# **Kmean Cluster Analysis**

1

#### **Learning Objectives**

- Understanding the kmean cluster analysis procedure.
- Understanding the methods used to determine the optimal number of clusters.
- Managing data for the sake of conducting cluster analysis.
- Conducing kmean cluster analysis using R
- Understanding the concept of dietary patterns

#### **Learning Objectives**

- Connecting cluster analysis results with other features of individuals
- Learning to conduct cross-tabulation analysis in R.

#### **Road Map**

- An introduction to cluster analysis and kmean cluster analysis.
- A simple example of kmean clustering.
- Issues to consider in conducting cluster analysis.
- Kmean cluster analysis in practice: the case of dietary patterns
  - Dataset and data management
  - Optimal number of clusters
  - Identifying the clusters
  - Means, frequencies and cross-tabulation

#### Machine Learning (ML)

- ML refers to methods and algorithms looking for patterns in a dataset by learning from the data itself.
- The machine, learns from data by conducting the same tasks several times until repeating the task does not improve a predefined criteria.
  - (mean squared error in linear regression or percentage of correct predictions in logistic regression)

#### Machine Learning (ML)

- There are two types of ML methods
  - Supervised ML: where the researcher defines features (variables) of the model (e.g. random forest and support vector machine)
  - Unsupervised ML: the researcher lets an algorithm to look for specific pattern(s) without determining what variables could possibly determine the pattern (e.g. cluster analysis and principle component analysis)

- CL refers to a series of methods aimed at finding the NATURAL GROUPS (CLUSTERS) in a dataset.
- There are two types of clustering methods
  - Hierarchical: refers to methods used for natural grouping in datasets that are in a top-bottom order (e.g. folders and files in your computer)
    - Hierarchical clustering is time consuming and proper for small datasets.

- There are two types of clustering methods
  - Hierarchical: refers to methods used to natural grouping in dataset ordered hierarchically (folders and files in your computer are ordered hierarchically)
  - Partitioning clustering: refers to the methods group the data into clusters that are not overlapping (kmean, kmedian)

- There are two types of clustering methods
  - Hierarchical: refers to methods used to natural grouping in dataset ordered hierarchically (folders and files in your computer are ordered hierarchically)
  - Partitioning clustering: refers to the methods group the data into clusters that are not overlapping (kmean, kmedian)
    - These methods can be used for large datasets and large sets of variables

- Among the methods, kmean clustering is highly popular.
- Kmean is employed in several subjects such as biology, physics marketing and nutrition studies.
- The popularity of kmean method is due to its ability in finding the patterns in data.
- For instance in **marketing** kmean CL can be used to find the shopping or expenditure patterns. In **nutrition** kmean clustering can be used to find food consumption patterns.

- Lets assume we have a dataset including the expenditures of 22 households on two different types of books: fiction books and kids' books.
- We would like to know if we can distinguish between the households based on their patterns of expenditures on these two types of books.
- We use kmean CL to find the clusters.

- Kmean CL find the natural groupings based on an iterative process.
- We have to tell the kmean clustering what are the variables that it should explores and how many groups we think exist in the dataset.
- For our dataset we tell kmean there are two variables: expenditures on fiction books and expenditures on kids' books.
- We also tell kmean that we think there are three groups of households based on their expenditures on these books.

• First: kmean choose 3 random values in the data set (blue diamonds)





- Second: kmean makes three groups of observations based on their distance to the randomly assigned values (blue diamonds).
  - So the closer data points to each random value, will be grouped into one cluster (inside the curves).



- Third: the mean part of kmean CL kicks in. So, the mean values of data points in each group are calculated (yellow diamonds).
- In our case we have now three mean values that are the mean of data points (red circles) in each group.



- Fourth: three new groups are determined based on their proximity to the mean values (yellow diamonds).
- The new mean values (yellow diamonds) play the same role as the random numbers in the first stage (blue diamonds).



- Fifth: this process is repeated
  - new mean values are calculated.
  - new groups are identified.



- Six: this process is repeated and repeated again
  - new mean values are calculated.
  - new groups are identified.
- Until: no changes are observed in the mean values
- In this stage the final clusters are identified.





- There are five important points that should be taken into account:
  - 1) Kmean CL can only be used to find the natural groups among continuous variables (MEAN!!)

2) The units of variables should not be necessary the same

- Example: We can include expenditures on books, number of hours spent on family gathering, number of social connections and so on.

- However, we should standardize all the variables that is we should put different variables on the same scale.

Zx = <mark>[observation i of var x]</mark> – [mean of var x] / <mark>[standard</mark> deviation]

3) Kmean CL is highly sensitive to the presence of outliers (MEAN!!)

- Usually we should drop the outliers
- Otherwise, the results will be misleading (extra clusters or non-natural groupings)

4) we can evaluate kmean CL results (remember natural grouping is the primary task of kmean clustering.

- If our CL performs well, we will be able to find patterns that are consistent with theories or our expectations.

# 5) The most important point is to determine the optimal number of clusters.

- Remember in the first step we have to tell kmean that how many random numbers and consequently groups should it work with.
- There are several methods used to determine the optimal number of clusters

- **The main idea**: we conduct several cluster analysis where k (that is the number of clusters) increases from 2 to an arbitrary number.
- The maximum number of k is the number of observations where each observation is considered as one cluster.

- 1) Scree plot (Elbow Method)
- We need to review a few concepts to understand the method.
  - Total Sum of Squares (TSS)
  - Within Clusters Some of Squares (WCSS)

- Total sum of square:  $\sum_{i=1}^{\infty} (x_i \overline{x})^2$
- Each sets of observations have a mean value.
- We calculate the difference between each observation and the mean and square the differences.
- We sum the values and we will get TSS

- Lets say we have 5 observation: c(5, 9, 2, 10, 4)
- The average of these 5 observations is equal to 6.
- TSS= 46=  $(5-6)^2 + (9-6)^2 + (2-6)^2 + (10-6)^2 + (4-6)^2$

- WCSS measures the variability of observations within a cluster
  - Each cluster contains a series of observations.
  - Each set of observations has a mean value.
  - Total sum of square for each cluster is WCSS.
  - For two clusters with the same number of observations, the smaller WCSS means the observations are closer together
- Sums of WCSS is the primary measure used to determine the optimal number of clusters in elbow method.
- So we conduct several cluster analysis for a same datasets.
- For the book expenditures datasets, we assumed 3 clusters.
- Now lets use R to use elbow method to determine optimal number of clusters.

- We need to install and load the following packages:
- library(tidyverse) # data manipulation
- library(cluster) # clustering algorithms
- library(factoextra) # clustering algorithms & visualization
- library(NbClust) #a very good package for determining the optimal number of clusters.

We start with k=1 (no cluster) and go up till k=5, and record WCSSs

*set.seed(123)* #because kmean starts with a random procedure we use set.seed() to reproduce the results if necessary.

*book\_k <- kmeans(book, k, nstart = 25) #* we store the results in **book\_k**. Kmean is the function. **book** is the dataset name. **K** is the number of clusters. Finally nstart=25 is related to the initial stage of clustering. We tell the function to start with 25 initial points in the datasets (remember the blue diamonds) and choose the best ones.

Now we record the results (WCSS) as k goes from 1 to 5

- K=1: <mark>23.16</mark>
- K=2: 5.07 +4.78 = <mark>9.85</mark>
- K=3: 3.3 + 0.53 + 1.07= 4.9
- K=4: 1.5 + 0.53 + 0.12 + 0.5 = 2.75
- K=5: 0.23 + 0.05 +0.13 + 0.7 + 0.53 = 1.8



- We can see two elbows (kinks) at k=2 and k=3.
- If we are uncertain about the number of clusters we should use other methods.
- Milligan and Cooper (1985) tested 30 methods used to determine the optimal number of clusters.
- They conclude Calinski and Harabasz (CH) method outperforms other methods.

- CH method is more straight forward than scree plot.
- The formula for CH method also contains the information about WCSS.
- However, the optimal number of clusters is determined based on the the highest CH index.

- The R code for getting CH index is:
  - ch <- NbClust(book, min.nc=2, max.nc=5, method = "kmean", index = "ch")
  - The results are stored in *ch* <-. *NbClust* is the function under a package with the same name. (*book*, is the name of the dataset. *min.nc=2*, *max.nc=5* are two parameters telling NbClust function to report CH index for k=2 to k=5. *method= "kmean"* is the CL method. *index="CH"*) determine the calculation method that is CH.

print(ch) renders the following results



We can also plot the CH index for different numbers of clusters using the following code:

*plot(ch\$All.index, type="b").* So ch stores a series of information one of which is All.index that shows the number of clusters and their corresponding CH index.
 ch\$All.index calls for that component. type="b" means the plot type is both line and point.



#### Just in Case you want to use ggplot

ch\_all <- as.data.frame(ch\$All.index)</pre>

ggplot(ch\_all, aes(c(2:5), ch\_all\$`ch\$All.index`, fill=factor(c(2:5)))+
geom\_col()



- CH index tells us 5 is the best number of clusters (CH index for k=5 is equal to 37.8)
- However, the CH index for k=3 is equal to 35.2
- CH index for both k=3 and k=5 are close to each other.
- Considering the results of elbow method, we go with 3 clusters because both CH and elbow methods point to k=3.
- If we wanted to follow only one method, CH method is preferred

#### •So the final command is:

-book3 <- kmeans(cluster, 3, nstart = 25) -print(book3) -Cluster means: (the following table shows the mean of expenditures on both types of books across 3 clusters).



- We can also plot the CL using the following command:
- ggplot(book, aes(kids, fiction, colour=factor(book3\$cluster))) + geom\_point(aes(size=5))
- We simply ask R to use ggplot to render scatter plot for two variables of kids and fiction from book dataset. However, the colouring should be based on the cluster component of book3 (book3\$cluster) where we stored the results of CL with k=3.





- We now are going to extract dietary patterns of Canadian adults.
- You should be familiar with most of the coding in this part.
- We first need to look at the main dataset

- The dataset includes information about foods intakes, nutrients intakes and socioeconomic status of adults in Canada.
- •We use 9 variables indicating the intakes of 9 different food groups in servings for CL.
- •The food intakes variables are adjusted for 2000 Kcal of energy intake (that is if an adult eats 6 servings of grains and his/her energy intake is 2500 Kcal, he/she eats 4.8 =(2000\*6)/(2500) servings of grains per 2000 Kcal of energy (a solution for outliers)

• The dataset called "cluster\_data" includes all information. However, for CL we make a new data frame that includes the food intakes only.

#### •new <-cluster\_data</pre>

- •we tell R to store *cluster\_data* to *"new"*.
- •We use the dataset called "*new*" and make the food intakes dataset from it.
- food <- new %>%
  - select(starts\_with("adj"))
- •This command use "*new*" dataset and and then (%>%)select only those variables whose name *start with "adj*"

- Looking at "food" dataset, we have three variables including adj\_fruits, adj\_veg\_nopot\_and adj\_fruitveg. The variable adj\_fruitveg is the sum of two other variables, so we tell R to drop the other two variables
- *food <- new %>%*

select(-adj\_fruits, -adj\_veg\_nopot)

- •Optimal number of clusters using *fviz\_nbclust* function (elbow method)
- •fviz\_nbclust(food, kmeans, method = "wss", k.max = 7) + labs(subtitle = "Elbow Method") +

scale\_y\_continuous(breaks = scales::pretty\_breaks(n = 15)) +
theme(

axis.title.x = element\_text(size = 20), axis.text.x =
element\_text(size = 15), axis.text.y = element\_text(size = 20),
axis.title.y = element\_text(size = 20), plot.title =
element\_text(size=18))

- We use *fviz\_nbclust* and tell it to choose *food* dataset, conduct *kmean* CL, and find optimal number of clusters using method= "wss" (within cluster sum of squares).
- The rest of commands were discussed in GVC lecture and are only for better visibility of graph.



- We also going to use CH index to confirm the results of elbow method.
- ch\_ind <- NbClust(food, min.nc=2, max.nc=7, method = "kmean", index = "ch")
- print(ch\_ind)
- plot(ch\_ind\$All.index)
- We use *NbClust* function and tell it to use food dataset, with *minimum* number of clusters *=2* and *maximum=7*, the CL method is kmean. We also tell the function that we want the "*CH*" index

- Printing the results we see optimal k=3:
  - \$All.index
  - 2 <mark>3</mark> 4 5 6 7
  - 3913.67 <mark>3920</mark> 3276.4 3162.07 2882.5 2703.6
  - \$Best.nc
  - Number\_clusters Value\_Index
  - <u>3.00</u><u>3920</u>

- we choose k=3. We tell R to store kmean CL results in *food\_cl*.
- We use *set.seed(#)* in case we want to reproduce the results.
- We tell R to conduct *kmean* CL for dataset of food where *k* = 3 and *nstart*=40. Finally we want only to see the main results (next page) by printing the centres only
- *set.seed(1234)*
- *food\_cl <- kmeans(food, 3, nstart = 40)*
- print(food\_cl\$centers)

	Cluster 1	Cluster 2	Cluster 3
	Medium Quality	High Quality	Low Quality
Whole Grain	1.6	1.3	0.4
Refined Grain	2.8	3.3	7.8
Dairy Product	1.7	1.5	1.4
Red Meat	0.7	0.6	0.6
White Meat	0.8	1.1	0.6
Pulses and Nuts	0.5	0.5	0.3
Eggs	0.3	0.3	0.2
Processed Meat	0.3	0.2	0.3
Fruits and Vegetables	2.7	10.3	2.5

- Now we can evaluate the CL results by examining the prevalence of few socioeconomic characteristics across clusters
- To make everything a bit easier we add the cluster results to the "new" datasets
- so we tell R to make a new variable called *cluster3* in the "new" dataset whose values are the values of "*cluster*" variables in *food\_cl* where we stored kmean CL results.
- new\$cluster3 <- food\_cl\$cluster</li>

- Cross Tabulation
- We use "epiDisplay" package
- The following lines of codes tell R to use function *tab1* to report the distribution of *clusters* in the *"new"* data set. It also prints the results in a graph where the *bar values* are *percent* values
  - tab1(new\$cluster3, bar.values = "percent")

#### Distribution of new\$cluster3



%

3

- We also would like to know the prevalence of males, immigrants, and those with university degrees across clusters identified.
- This is called cross tabulation and we use package "descr".

male <- crosstab(new\$male, new\$cluster3,</pre>

expected = F, prop.r = T, prop.c = F, prop.t = F,prop.chisq = F, chisq = T, missing.include = F, format = "SPSS". dnn = "label", xlab = "Clusters", ylab = "Male", main = "", plot = getOption("descr.plot"))

We tell R use crosstab function, put male on column and cluster3 on rows, reporting expected values is FULSE (F), row percentages is True (T), column percentages is F, total percentages is F, chi square of proportion is F, chi square value is T, including missing values is F, table format is the same as SPSS, the rest of codes are related to the plot shown in viewer pane.
	new\$ma]	le	
new\$cluster3	0	1	Total
1	2578	2667	5245
1	49.2%	50.8%	48.3%
2	1165	567	1732
	67.3%	32.7%	15.9%
3	1946	1942	3888
5	50.1%	49.9%	35.8%
Total	5689	5176	10865
Statistics for	All Tabl		20
Statistics for	ALL TUD	Le Fucto	15
Pearson's Chi-	squared	test	
Chi^2 = 184.17	'2 d	.f. = 2	p <2e-1

- University Degree across Clusters:
- We first make a dummy variable called uni\_degree. it takes value of 1 if edu\_res4 is equal to 4.
- edu\_res4 is a variable includes 4 levels of education where the edu\_res4= 1 is high school drop out, =2 is high school diploma = 3 is trade diploma and finally =4 is university degree

• *new <- new %>%* 

mutate(uni\_degree = as.numeric(edu\_res4 == 4))

- We tell R that dataset to store new variable is "new" (new<-). We also tell R to use the "new" and then (%>%) create (mutate) a new variable called "uni\_degree".
- uni\_degree takes value of 1 if edu\_res==4

	new\$un	i_degree	<b>T</b> 1	
new\$cluster3		1	lotal	
1	3858	1359	5217	
	74.0%	26.0%	48.3%	
2	1112	611	1723	
	64.5%	35.5%	16.0%	
3	2946	911	3857	
	76.4%	23.6%	35.7%	
Total	7916	2881	10797	
Statistics for	All Tab	le Facto	rs	
Pearson's Chi-s	quared	test		
Chi^2 = 87.4440	2 0	d.f. = 2	p <2	2e-2

- The mean value of continuous variables over clusters
- We tell R to use the "new" dataset and then (%>%) group by cluster3 and then summarise (average) variable fsddekc (daily energy intakes in Kcal).
- *energy <- new %>%*

group\_by(cluster3) %>%
summarise(mean(fsddekc))

- We can also use ggplot to plot what we did in the previous stage
- So we tell R to use ggplot to make a column graph with the use of dataset energy where x is cluster3 and y is the mean of energy intake and the columns should be filled based on cluster3. The final line add values on the top of columns with geom\_text.
- ggplot(energy, aes(x=cluster3, y=`mean(fsddekc)`, fill=factor(cluster3)))+

geom\_col() +

geom\_text(aes(label = round(`mean(fsddekc)`,digits = 1), vjust = 0.5))
78



# The End!!!