Copyright 2006-2015 Togaware Pty Ltd.
Rattle is a registered trademark of Togaware Pty Ltd.
Rattle was created and implemented by Graham Williams.

Select Appropriate Data Set & Variable & Execute

The screenshot shows the R Data Miner application window. The title bar reads "R Data Miner - [Rattle (cluster_children.csv)]". The menu bar includes Project, Tools, Settings, Help, and a Rattle Version 4.1.0 link. The toolbar has icons for Execute, New, Open, Save, Report, Export, Stop, Quit, and Connect R. Below the toolbar is a tab bar with Data, Explore, Test, Transform, Cluster, Associate, Model, Evaluate, and Log, with Data selected. A Source dropdown shows Spreadsheet selected. The filename is cluster_children.... with separator comma and decimal dot. A Header checkbox is checked. Partition is set to 70/15/15, Seed to 42, with View and Edit buttons. The Target Data Type is set to Auto. The main area displays a table of variables:

No.	Variable	Data Type	Input	Target	Risk	Ident	Ignore	Weight	Comment
1	X.U.FEFF.Name	Categoric	<input checked="" type="radio"/>	Unique: 20					
2	age	Numeric	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Unique: 13
3	mem_span	Numeric	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Unique: 11
4	iq	Numeric	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Unique: 13
5	read_ab	Numeric	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Unique: 11

Data File: cluster_children

K Means Cluster – Rattle Output

The screenshot shows the Rattle software interface for data mining. The top menu bar includes Execute, New, Open, Save, Report, Export, Stop, Quit, and Connect R. Below the menu is a tab bar with Data, Explore, Test, Transform, Cluster, Associate, Model, Evaluate, and Log. The Cluster tab is selected. A Type dropdown shows KMeans selected. Below it, Number of clusters is set to 3, Seed to 42, and Runs to 1. A checkbox for Re-Scale is checked. Other options like Use HClust Centers and Iterate Clusters are unchecked. Buttons for Stats, Plots, Data, Discriminant, and Weights are present.

Cluster sizes:

```
[1] "9 2 3"
```

Data means:

```
age mem_span iq read_ab
0.6850649 0.5413534 0.4880952 0.6000000
```

Cluster centers:

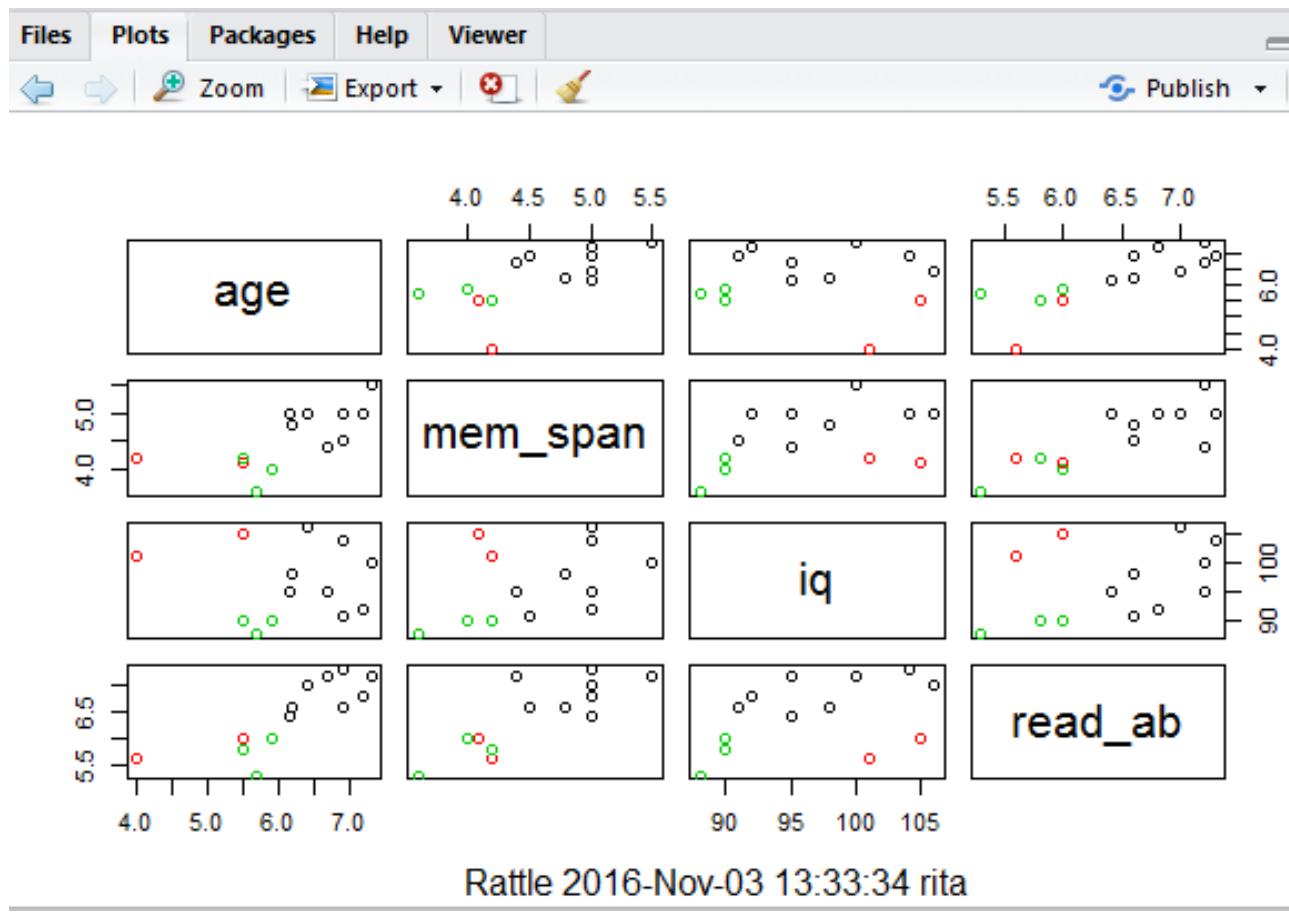
```
age mem_span iq read_ab
1 0.8434343 0.7192982 0.54938272 0.8111111
2 0.2272727 0.2894737 0.83333333 0.2500000
3 0.5151515 0.1754386 0.07407407 0.2000000
```

Within cluster sum of squares:

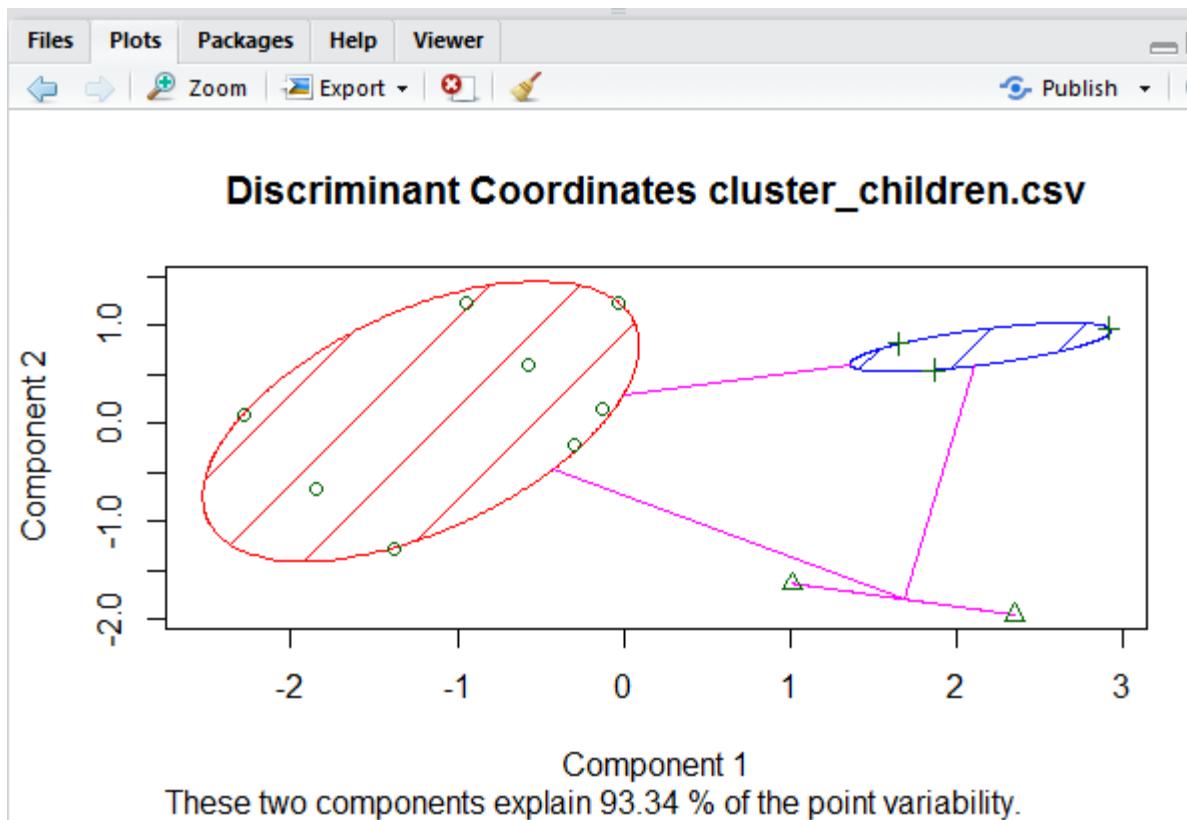
```
[1] 1.3352486 0.1493822 0.1322849
```

Time taken: 0.00 secs

K Means Cluster – Plots



K Means Cluster – Plots



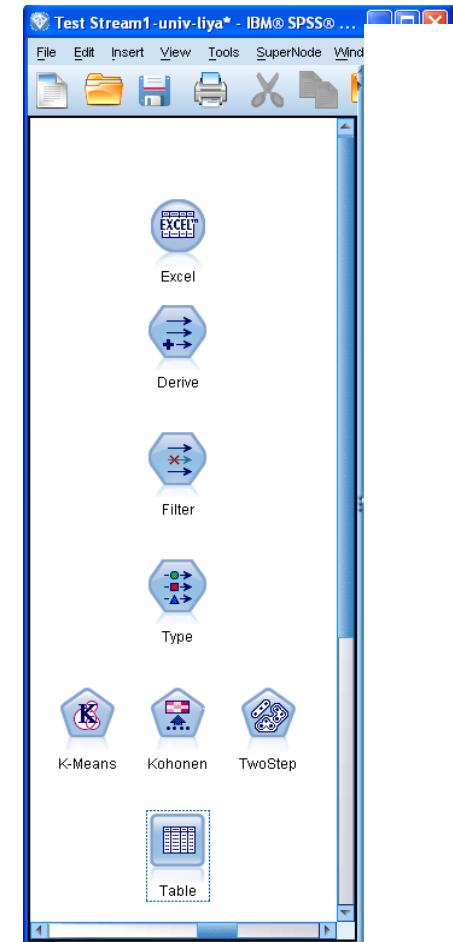
Cluster Analysis Using SPSS Modeler

Data Set: UScolleges

Clustering using SPSS Modeler

Nodes involved have following operations:

- **Source:** Access dataset e.g. Excel
- **Derive:** Generate field based on existing ones
- **Filter:** Rename or remove field for next step
- **Type:** Specify the type (continuous/nominal/ ordinal) of fields
- **Modeling** – Clustering algorithms : K-means, Kohonen, TwoStep etc.
- **Table:** Final result capture & validation



Data File: US Colleges

- Dataset: American College and University
 - 1302 records with 23 fields:
 - College name, State, Public/private,
 - Math SAT, Verbal SAT, Tuition fee,
 - No. of applicants received,
 - No. of applicants accepted, ...

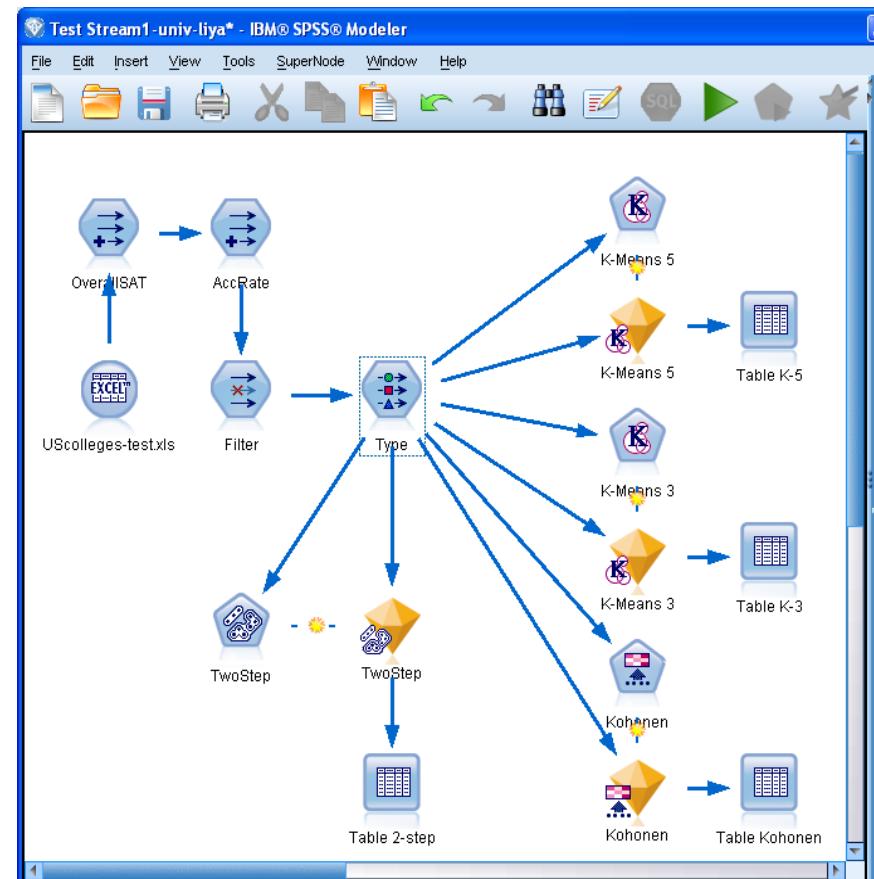
(source: G. Shmueli , N. R. Patel, P. C. Bruce, "Data Mining for Business Intelligence", Wiley, 2010)



- Application
 - An education agent to recommend US universities to applicants (outside of US) given his/her SAT scores and financial capability
 - The company needs an optimal overall plan for resource allocation to cope with potential demands to different groups of colleges, such as private, public, top school, ...

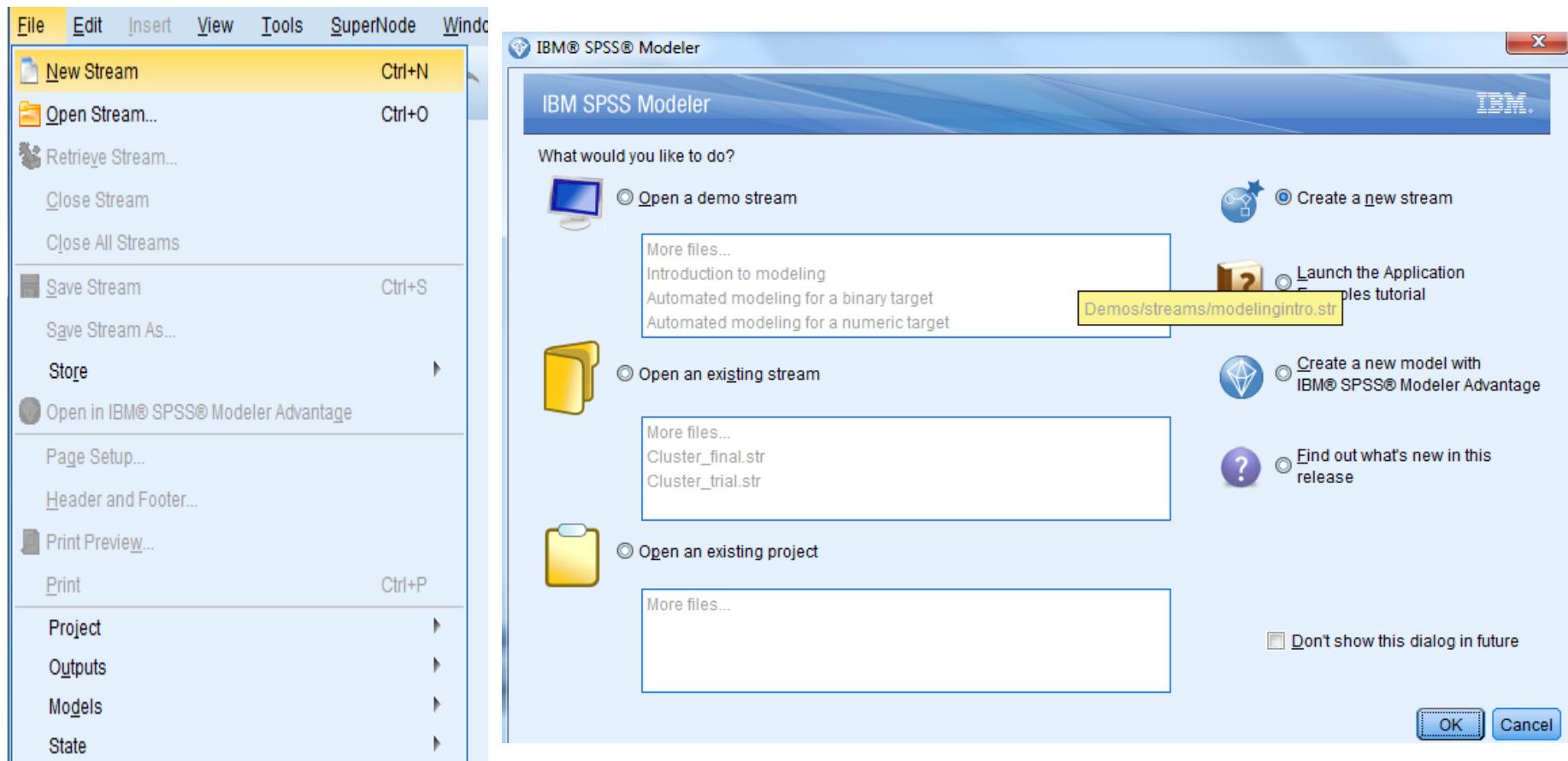
Case Study: US Colleges (cont.)

Demo Stream :
**One needs to
create a stream
to do any
analysis**



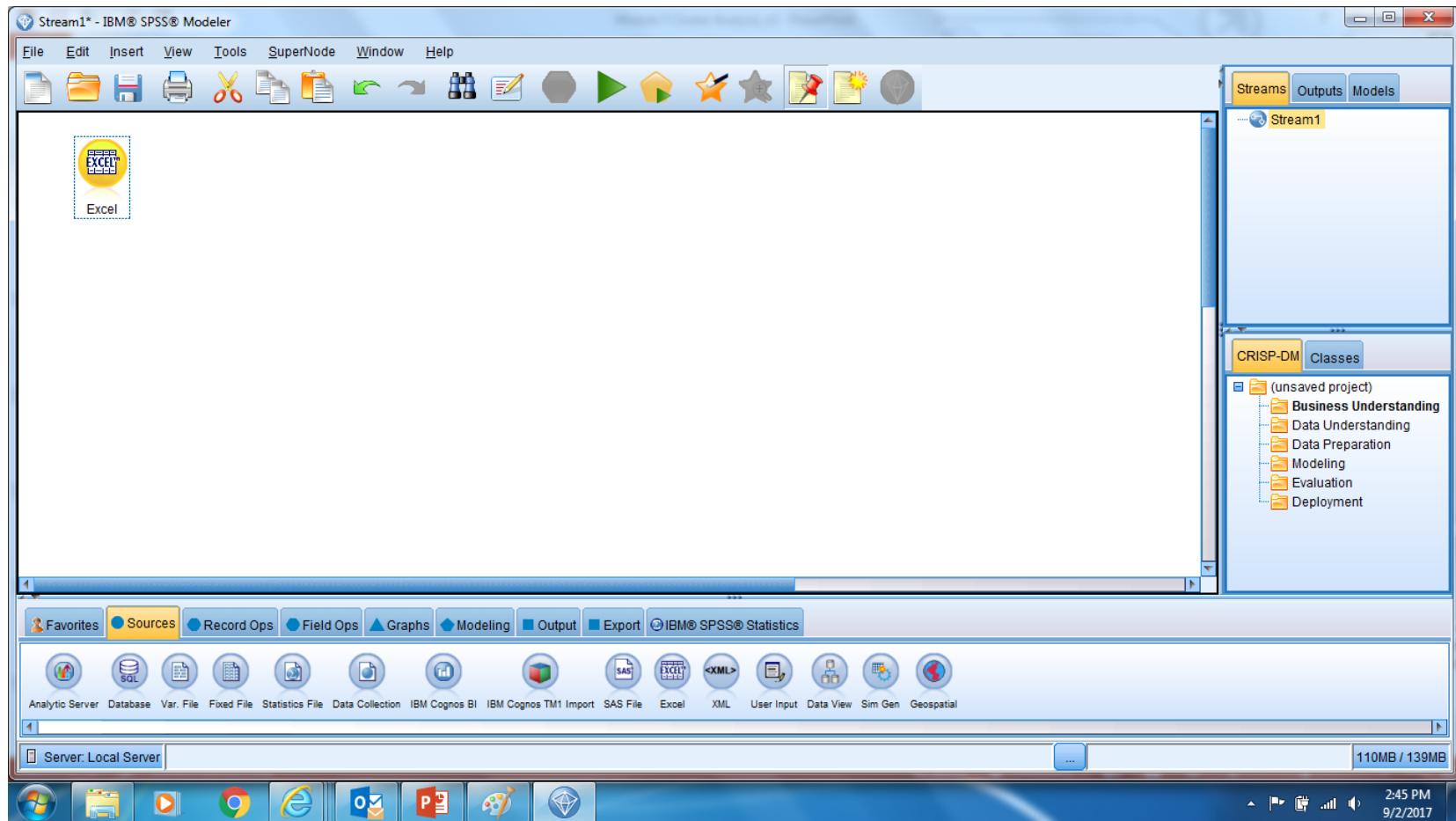
Creation of Stream

Launch SPSS Modeler – start a new stream



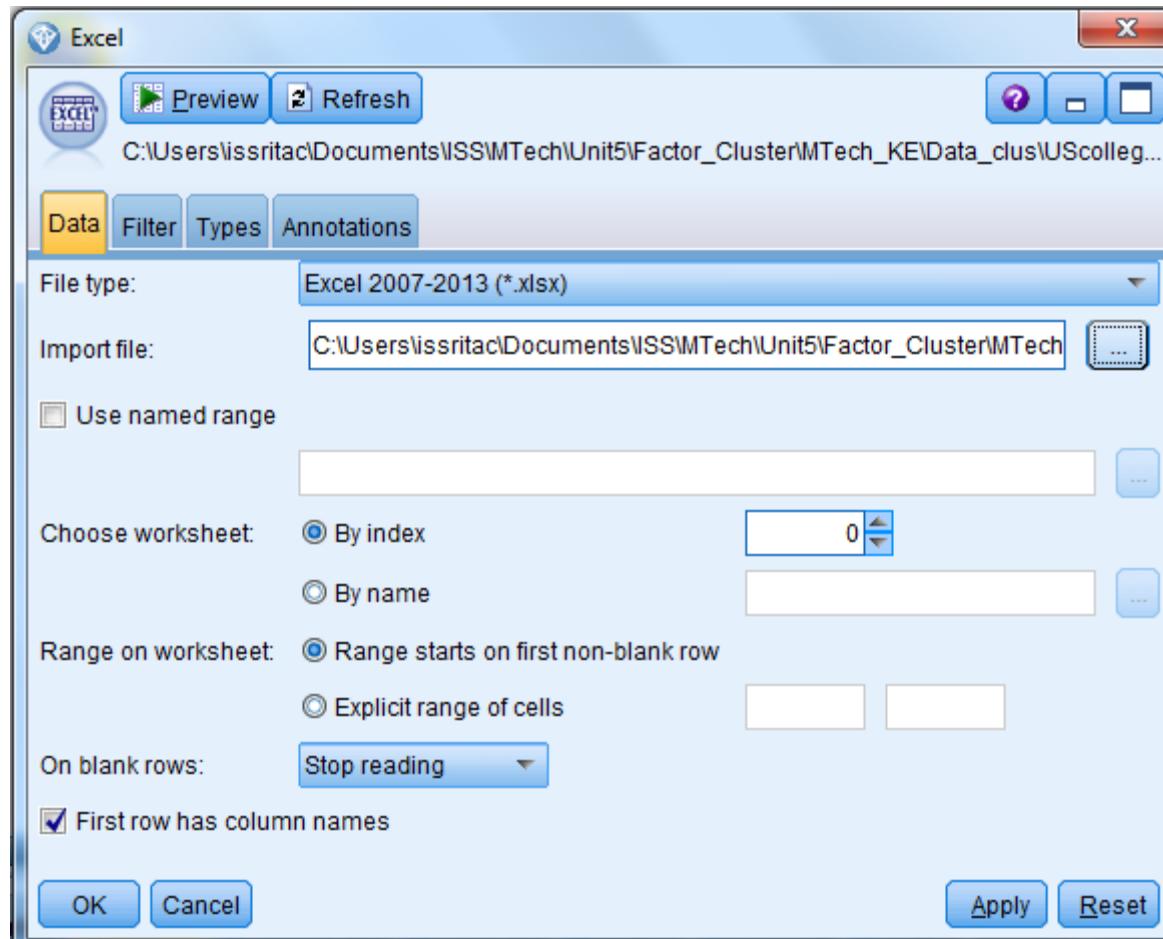
Creation of Stream

Drag & drop necessary icons on the Stream screen and link them up



Creation of Stream

Import file and start various operations



US Colleges: Data Preparation

New variable creation

The screenshot shows the KNIME Data Preparation interface. On the left, there is a file icon labeled "UScolleges_1.xlsx" and a "Derive" node icon. A blue arrow points from the file icon to the "Derive" node. The main window displays the "Derive" node configuration and an open "Expression Builder" dialog.

Derive Node Configuration:

- Derive field:** OverallSAT
- Derive as:** Formula
- Field type:** Continuous
- Formula:** 1

Expression Builder Dialog:

The formula entered is: 'Math SAT' + 'Verbal SAT'

General Functions:

Function	Return
is_integer(ITEM)	Boolean
is_real(ITEM)	Boolean
is_number(ITEM)	Boolean
is_string(ITEM)	Boolean
is_date(ITEM)	Boolean
is_time(ITEM)	Boolean
is_timestamp(ITEM)	Boolean
is_datetime(ITEM)	Boolean
to_integer(ITEM)	Integer
to_real(ITEM)	Real

Fields:

Type	Field	Storage
String	College Name	String
String	State	String
Real	Public (1)/ Priv...	Real
Real	Math SAT	Real
Real	Verbal SAT	Real
Real	ACT	Real
Real	# appl. rec'd	Real
Real	# appl. accept...	Real
Real	# new stud. en...	Real
Real	% new stud. fr...	Real

Check Expression Before Saving:

Buttons: OK, Cancel, Check, Help

US Colleges: Data Preparation

Data source: Uscolleges_1.xls

Derived fields (using derive nodes)

- Overall SAT = Math SAT + Verbal SAT
- AccRate = appl. accepted / appl. rec'd

The image shows two side-by-side dialog boxes from a data preparation tool. Both boxes have a blue header bar with icons for Preview, Help, Minimize, and Close. The left box is titled 'OverallSAT' and the right box is titled 'AccRate'. Both boxes have tabs for 'Settings' (which is selected) and 'Annotations'. Under 'Settings', there is a 'Mode' section with radio buttons for 'Single' and 'Multiple', and a 'Single' button is selected. Below that is a 'Derive field:' input field containing 'OverallSAT'. Underneath it is a 'Derive as:' dropdown set to 'Formula', a 'Field type:' dropdown set to 'Continuous', and a 'Formula:' input field containing "'Math SAT' + 'Verbal SAT'" with a small icon to its right. At the bottom are four buttons: 'OK', 'Cancel', 'Apply', and 'Reset'. The right box has a similar layout but with 'AccRate' in the title and 'AccRate' in the 'Derive field:' field, with the formula '# appl. accepted' / '# appl. rec'd' in the 'Formula:' field.

US Colleges: Data Preparation (cont.)

Fields selected for clustering (using filter node)

- Public/private
- % new stud. from top 10%
- % new stud. from top 25%
- In-state tuition
- Stud./fac. ratio
- **OverallSAT**
- **AccRate**

Not for clustering,
but keep
for result
evaluation

The screenshot shows the KNIME 'Filter' node configuration window. The 'Filter' tab is selected. A red box highlights the first two rows: 'College Name' and 'State'. A yellow box highlights the last three rows: 'Graduation rate', 'OverallSAT', and 'AccRate'. A red arrow points from the text box above to the 'Graduation rate' row in the table.

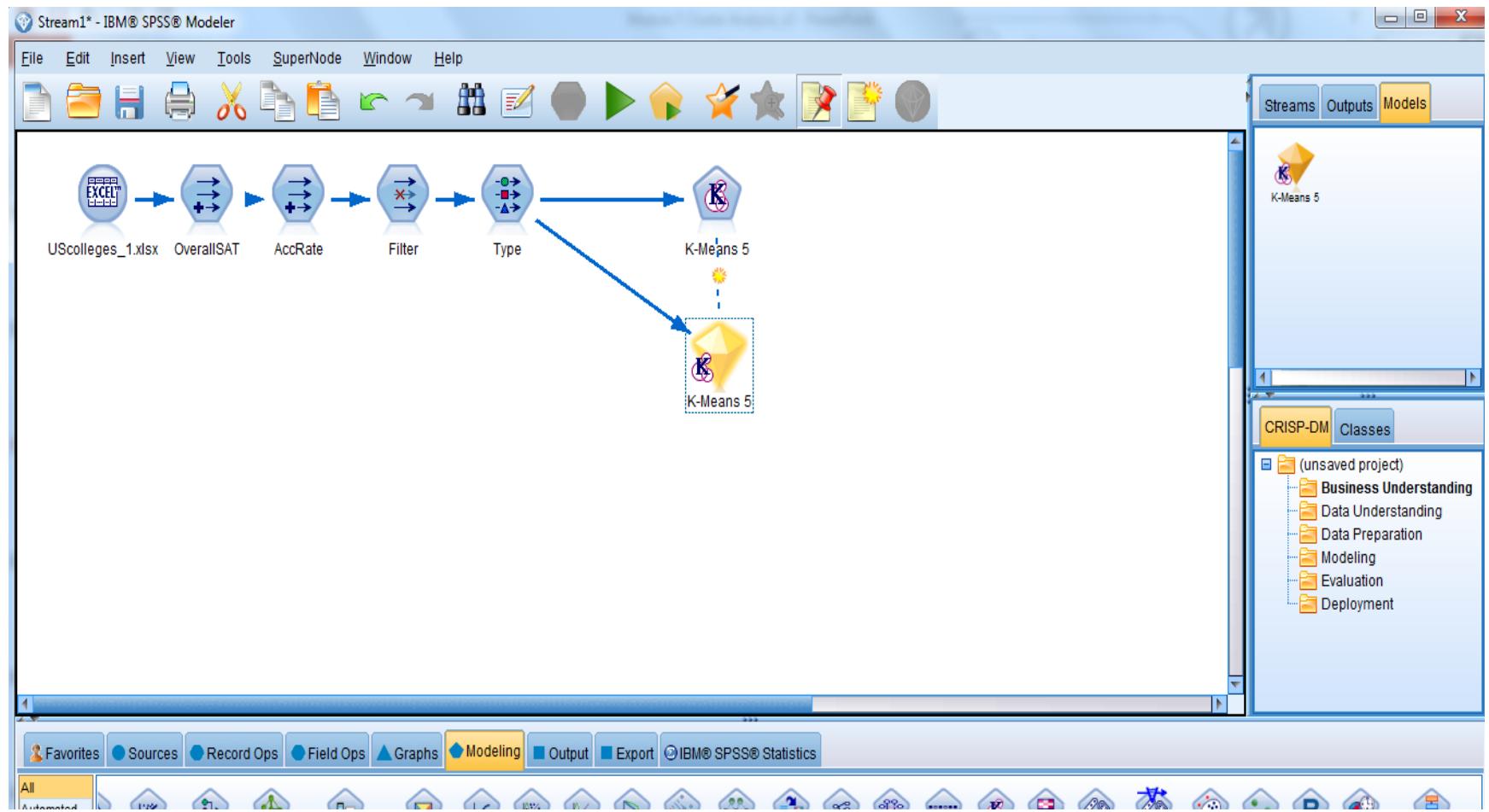
Field	Filter	Field
College Name	→	College Name
State	→	State
Public (1)/ Private (2)	→	Public (1)/ Private (2)
Math SAT	✗ →	Math SAT
Verbal SAT	✗ →	Verbal SAT
ACT	✗ →	ACT
# appli. rec'd	✗ →	# appli. rec'd
# appli. accepted	✗ →	# appli. accepted
# new stud. enrolled	✗ →	# new stud. enrolled
% new stud. from top 10%	→	% new stud. from top 10%
% new stud. from top 25%	→	% new stud. from top 25%
# FT undergrad	✗ →	# FT undergrad
# PT undergrad	✗ →	# PT undergrad
in-state tuition	→	in-state tuition
out-of-state tuition	✗ →	out-of-state tuition
room	✗ →	room
board	✗ →	board
add. fees	✗ →	add. fees
estim. book costs	✗ →	estim. book costs
estim. personal \$	✗ →	estim. personal \$
% fac. w/PHD	✗ →	% fac. w/PHD
stud./fac. ratio	→	stud./fac. ratio
Graduation rate	✗ →	Graduation rate
OverallSAT	→	OverallSAT
AccRate	→	AccRate

Fields: 25 in, 16 filtered, 0 renamed, 9 out

View current fields View unused field settings

US Colleges: Clustering Method Selection

Select Modeling Method

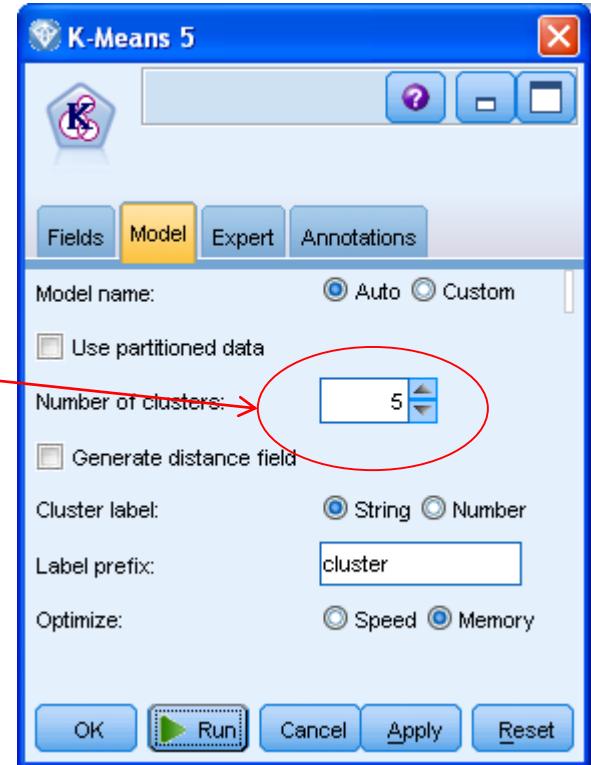


US Colleges: k -Means

For demo purpose, we have two k -Means nodes

- $k = 5$
- $k = 3$

We first run with $k = 5$



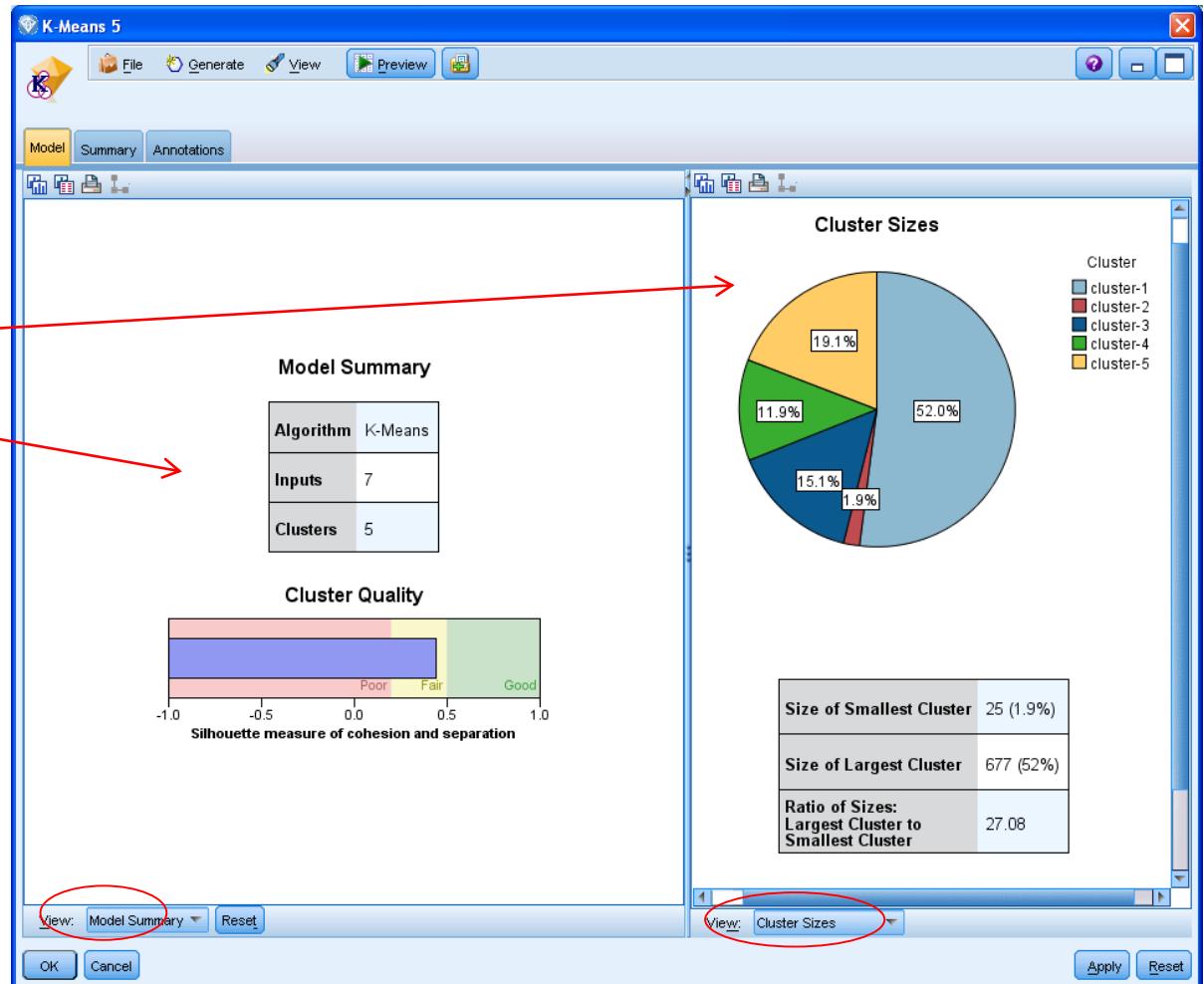
US Colleges: k -Means ($k = 5$)

Double click the nugget

- “K-means 5”

Cluster size

Model summary



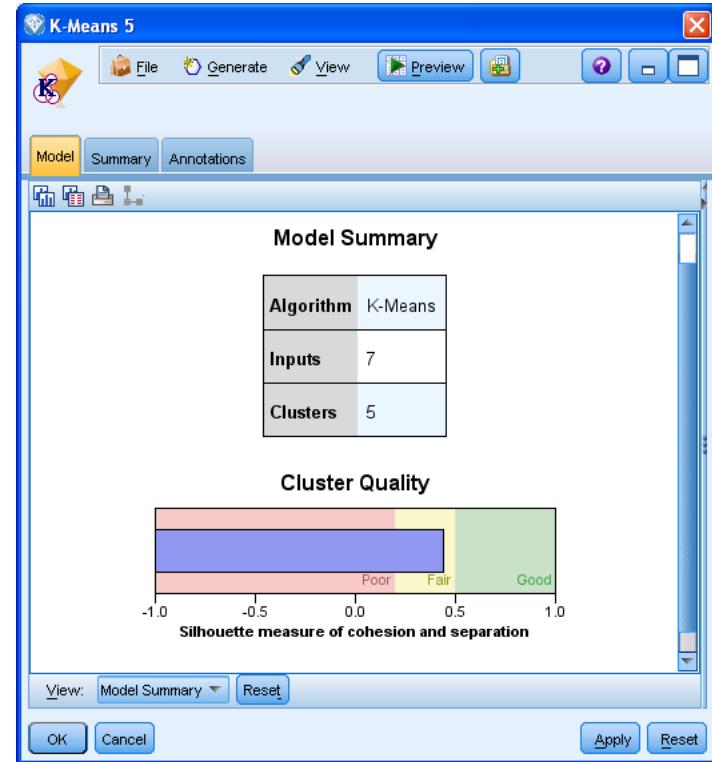
Cluster Quality

Silhouette coefficient is a measure of both cohesion and separation. For each element in a cluster, calculate :

- the average distance to all other elements in its cluster, d_{ave} &
- the average distance to all elements in each of the other clusters, D_{ave}

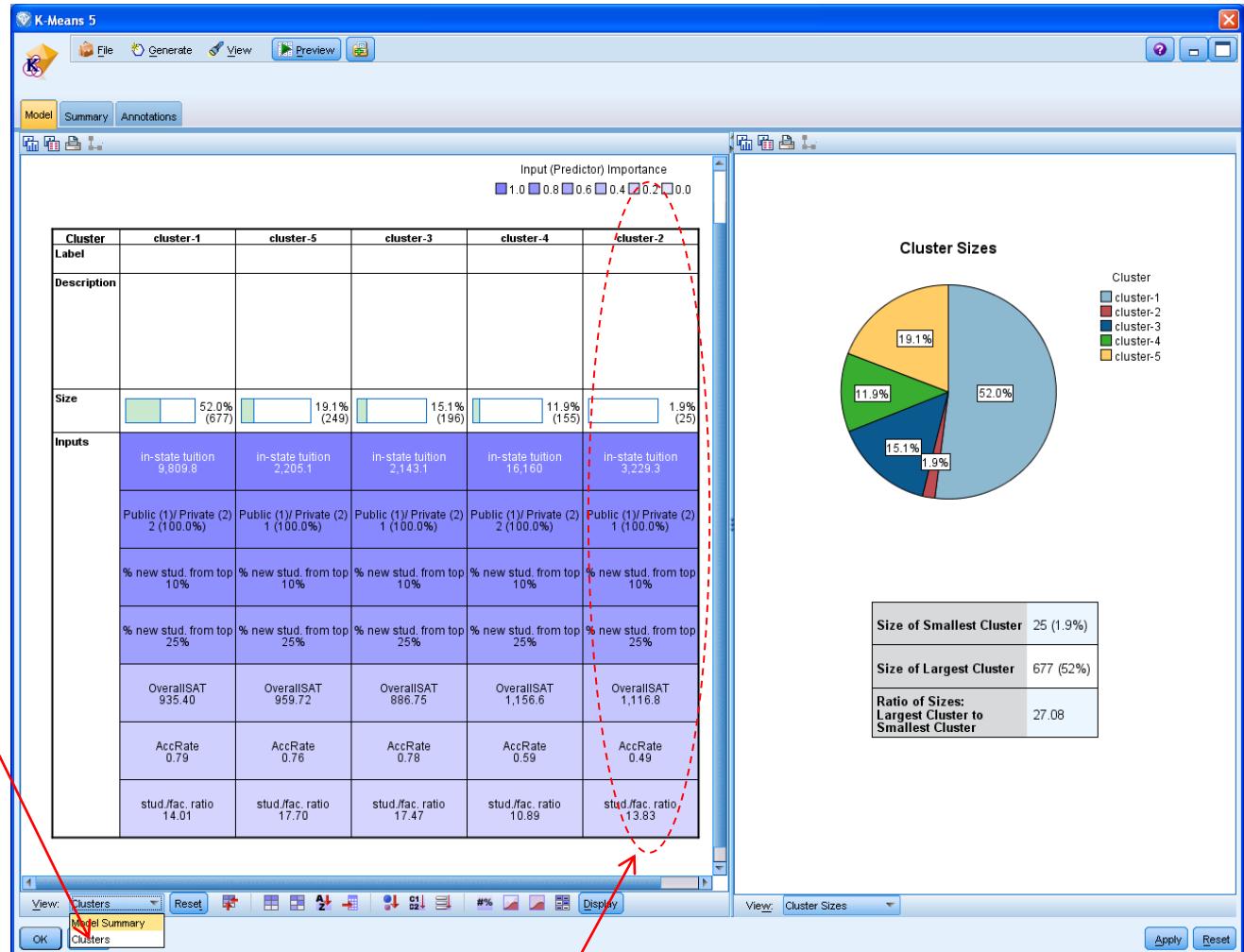
Silhouette =

$$\cdot (D_{ave} - d_{ave}) / \text{Max}(D_{ave}, d_{ave})$$



US Colleges: k -Means ($k = 5$) (cont.)

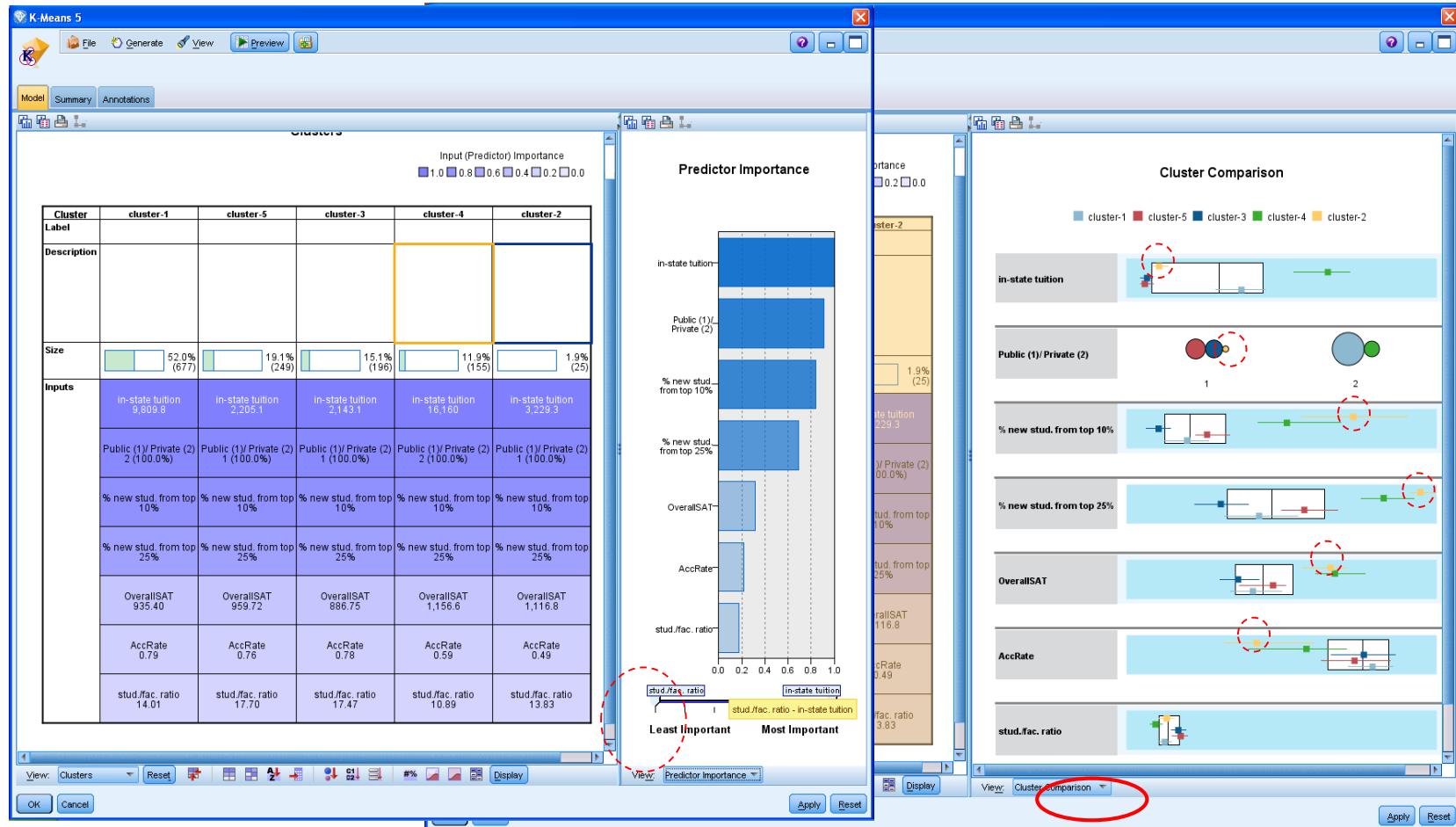
Choose
“cluster”



The smallest cluster

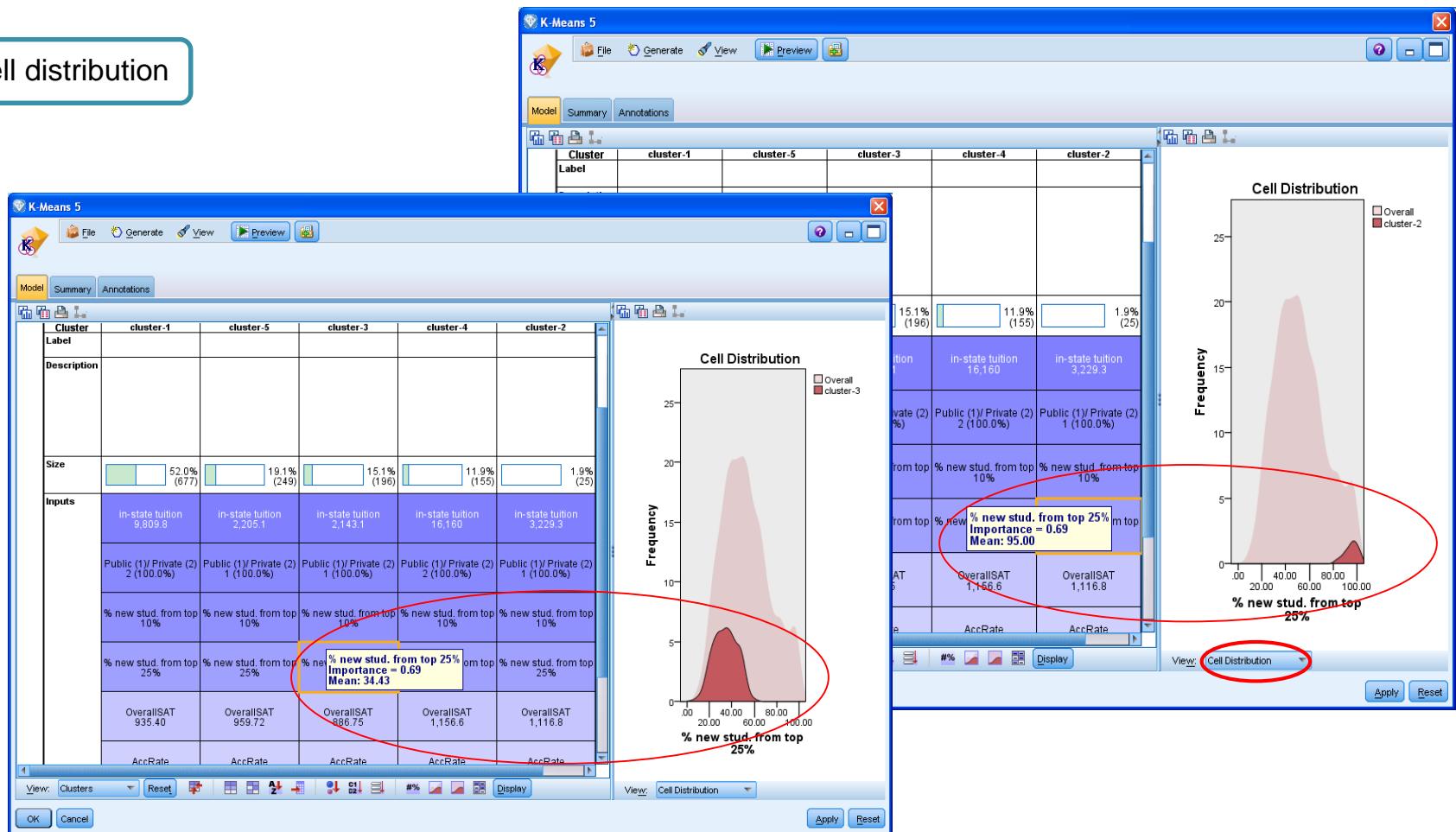
US Colleges: k -Means ($k = 5$) (cont.)

Cluster comparison (cluster-2, the smallest cluster)



US Colleges: k -Means ($k = 5$) (cont.)

Cell distribution



US Colleges: k -Means ($k = 5$) (cont.)

Run the Table K-5

- In the cluster-2
 - We found some famous top public universities

Table K-5 (10 fields, 1,302 records)

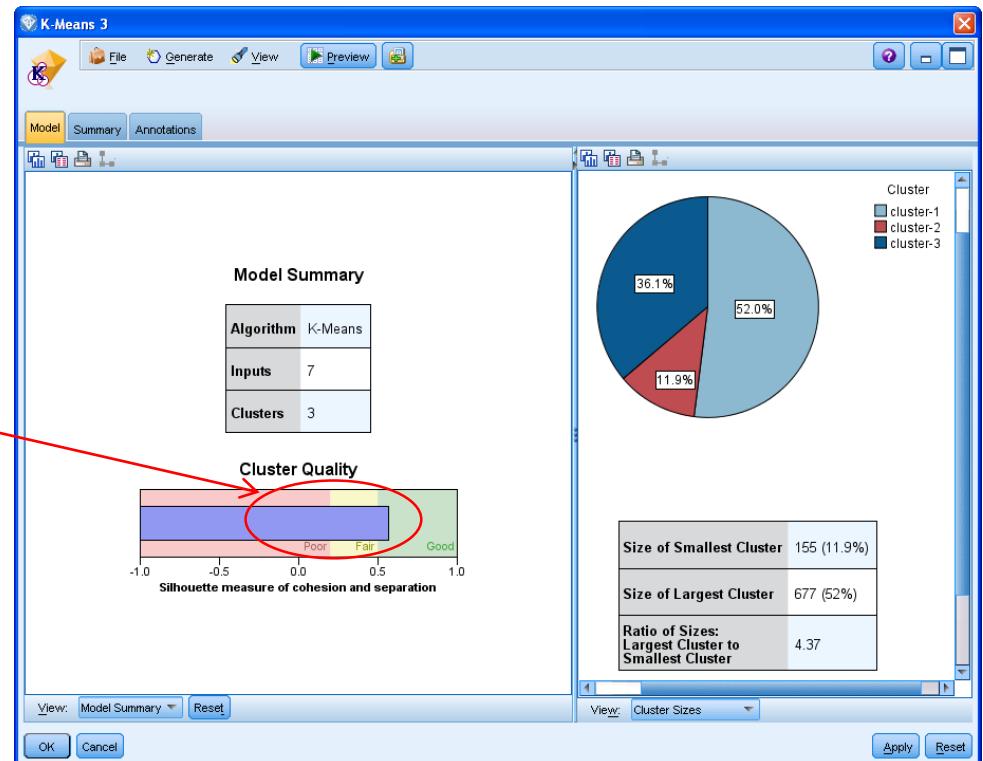
	College Name	State	Public (1)/ Private (2)	% new s...	% new stud. from top 25%	... stud./fac. ratio	OverallSAT	AccR...	\$KM-K-Means
98	Stanford University	CA	2.000	88.000	97.000	... 9.700	1401.000	0.215	cluster-4
99	University of California at Berkeley	CA	1.000	95.000	100.000	... 15.800	1218.000	0.415	cluster-2
100	University of California at Davis	CA	1.000	95.000	100.000	... 14.500	1070.000	0.697	cluster-2
101	University of California at Irvine	CA	1.000	85.000	100.000	... 16.100	1029.000	0.686	cluster-2
102	University of California at Los Angeles	CA	1.000	93.000	100.000	... 11.800	1141.000	0.471	cluster-2
103	University of California at Riverside	CA	1.000	80.000	100.000	... 13.400	982.000	0.772	cluster-2
104	University of California at San Diego	CA	1.000	95.000	100.000	... 13.900	1131.000	0.590	cluster-2
105	University of California at Santa Barbara	CA	1.000	90.000	100.000	... 18.900	1006.000	0.849	cluster-2
106	University of California at Santa Cruz	CA	1.000	94.000	100.000	... 21.300	1059.000	0.796	cluster-2
107	University of Redlands	CA	2.000	26.000	63.000	... 8.900	\$null\$	0.814	cluster-1
108	University of San Francisco	CA	2.000	23.000	48.000	... 13.600	977.000	0.746	cluster-1
109	Santa Clara University	CA	2.000	40.000	72.000	... 63.0	12.000	\$null\$	cluster-4

US Colleges: k -Means ($k = 3$)

Now run k -means with $k = 3$

Double click the nugget “K-means 3”

- The Silhouette measure is higher than that of $k = 5$
- Is the clustering more meaningful than that of $k = 5$?



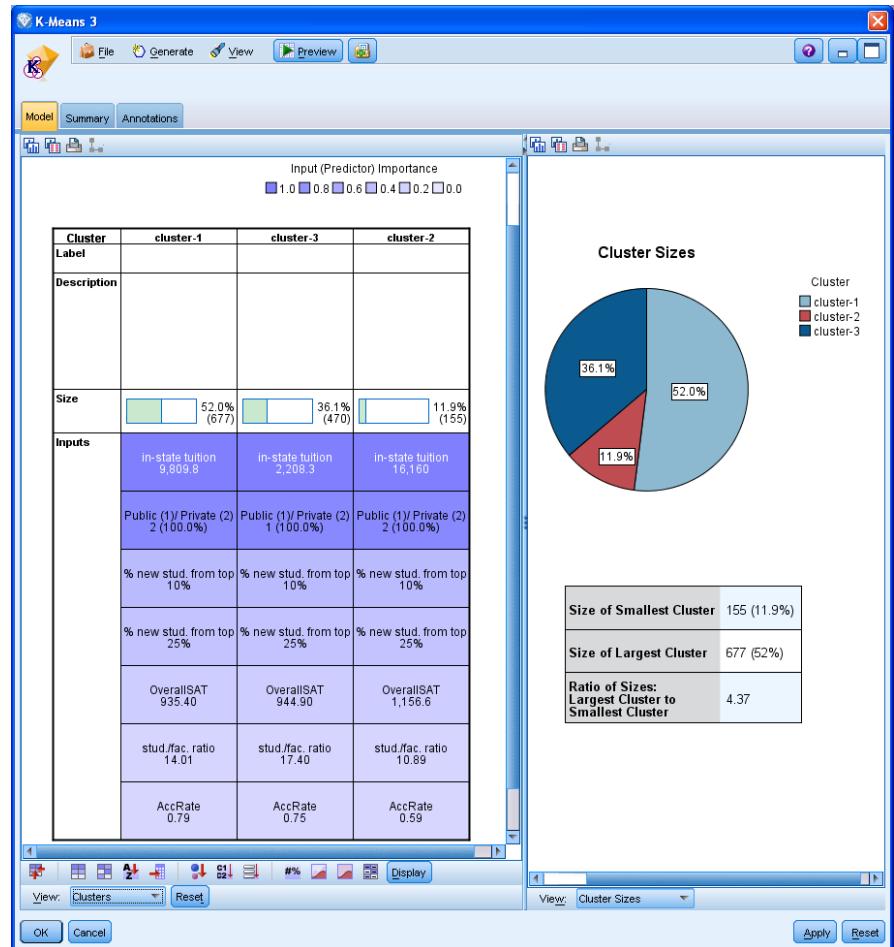
US Colleges: k -Means ($k = 3$) (cont.)

Result analysis

- Cluster 1: private schools
 - lowest overall SAT & highest acceptance rate
- Cluster 2: private schools
 - highest overall SAT & lowest acceptance rate
- Cluster 3:
 - includes all public schools

Compare with $k = 5$

- Three previous clusters of
 - public schools are now in
 - one cluster
- **Which one is more useful?**



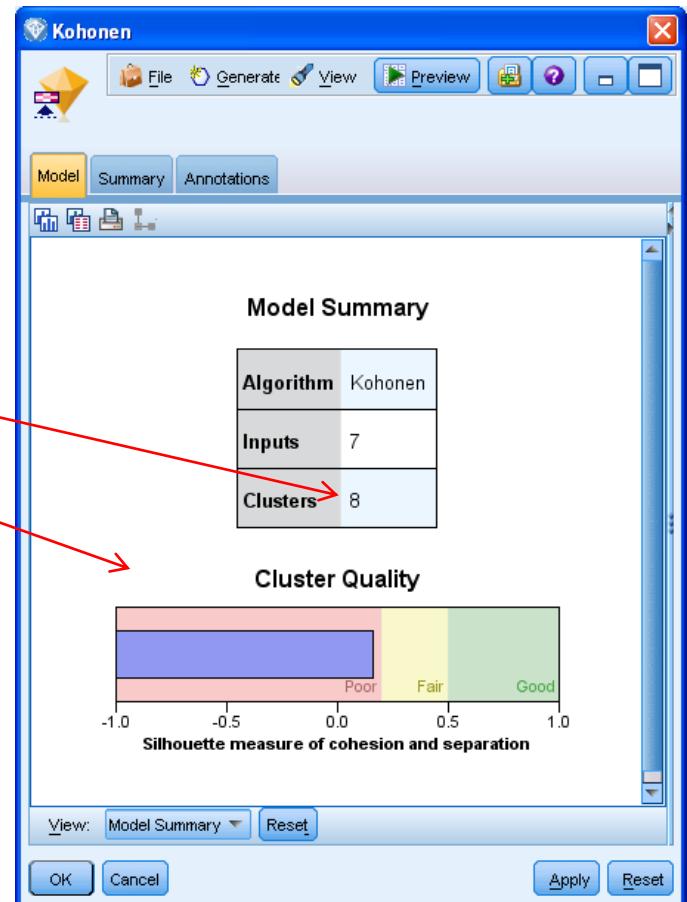
US Colleges: Kohonen clustering

Run Kohonen clustering

Open the nugget “Kohonen”

- 8 clusters formed
- The cluster quality seems poorer than that of *k*-means ?

Further analysis is needed



TwoStep Clustering Mechanism

Step 1:

- Pre-clustering to make little clusters

Step 2:

- Clustering of pre-clusters

TwoStep clustering algorithm offers function to exclude outliers

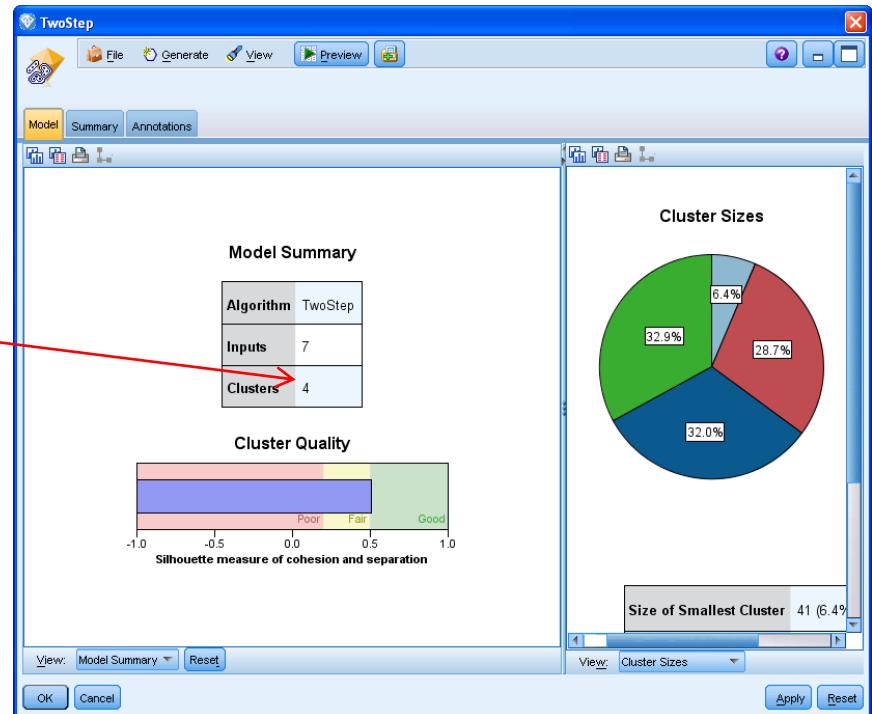
US Colleges: TwoStep clustering

Run TwoStep clustering

Open the nugget “TwoStep”

- 4 clusters formed(using Euclidean distance)

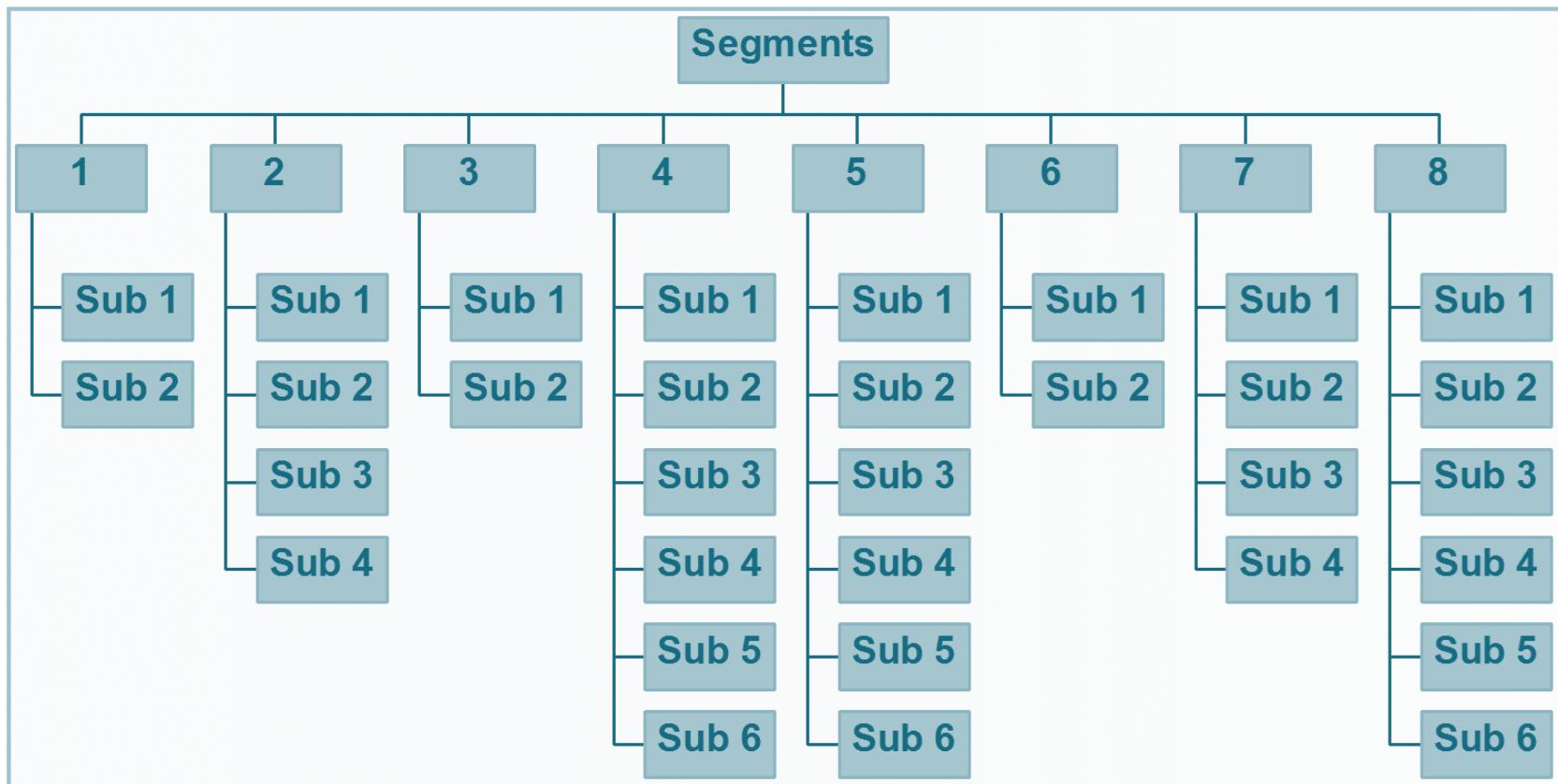
Further analysis is needed



REAL LIFE EXAMPLE OF TWO STEP CLUSTERING

First Level Segmentation Variables : Revenue, Utilization, Credit Risk, Attrition Risk

Second Level Segmentation Variables : Please see slide 96



REAL LIFE EXAMPLE OF TWO STEP CLUSTERING

Note how categorical variables are treated

Variable	Definition
Gender	Female vs. others.
Marital status	Single vs. others.
Age as of latest month	Older (>45 years) vs. others.
Home ownership	Buying or owning vs. others.
Occupation: “White-collar, educated jobs” vs. others.	White-collar educated jobs include: Snr Govt Officer, Bus Mngr./Exec., Cler/Sales Mngr., Health Prof., Build Prof, Bus. Prof., Social Prof., Tech/Skill, Misc Prof, Nurse/Med Tech, Build Tech, Para-Prof.
Spend	Monthly average of customer sales amount in the last 12 months.
Overseas transactions as % of total transactions. The following spending categories are used:	Proportion of average transactions in a particular merchant type as % of average number of total transactions in the last 12 months.
Travel	Airlines, auto rental, hotel, travel services, and other transportation.
Restaurant & Entertainment	Restaurant, bar, amusement entertainment
Shopping	Department store and clothing/apparel stores. retail stores
Day to Day	Telecommunications, utilities, gas station, grocery, insurance payments.
Cash	Cash advance from manual disbursement and ATM.
Home	Home improvement, electronics appliances, furniture

Validating Clusters

As described the method in slide 12, in order to apply clustering result for decision, validating is necessary.

It may be reviewed from

- Cluster interpretability (through human examination)
 - Is the interpretation of the resulting clusters reasonable?
 - E.g.: in the US College case study, k-means with k=5 and k=3
- Cluster stability (through experimental test)
 - Partition the data and see how well clusters that are formed based on one part apply to the other part
- Cluster separation (through statistical test)
 - Statistical tests to examine the ratio of between-cluster variation to within-cluster variation (but their usefulness is somewhat controversial)

Data Sets For Experimentation



Manufacturing Industry Example

http://www.jmp.com/en_us/events/ondemand/mastering-jmp/pinpointing-and-reducing-defects.html

Thank You!