

# UNCOVERING RETAIL CUSTOMER SEGMENTATION FROM LARGE TRANSACTION RECORDS

---

A NUANCED COMPARISON OF CLUSTERING ALGORITHMS  
USING ROUGH SET REDUCED DATASET

SYED AHMAD ZAKI BIN SYED SAKAF AL-ATTAS



# PROBLEM STATEMENT

Traditional clustering methods does not take into account inherent inconsistencies in data, hence the need to explore and compare alternative clustering algorithms



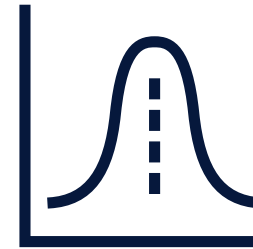
**TRADITIONAL  
K-MEANS**

Clustering Comparison



**FUZZY  
K-MEANS**

Clustering Comparison



**GAUSSIAN  
MIXTURE MODEL**

Clustering Comparison



**ROUGH  
SET**

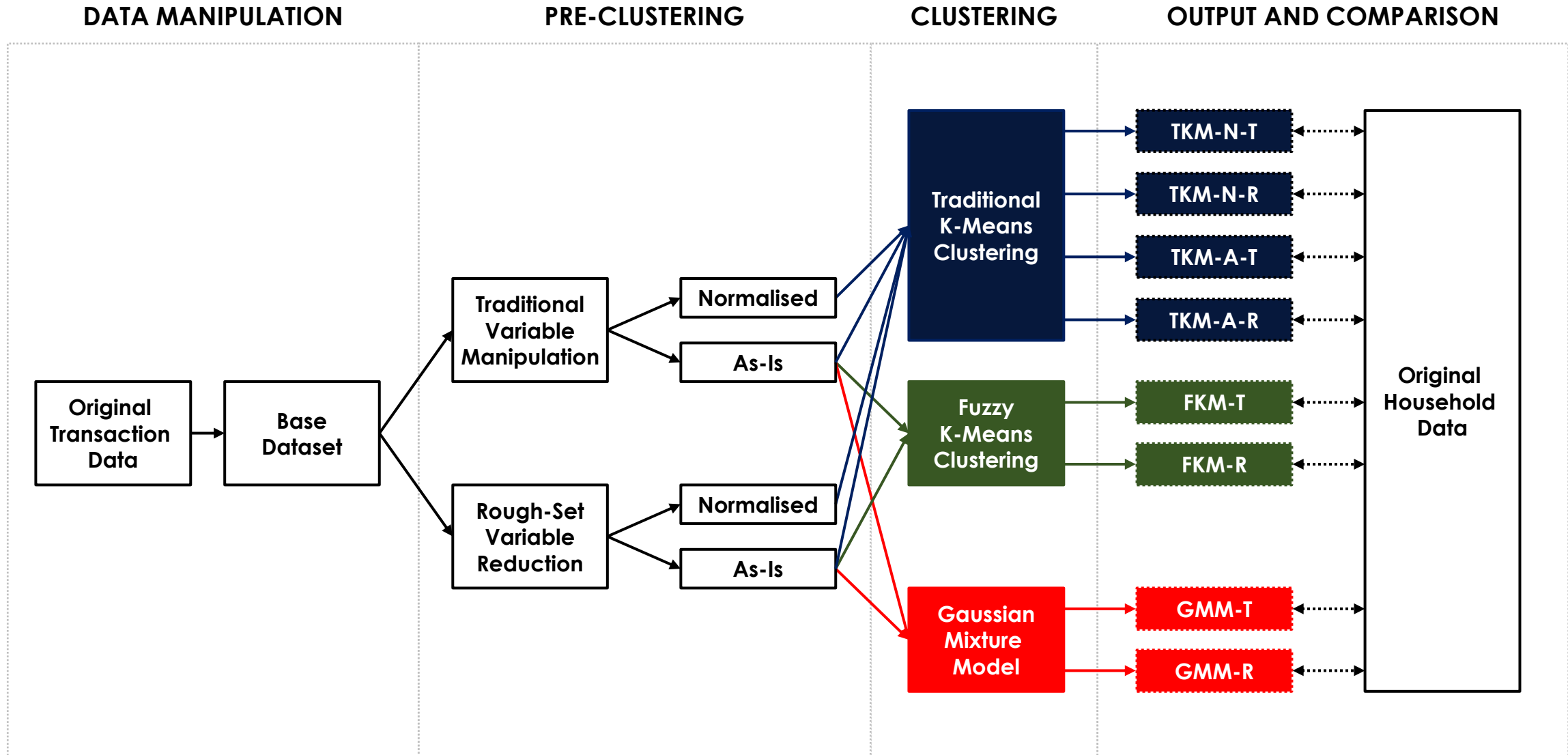
Feature Reduction

## OPEN-SOURCE INTEGRATION (R)

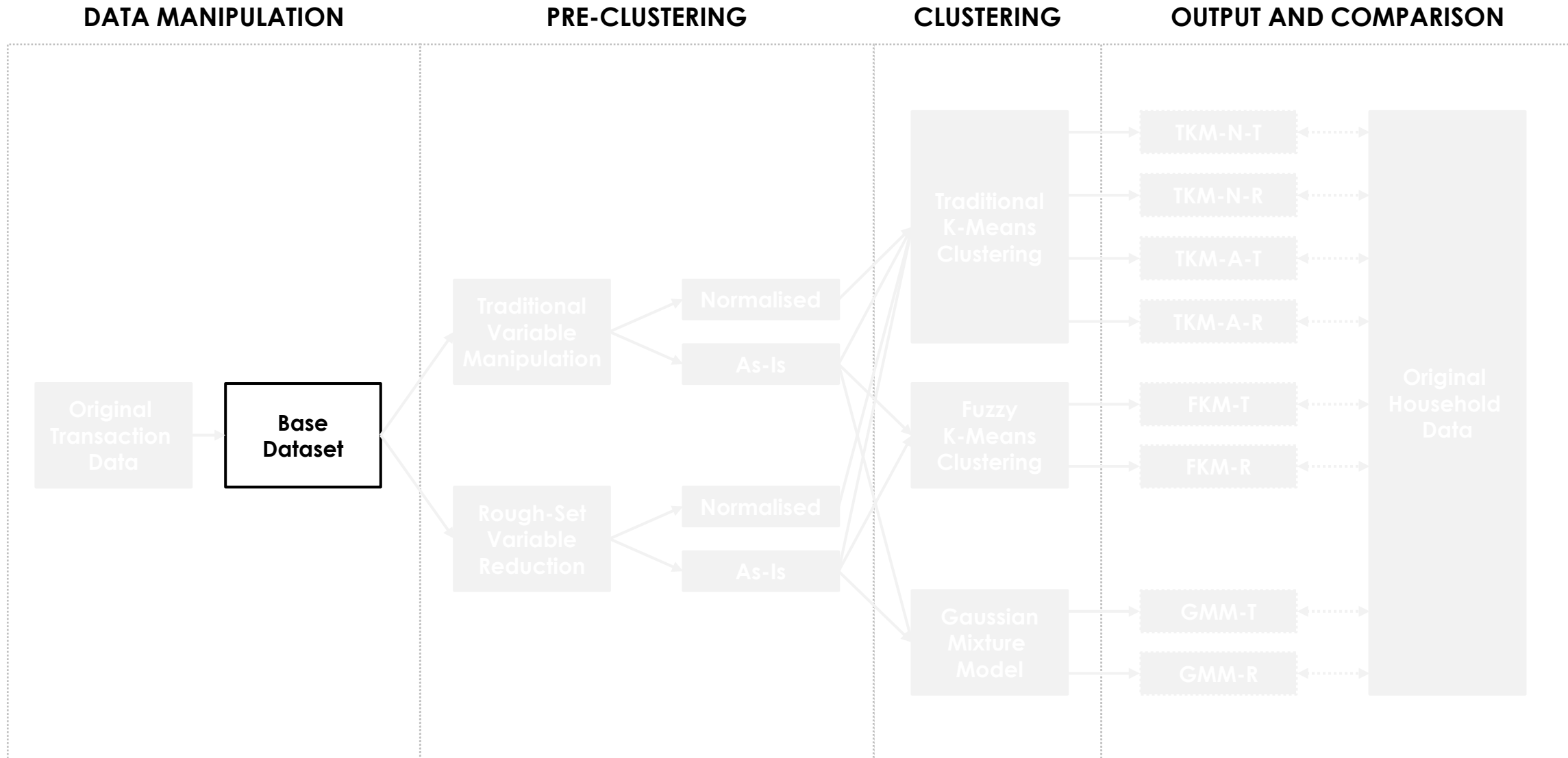
**SAS ENTERPRISE MINER 14.1**



# CAPSTONE WORKFLOW



# CAPSTONE WORKFLOW



# BASE DATASET (RFM MODEL INSPIRED)

## DUNNHUMBY'S RETAIL SHOPPING DATASET (TXNS)



### RECENCY

*(How long ago were their purchase?)*

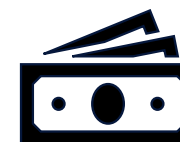
- Days Since First Order
- Days Since Last Order



### FREQUENCY

*(How often were the purchases made?)*

- Active Weekday Count
- Active Weekend Count
- Active Morning ( $6 \leq x < 12$ ) Count
- Active Afternoon ( $12 \leq x < 18$ ) Count
- Active Evening ( $18 \leq x < 0$ ) Count
- Active Late-Night ( $0 \leq x < 6$ ) Count
- Average Basket Count Per Active Weekday
- Average Basket Count Per Active Weekend
- Average Basket Count Per Active Morning
- Average Basket Count Per Active Afternoon
- Average Basket Count Per Active Evening
- Average Basket Count Per Active Late-Night



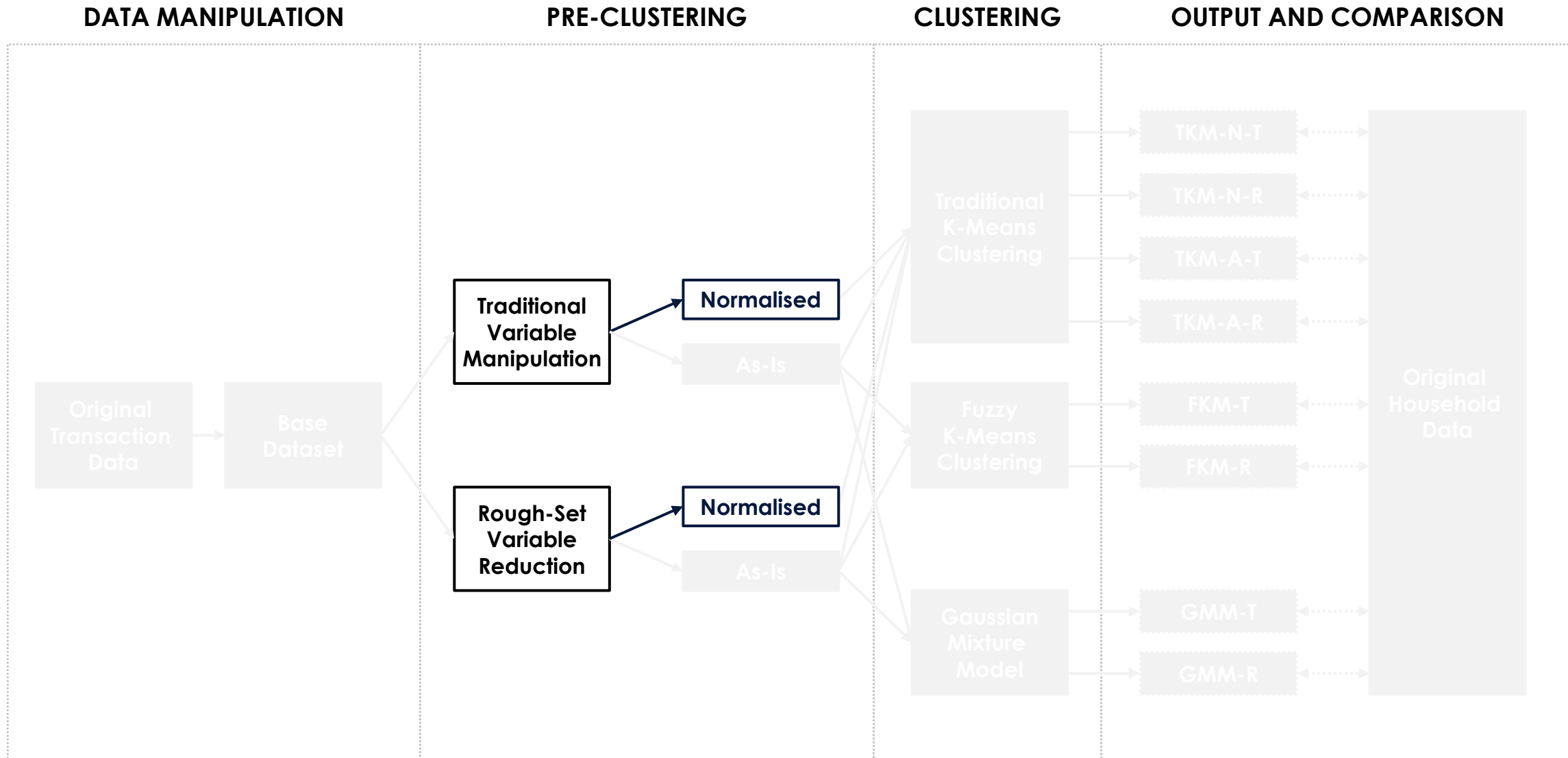
### MONETARY

*(How much money was spent?)*

- Total Spend
- Min Spend Per Basket
- Average Spend Per Basket
- Max Spend Per Basket
- Average Sales Value Per Qty
- Min Spend Per Active Week
- Average Spend Per Active Week
- Max Spend Per Active Week
- Average Spend Per Active Morning
- Average Spend Per Active Afternoon
- Average Spend Per Active Evening
- Average Spend Per Active Late Night
- Discount Used Per Active Weekday
- Discount Used Per Active Weekend

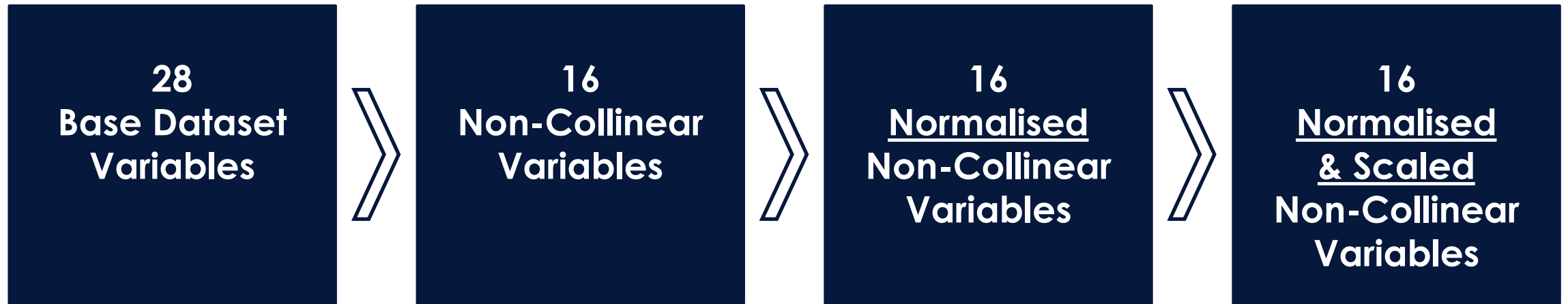
\*Active refers to time or periods when the transaction occurred

# CAPSTONE WORKFLOW



# PRE-CLUSTERING

## Traditional Variable Manipulation

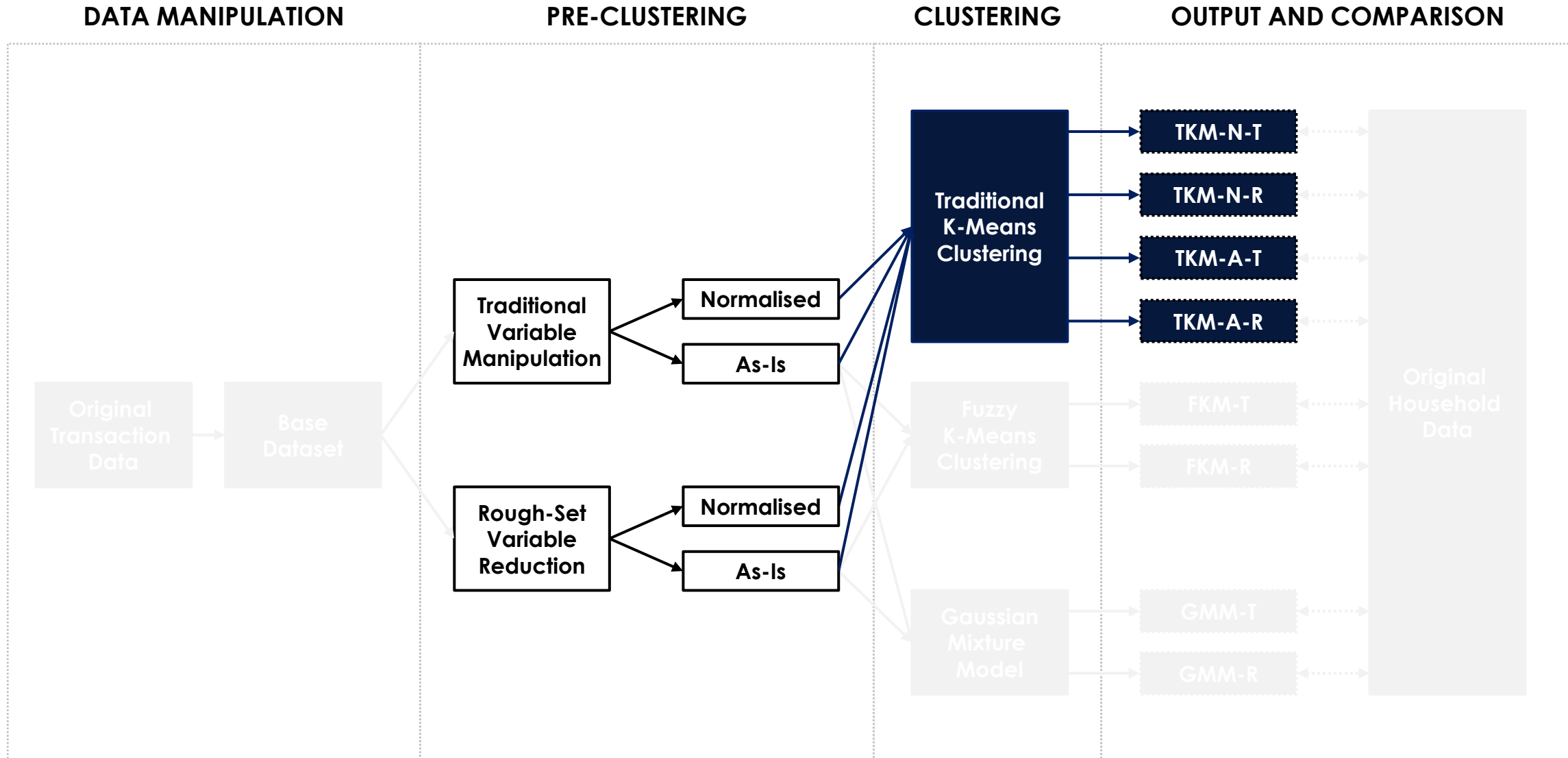


## Rough-Set Variable Reduction





# CAPSTONE WORKFLOW





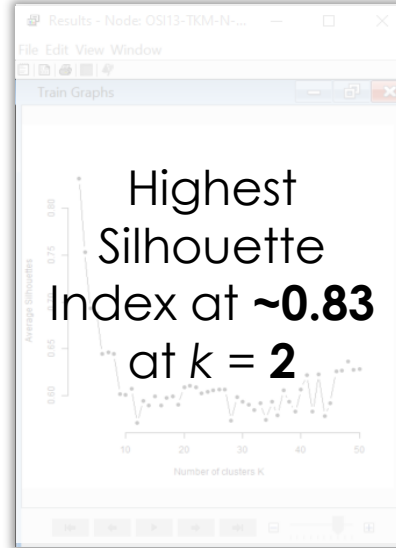
# ANALYSIS & RESULTS (TRAD. K-MEANS)

TKM-N-T

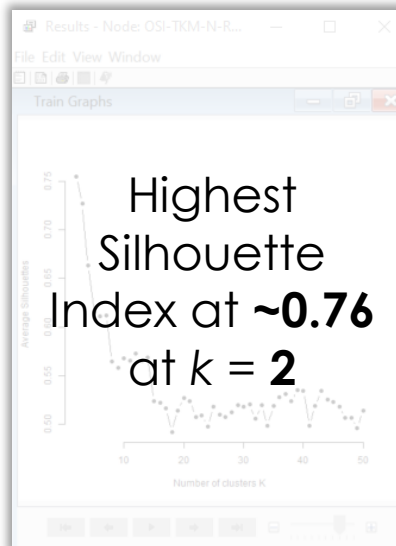


Average Silhouette Index  
( $k = 2$  to 50)

TKM-N-R



TKM-N-R3



TKM-A-T



Average Silhouette Index  
( $k = 2$  to 50)

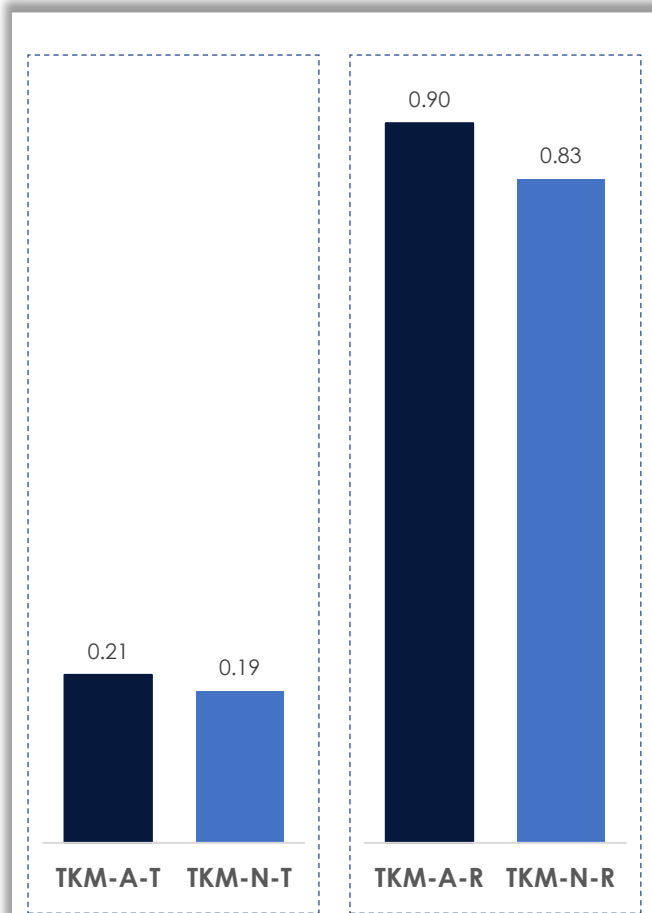
TKM-A-R



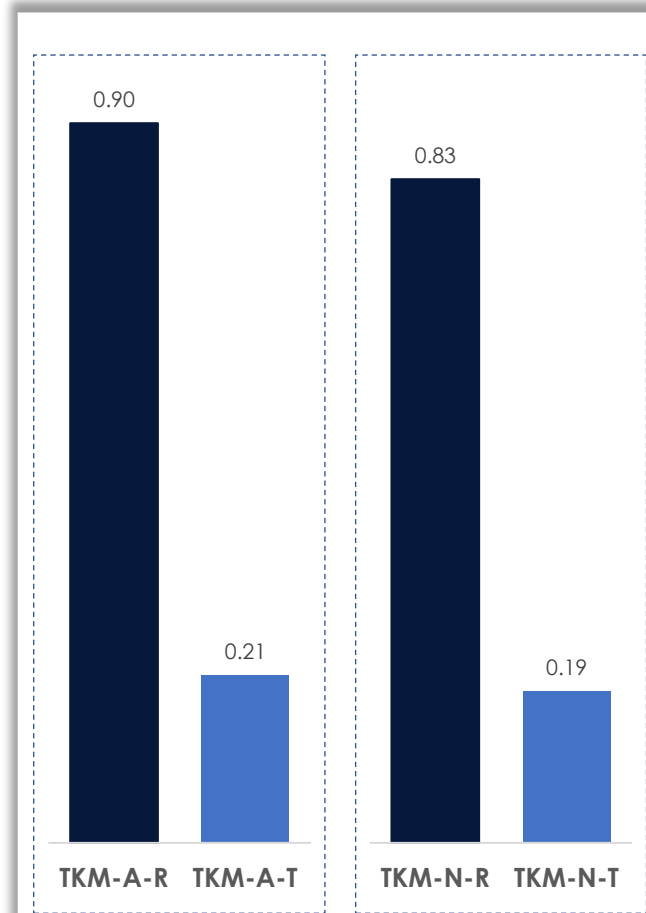
TKM-A-R3



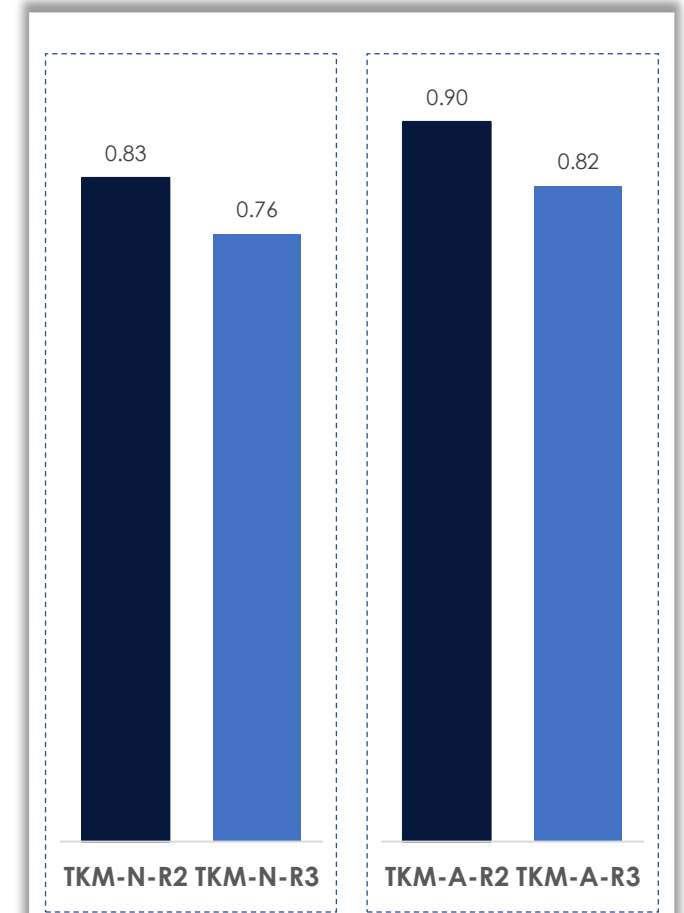
# SILHOUETTE COMPARISON (TRAD. K-MEANS)



AS-IS DATA HIGHER  
THAN **NORMALISED** DATA



REDUCT DATA HIGHER  
THAN **NON-REDUCT** DATA



REDUCT 2 DATA HIGHER  
THAN **REDUCT 3** DATA

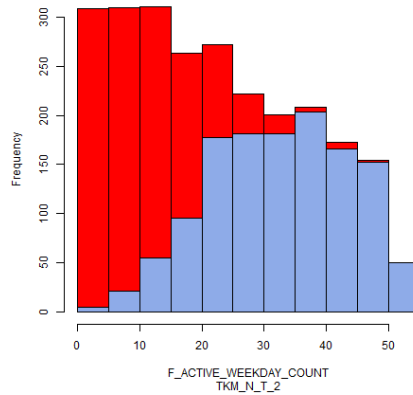
# CLUSTER COMPARISON (TRAD. K-MEANS)

Chosen One of 28 Transactional Attributes (F\_ACTIVE\_WEEKDAY\_COUNT) For Apple-to-Apple Comparison

## TKM-N-T

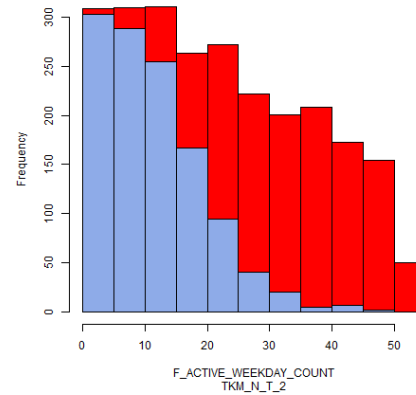
Cluster 1 (1,287 IDs)

Segment Profile of Overall(Red) vs. Cluster 1Blue



Cluster 2 (1,183 IDs)

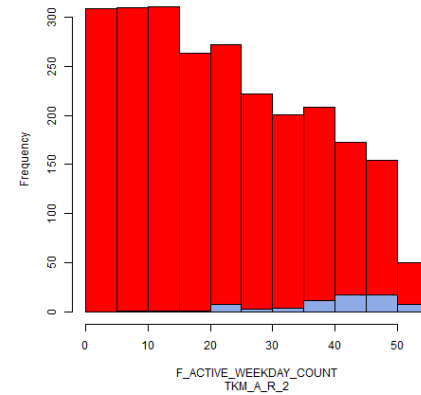
Segment Profile of Overall(Red) vs. Cluster 2Blue



## TKM-A-R

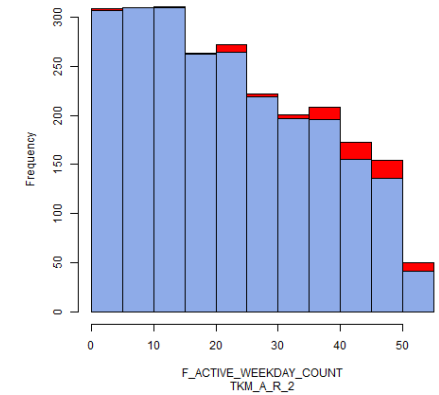
Cluster 1 (74 IDs)

Segment Profile of Overall(Red) vs. Cluster 1Blue



Cluster 2 (2,396 IDs)

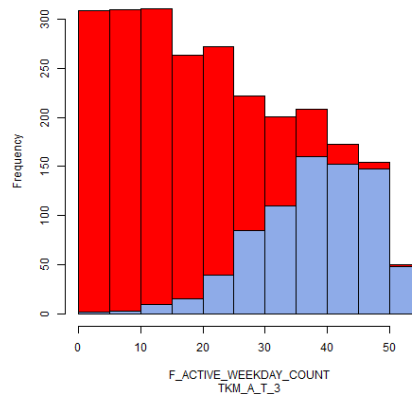
Segment Profile of Overall(Red) vs. Cluster 2Blue



## TKM-A-T

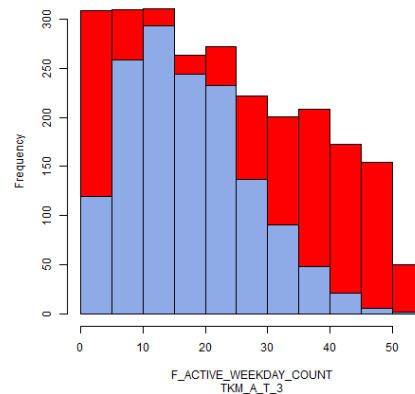
Cluster 1 (774 IDs)

Segment Profile of Overall(Red) vs. Cluster 1Blue



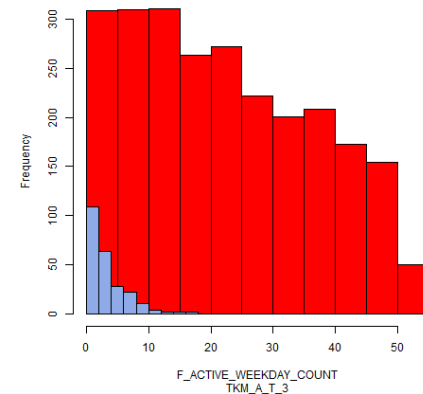
Cluster 2 (1,452 IDs)

Segment Profile of Overall(Red) vs. Cluster 2Blue



Cluster 3 (244 IDs)

Segment Profile of Overall(Red) vs. Cluster 3Blue



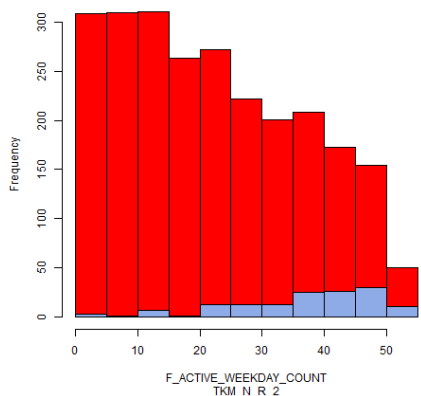
# CLUSTER COMPARISON (TRAD. K-MEANS)

Chosen One of 28 Transactional Attributes (F\_ACTIVE\_WEEKDAY\_COUNT) For Apple-to-Apple Comparison

## TKM-N-R

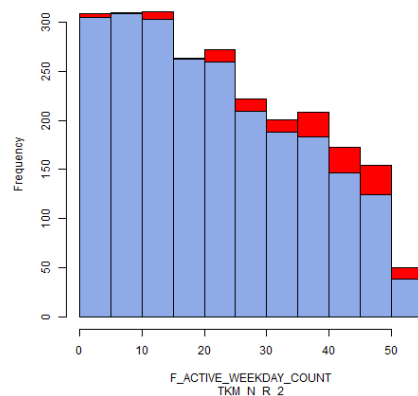
Cluster 1 (143 IDs)

Segment Profile of Overall(Red) vs. Cluster 1Blue



Cluster 2 (2,327 IDs)

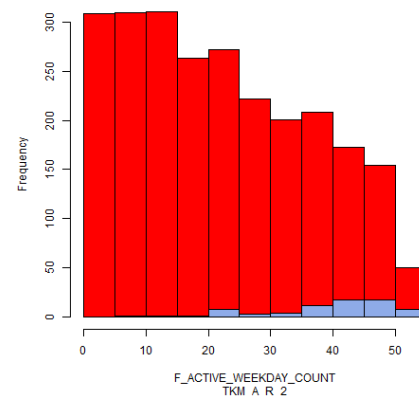
Segment Profile of Overall(Red) vs. Cluster 2Blue



## TKM-A-R

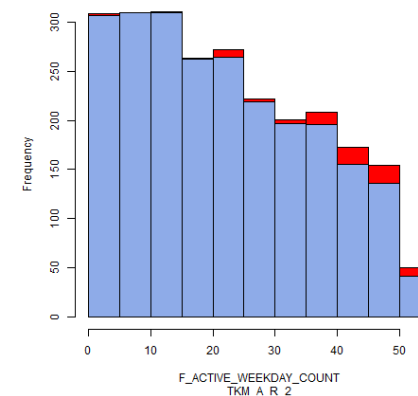
Cluster 1 (74 IDs)

Segment Profile of Overall(Red) vs. Cluster 1Blue



Cluster 2 (2,396 IDs)

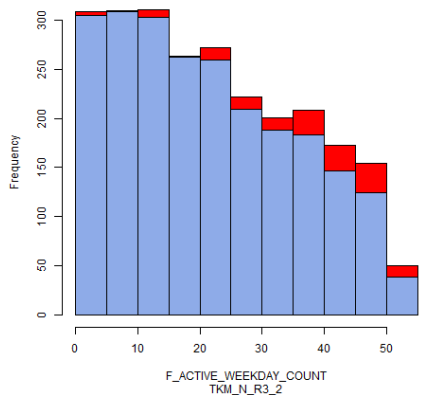
Segment Profile of Overall(Red) vs. Cluster 2Blue



## TKM-N-R3

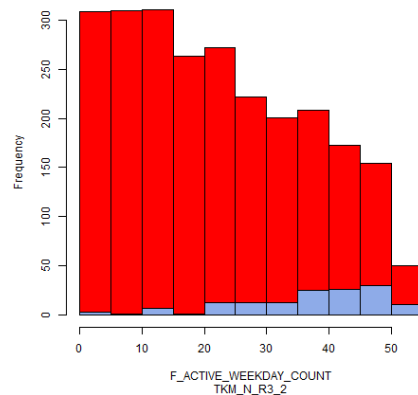
Cluster 1 (2,327 IDs)

Segment Profile of Overall(Red) vs. Cluster 1Blue



Cluster 2 (143 IDs)

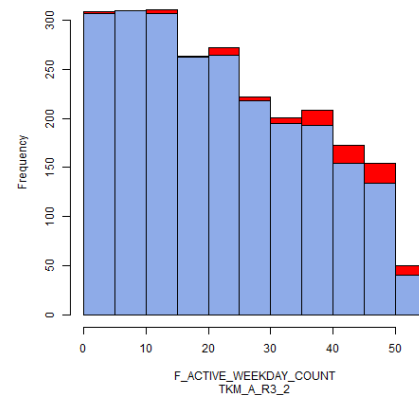
Segment Profile of Overall(Red) vs. Cluster 2Blue



## TKM-A-R3

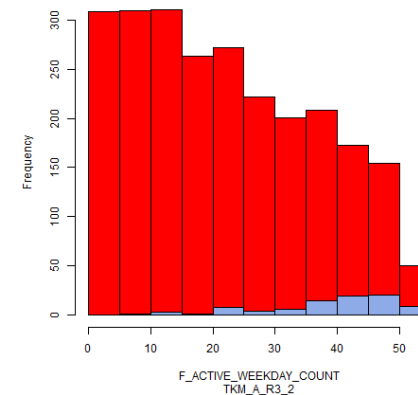
Cluster 1 (2,384 IDs)

Segment Profile of Overall(Red) vs. Cluster 1Blue



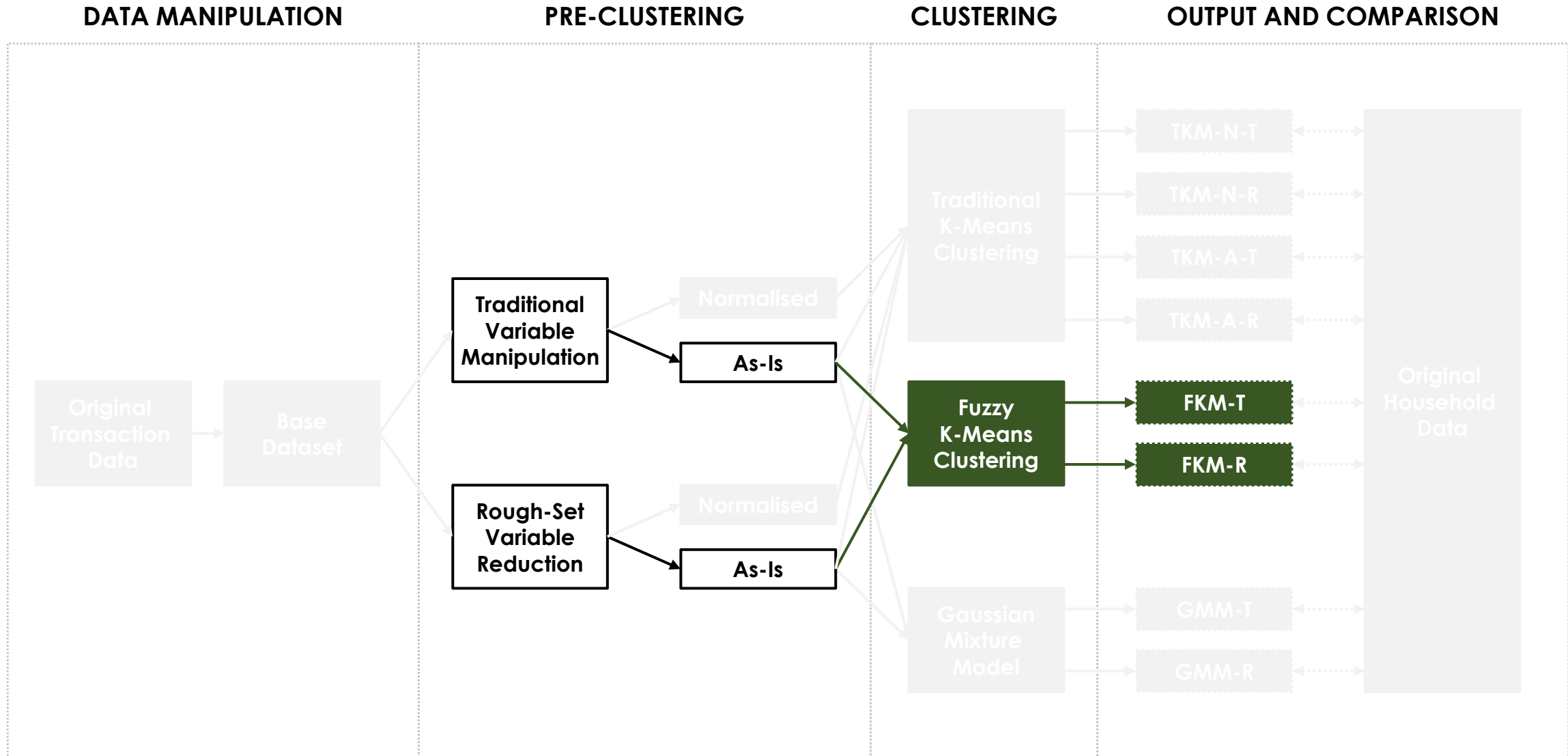
Cluster 2 (86 IDs)

Segment Profile of Overall(Red) vs. Cluster 2Blue





# CAPSTONE WORKFLOW



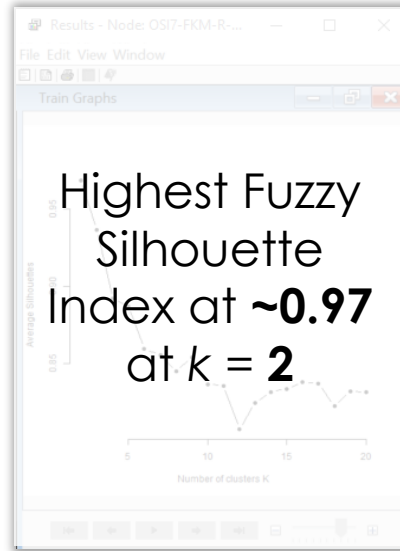
# ANALYSIS & RESULTS (FUZZY K-MEANS)

## FKM-T



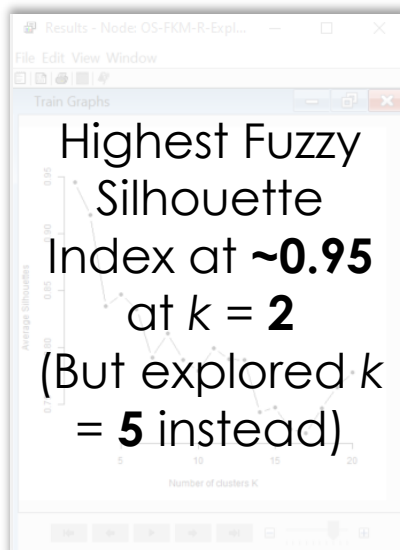
Average Silhouette Index  
( $k = 2$  to  $20$ )

## FKM-R

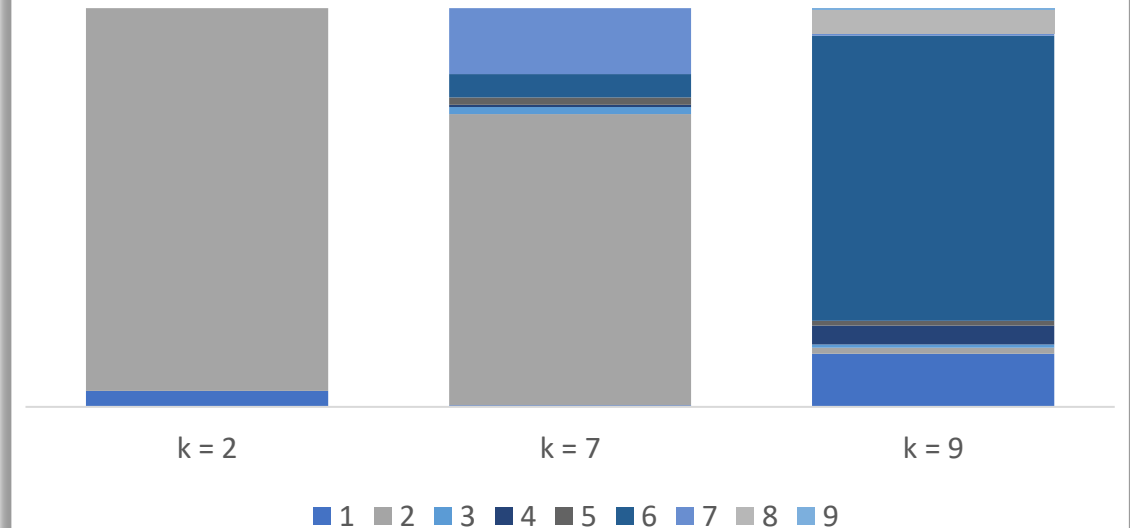


FKM-R had the highest silhouette index values amongst TKM and FKM. Possibility of exploring other cluster counts than just the one with the highest silhouette value

## FKM-R3



## CLUSTER PROPORTIONS (TOTAL ID = 2,470)



Highly uneven distribution, despite the high silhouette values

# CLUSTER COMPARISON (FUZZY K-MEANS)

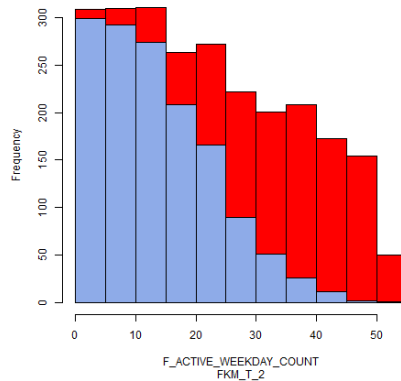
Chosen One of 28 Transactional Attributes (F\_ACTIVE\_WEEKDAY\_COUNT) For Apple-to-Apple Comparison

**FKM-T**

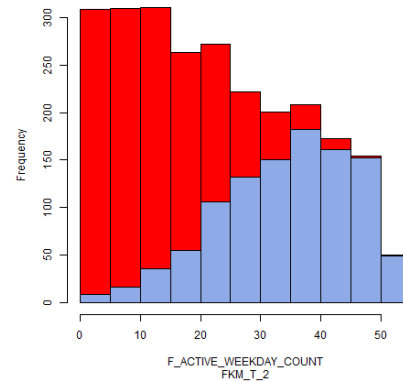
Cluster 1 (1,421 IDs)

Cluster 2 (1,049 IDs)

Segment Profile of Overall(Red) vs. Cluster 1Blue



Segment Profile of Overall(Red) vs. Cluster 2Blue

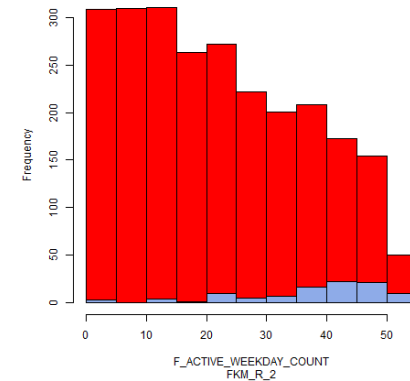


**FKM-R**

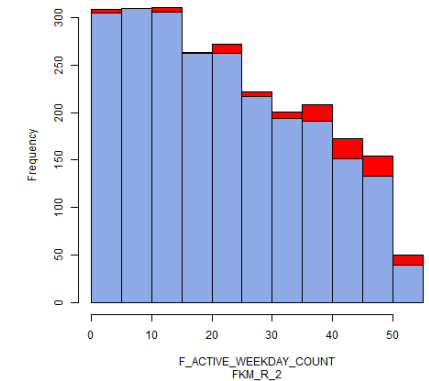
Cluster 1 (100 IDs)

Cluster 2 (2,370 IDs)

Segment Profile of Overall(Red) vs. Cluster 1Blue



Segment Profile of Overall(Red) vs. Cluster 2Blue



**FKM-R3**

Cluster 1 (57 IDs)

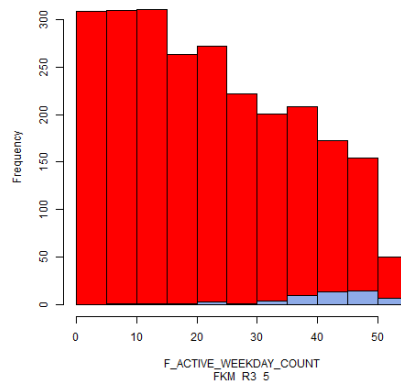
Cluster 2 (233 IDs)

Cluster 3 (59 IDs)

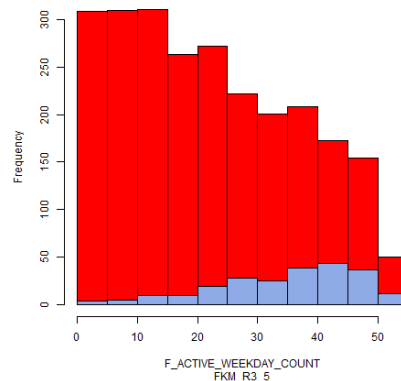
Cluster 4 (1,856 IDs)

Cluster 5 (265 IDs)

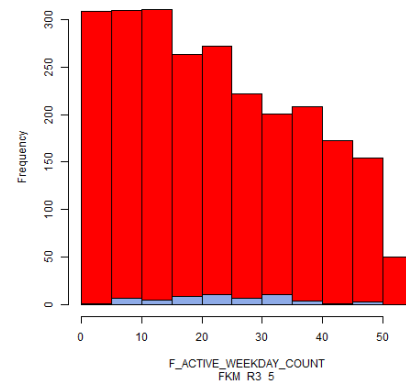
Segment Profile of Overall(Red) vs. Cluster 1Blue



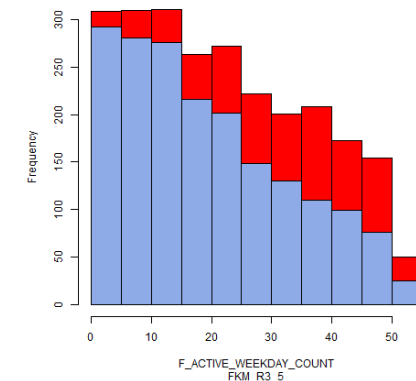
Segment Profile of Overall(Red) vs. Cluster 2Blue



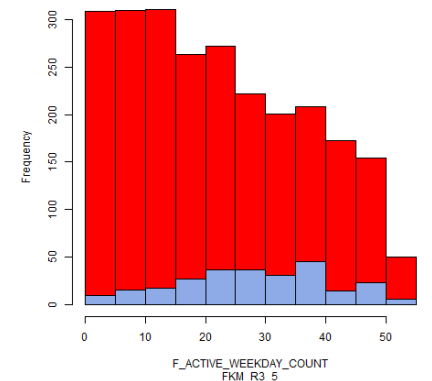
Segment Profile of Overall(Red) vs. Cluster 3Blue



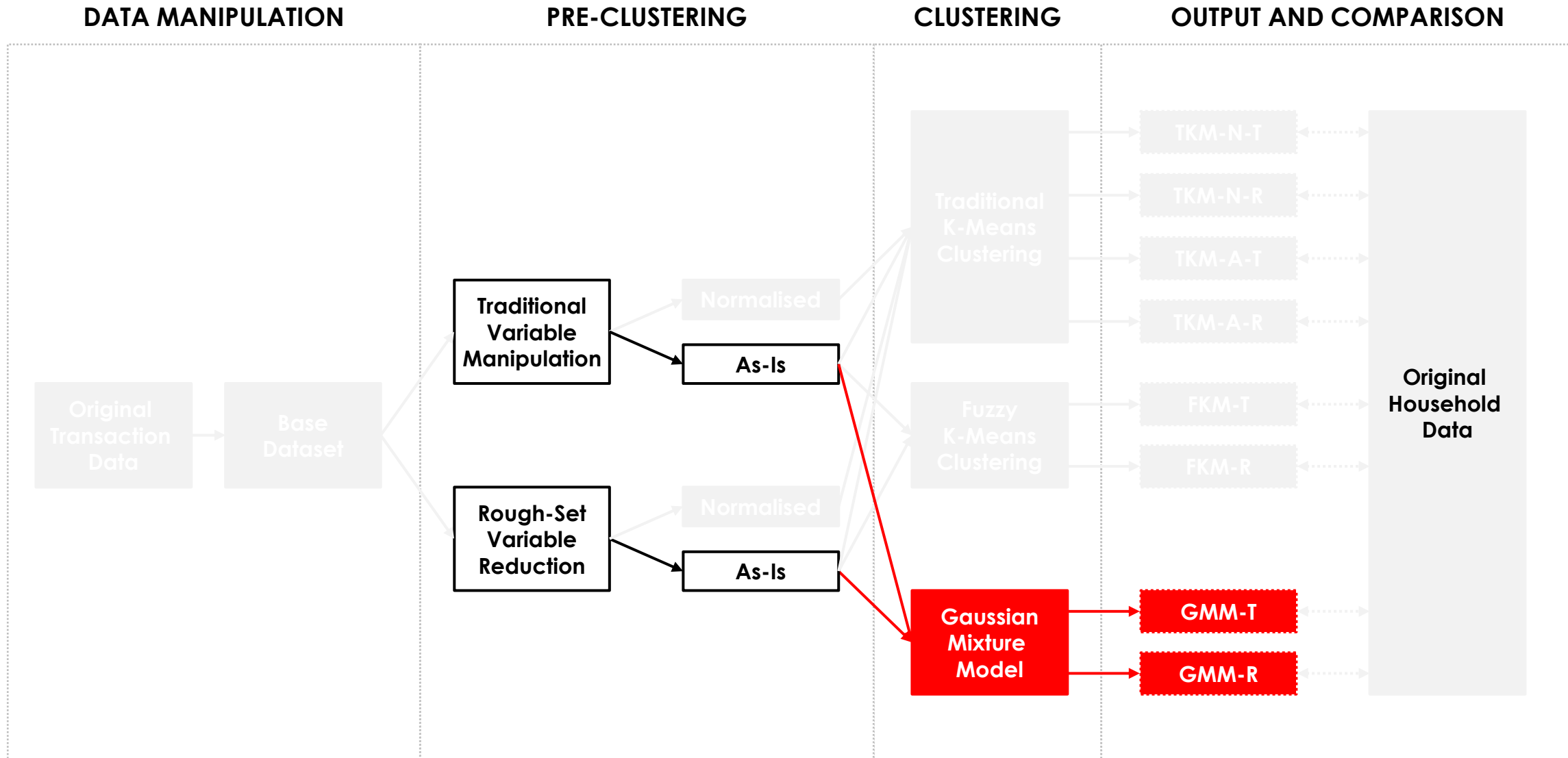
Segment Profile of Overall(Red) vs. Cluster 4Blue



Segment Profile of Overall(Red) vs. Cluster 5Blue



# CAPSTONE WORKFLOW

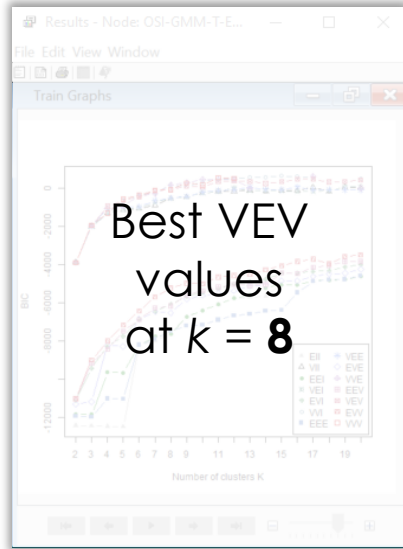






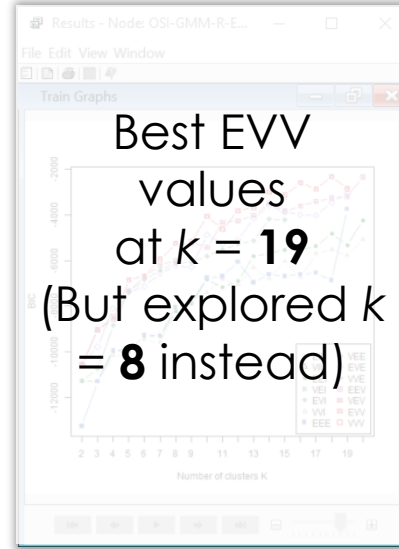
# ANALYSIS & RESULTS (GAUSSIAN)

## GMM-T



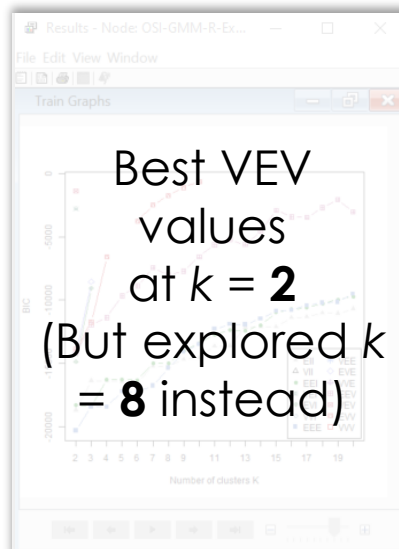
Average Silhouette Index  
(k = 2 to 20)

## GMM-R



GMM-R had optimal cluster count of 19, but it is too granular. Though sub-optimal, capstone super-imposed  $k=8$  onto GMM-R

## GMM-R3



## CLUSTER PROPORTIONS (TOTAL ID = 2,470)



GMM-T (8)

GMM-R (19)

GMM-R (8)

1 2 3 4 5 6 7 8 9 10  
11 12 13 14 15 16 17 18 19

Reduct dataset had highly uneven distribution

# CLUSTER COMPARISON (GAUSSIAN)

Chosen One of 28 Transactional Attributes (F\_ACTIVE\_WEEKDAY\_COUNT) For Apple-to-Apple Comparison

## GMM-T

Cluster 1 (283 IDs)

Cluster 2 (504 IDs)

Cluster 3 (298 IDs)

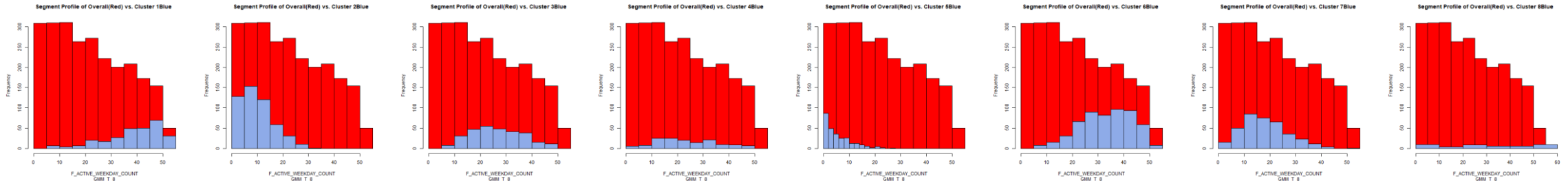
Cluster 4 (146 IDs)

Cluster 5 (274 IDs)

Cluster 6 (552 IDs)

Cluster 7 (368 IDs)

Cluster 8 (45 IDs)



## GMM-R (8)

Cluster 1 (78 IDs)

Cluster 2 (2,067 IDs)

Cluster 3 (231 IDs)

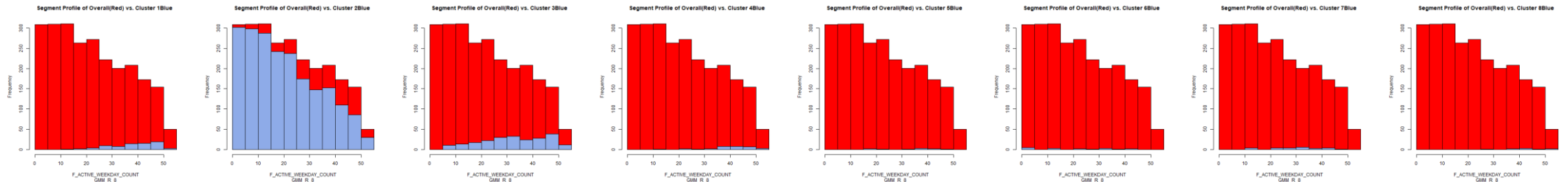
Cluster 4 (32 IDs)

Cluster 5 (10 IDs)

Cluster 6 (17 IDs)

Cluster 7 (26 IDs)

Cluster 8 (9 IDs)



## GMM-R3 (8)

Cluster 1 (54 IDs)

Cluster 2 (1,808 IDs)

Cluster 3 (103 IDs)

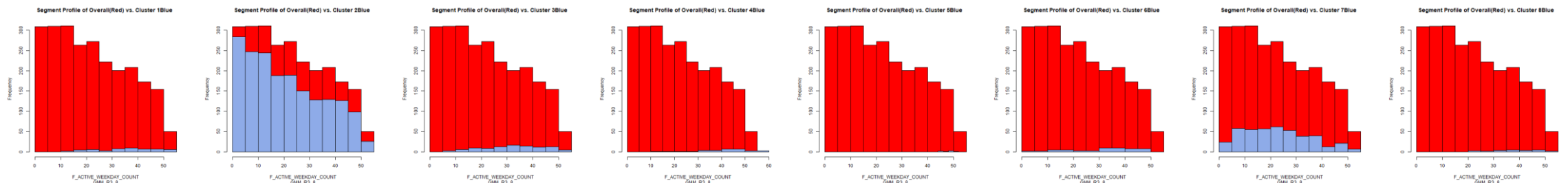
Cluster 4 (16 IDs)

Cluster 5 (7 IDs)

Cluster 6 (28 IDs)

Cluster 7 (429 IDs)

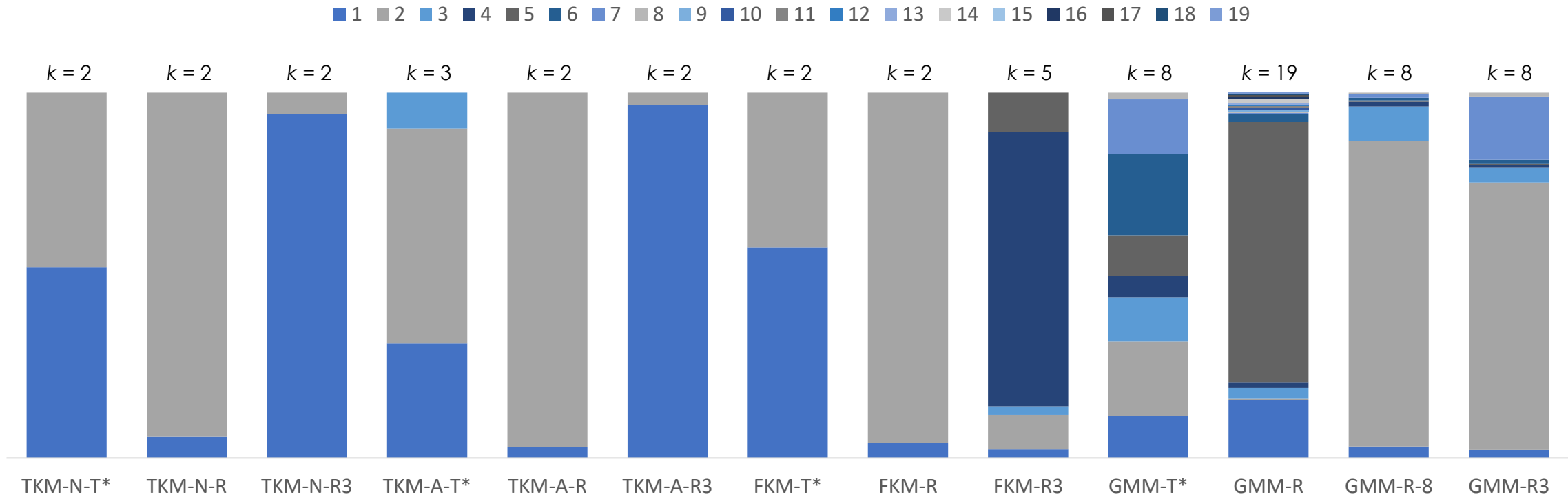
Cluster 8 (25 IDs)





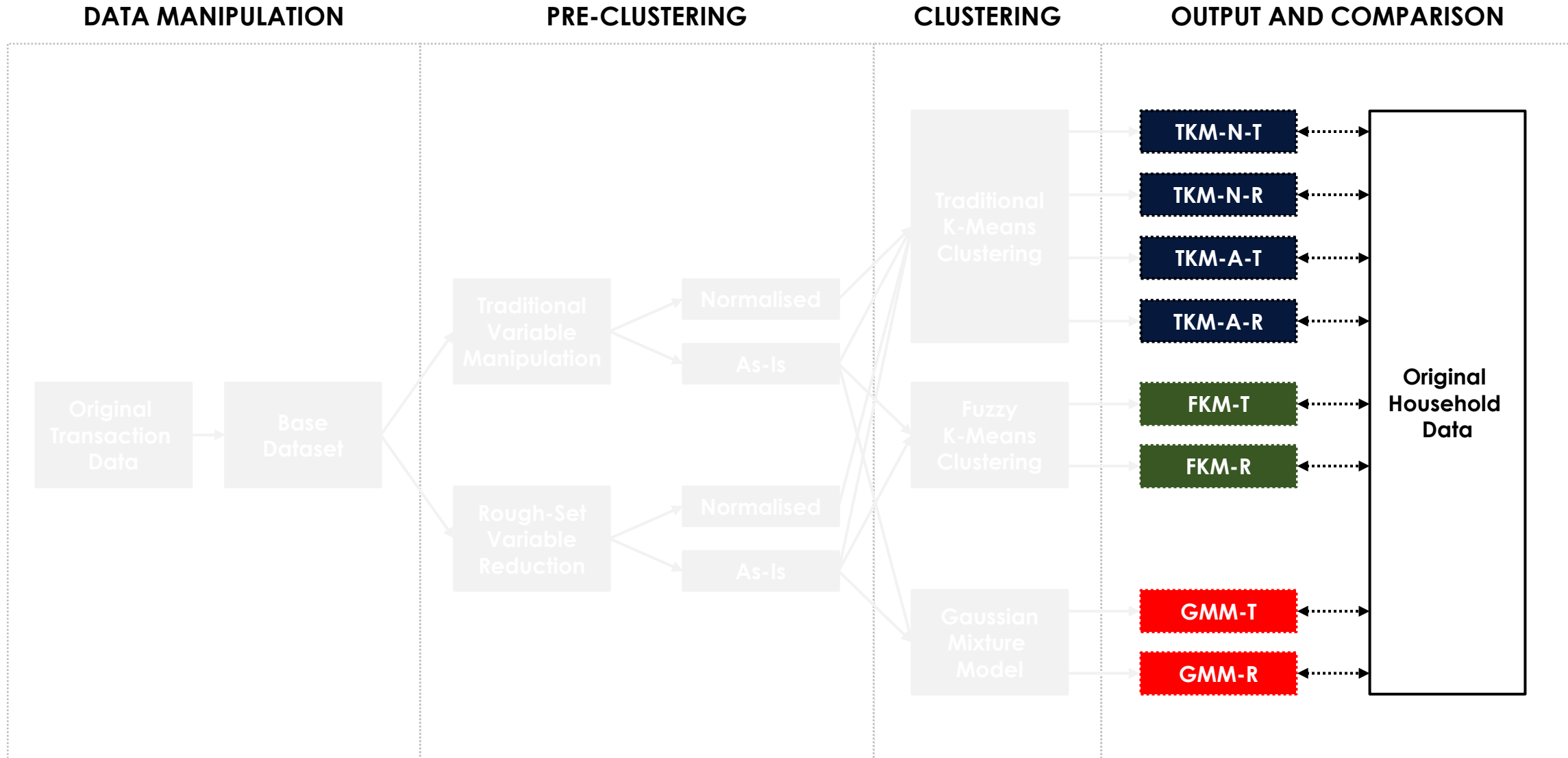
# ANALYSIS & RESULTS (OVERVIEW)

CLUSTER PROPORTIONS (TOTAL ID = 2,470)



- Optimised Cluster Counts Were Higher For GMM Than Other Clustering Algos (TKM And FKM)
- Non-Reduct Data\* Had More Evenly Spread Cluster Proportions Than Reduct Data
- Little Clustering Difference Between Reduct-2 and Reduct-3 Dataset

# CAPSTONE WORKFLOW



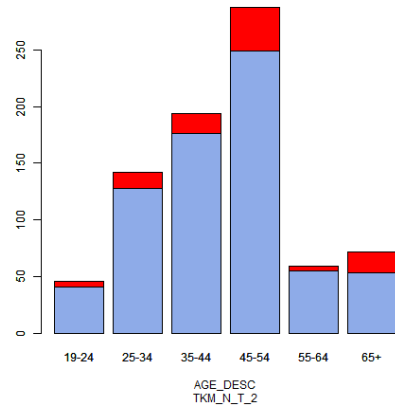
# CLUSTER COMPARISON (TRAD. K-MEANS)

Chosen One of 6 Demographic Attributes (AGE\_DESC) For Apple-to-Apple Comparison

## TKM-N-T

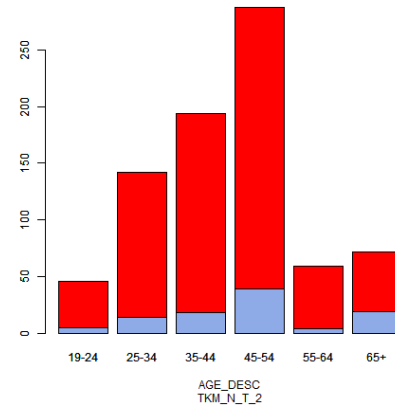
Cluster 1 (1,287 IDs)

Segment HH Profile of Overall(Red) vs. Cluster 1Blue



Cluster 2 (1,183 IDs)

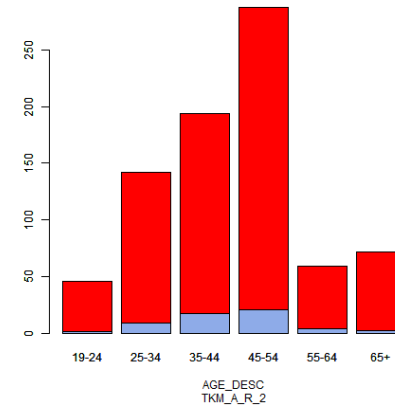
Segment HH Profile of Overall(Red) vs. Cluster 2Blue



## TKM-A-R

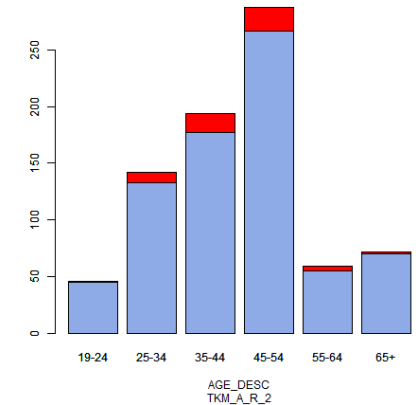
Cluster 1 (74 IDs)

Segment HH Profile of Overall(Red) vs. Cluster 1Blue



Cluster 2 (2,396 IDs)

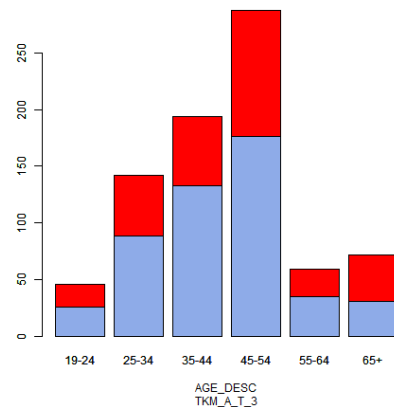
Segment HH Profile of Overall(Red) vs. Cluster 2Blue



## TKM-A-T

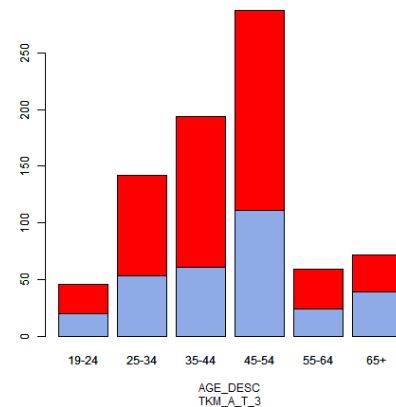
Cluster 1 (774 IDs)

Segment HH Profile of Overall(Red) vs. Cluster 1Blue



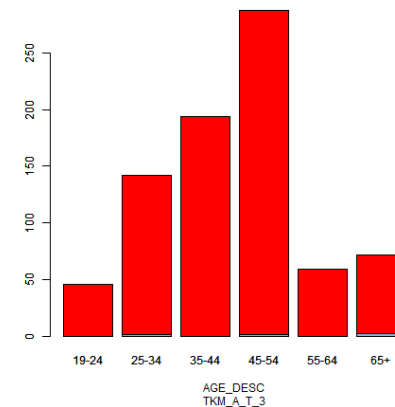
Cluster 2 (1,452 IDs)

Segment HH Profile of Overall(Red) vs. Cluster 2Blue



Cluster 3 (244 IDs)

Segment HH Profile of Overall(Red) vs. Cluster 3Blue



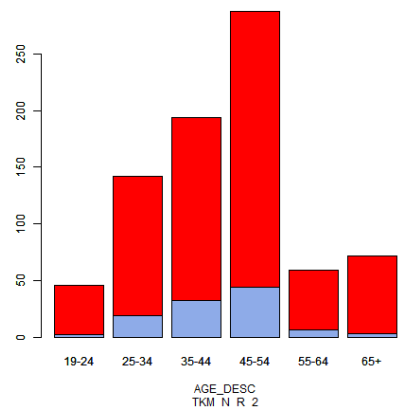
# CLUSTER COMPARISON (TRAD. K-MEANS)

Chosen One of 6 Demographic Attributes (AGE\_DESC) For Apple-to-Apple Comparison

## TKM-N-R

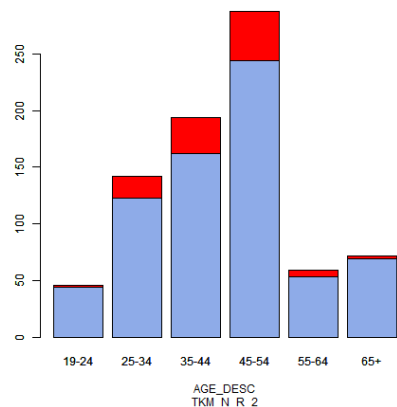
Cluster 1 (143 IDs)

Segment HH Profile of Overall(Red) vs. Cluster 1Blue



Cluster 2 (2,327 IDs)

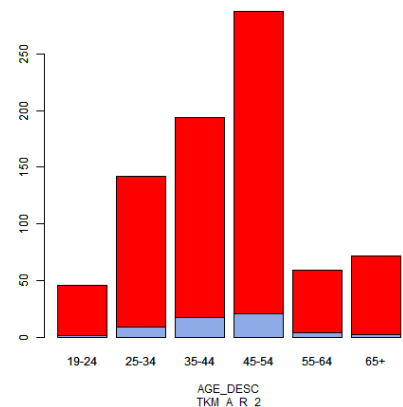
Segment HH Profile of Overall(Red) vs. Cluster 2Blue



## TKM-A-R

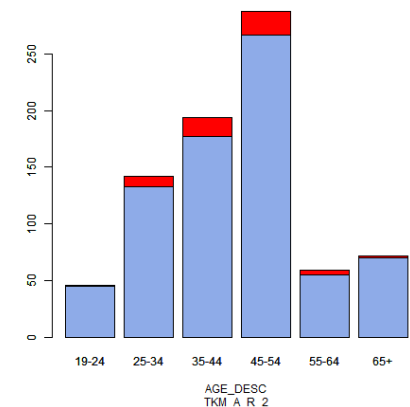
Cluster 1 (74 IDs)

Segment HH Profile of Overall(Red) vs. Cluster 1Blue



Cluster 2 (2,396 IDs)

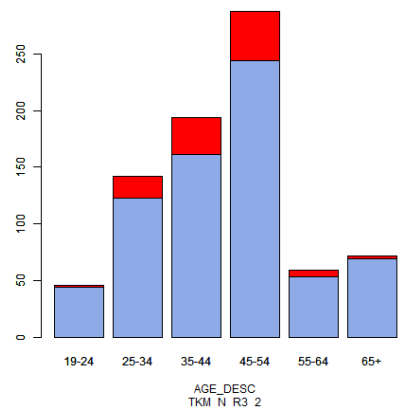
Segment HH Profile of Overall(Red) vs. Cluster 2Blue



## TKM-N-R3

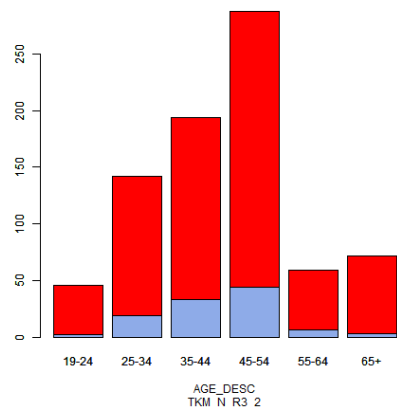
Cluster 1 (2,327 IDs)

Segment HH Profile of Overall(Red) vs. Cluster 1Blue



Cluster 2 (143 IDs)

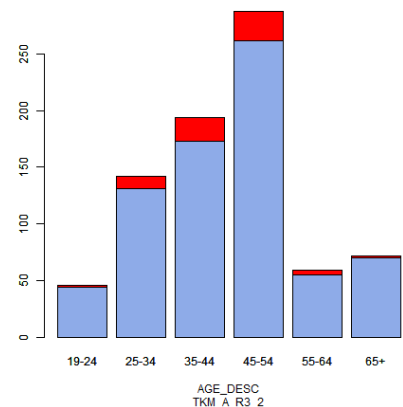
Segment HH Profile of Overall(Red) vs. Cluster 2Blue



## TKM-A-R3

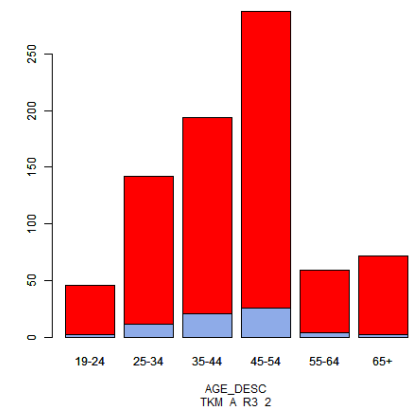
Cluster 1 (2,384 IDs)

Segment HH Profile of Overall(Red) vs. Cluster 1Blue



Cluster 2 (86 IDs)

Segment HH Profile of Overall(Red) vs. Cluster 2Blue



# CLUSTER COMPARISON (FUZZY K-MEANS)

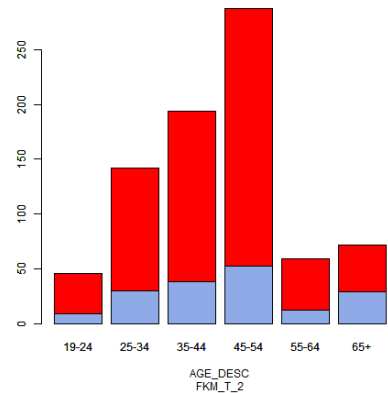
Chosen One of 6 Demographic Attributes (AGE\_DESC) For Apple-to-Apple Comparison

FKM-T

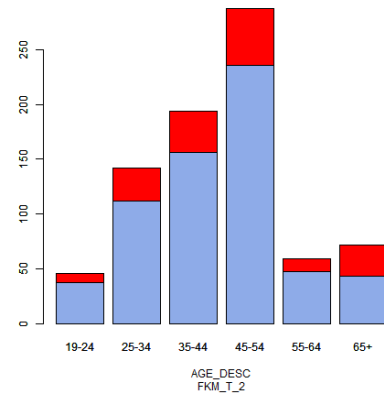
Cluster 1 (1,421 IDs)

Cluster 2 (1,049 IDs)

Segment HH Profile of Overall(Red) vs. Cluster 1Blue



Segment HH Profile of Overall(Red) vs. Cluster 2Blue

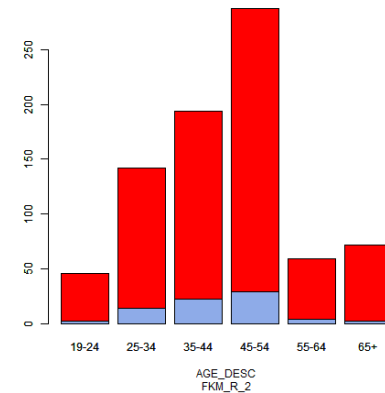


FKM-R

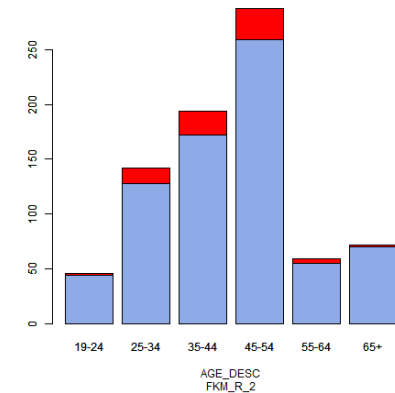
Cluster 1 (100 IDs)

Cluster 2 (2,370 IDs)

Segment HH Profile of Overall(Red) vs. Cluster 1Blue



Segment HH Profile of Overall(Red) vs. Cluster 2Blue



FKM-R3

Cluster 1 (57 IDs)

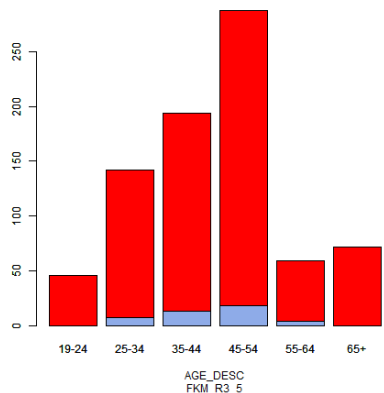
Cluster 2 (233 IDs)

Cluster 3 (59 IDs)

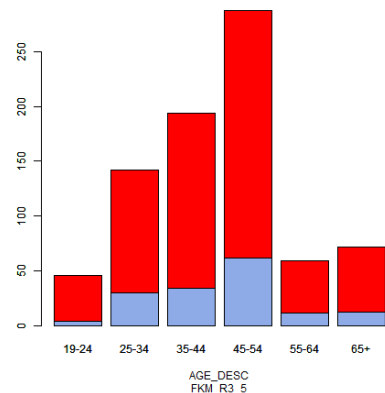
Cluster 4 (1,856 IDs)

Cluster 5 (265 IDs)

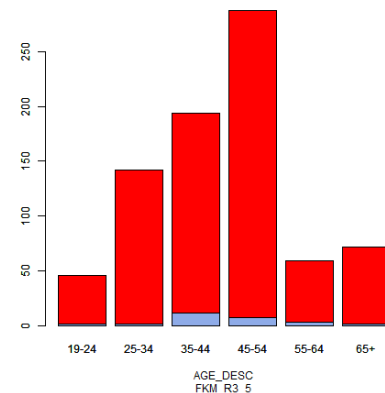
Segment HH Profile of Overall(Red) vs. Cluster 1Blue



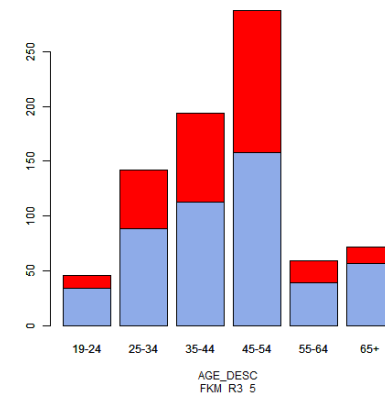
Segment HH Profile of Overall(Red) vs. Cluster 2Blue



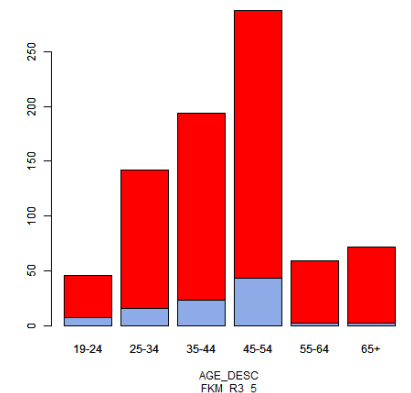
Segment HH Profile of Overall(Red) vs. Cluster 3Blue



Segment HH Profile of Overall(Red) vs. Cluster 4Blue



Segment HH Profile of Overall(Red) vs. Cluster 5Blue

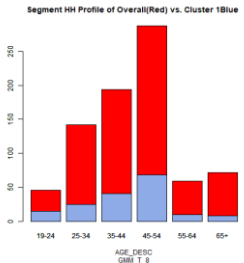


# CLUSTER COMPARISON (GAUSSIAN)

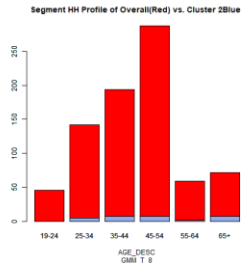
Chosen One of 6 Demographic Attributes (AGE\_DESC) For Apple-to-Apple Comparison

## GMM-T

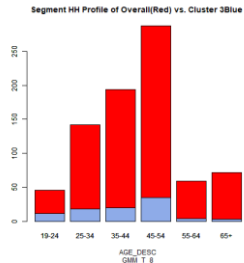
Cluster 1 (283 IDs)



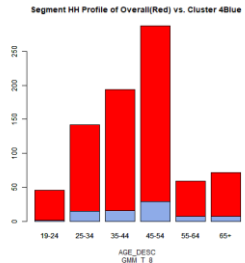
Cluster 2 (504 IDs)



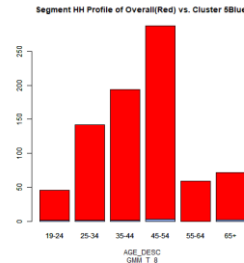
Cluster 3 (298 IDs)



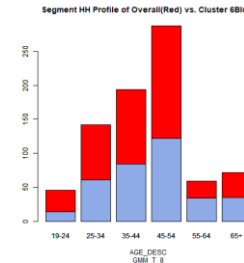
Cluster 4 (146 IDs)



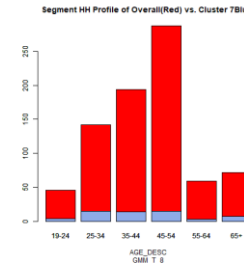
Cluster 5 (274 IDs)



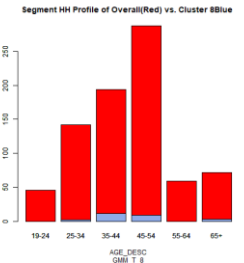
Cluster 6 (552 IDs)



Cluster 7 (368 IDs)

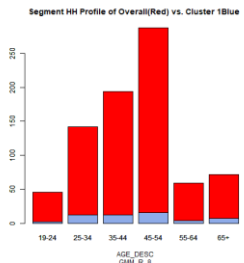


Cluster 8 (45 IDs)

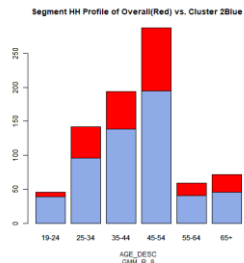


## GMM-R (8)

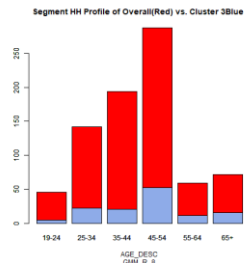
Cluster 1 (78 IDs)



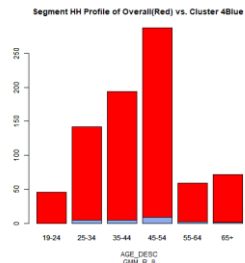
Cluster 2 (2,067 IDs)



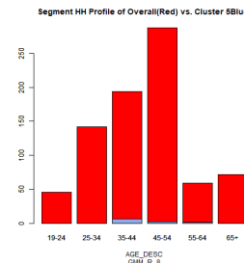
Cluster 3 (231 IDs)



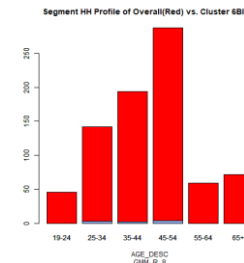
Cluster 4 (32 IDs)



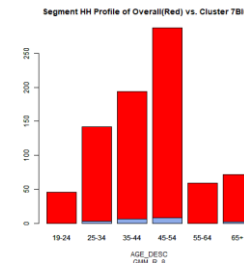
Cluster 5 (10 IDs)



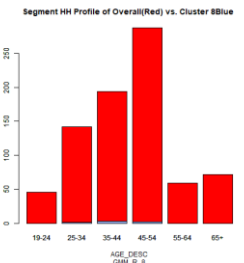
Cluster 6 (17 IDs)



Cluster 7 (26 IDs)

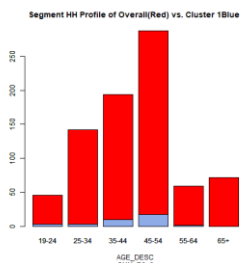


Cluster 8 (9 IDs)

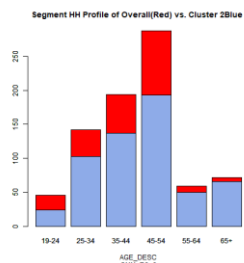


## GMM-R3 (8)

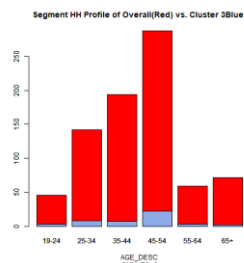
Cluster 1 (54 IDs)



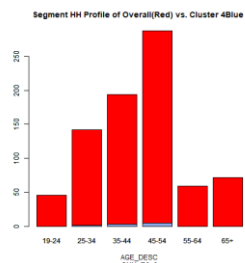
Cluster 2 (1,808 IDs)



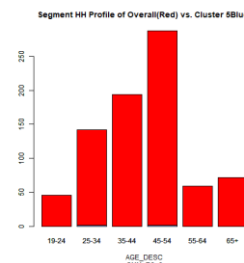
Cluster 3 (103 IDs)



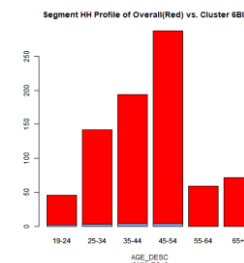
Cluster 4 (16 IDs)



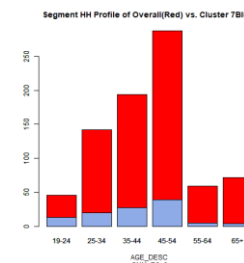
Cluster 5 (7 IDs)



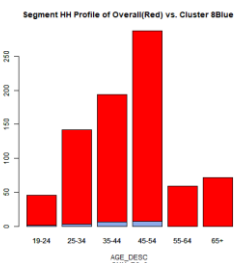
Cluster 6 (28 IDs)



Cluster 7 (429 IDs)



Cluster 8 (25 IDs)



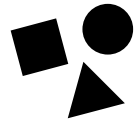




# IN SUMMARY

Data Type	Traditional k-means (TKM)	Fuzzy k-means (FKM)	Gaussian Mixture Model (GMM)
Non-Reduct Dataset	<ul style="list-style-type: none"><li>• Similar output to FKM</li><li>• Low optimal cluster counts of 2 and 3</li><li>• Even distribution</li><li>• Distinct cluster traits (transactional)</li><li>• Distinct cluster traits (demographic)</li></ul>	<ul style="list-style-type: none"><li>• Similar output to TKM</li><li>• Low optimal cluster count of 2</li><li>• Even distribution</li><li>• Distinct cluster traits (transactional)</li><li>• Distinct cluster traits (demographic)</li></ul>	<ul style="list-style-type: none"><li>• High optimal cluster counts of 8</li><li>• Even distribution</li><li>• Somewhat distinct cluster traits (transactional)</li><li>• Somewhat distinct cluster traits (demographic)</li></ul>
Reduct Dataset	<ul style="list-style-type: none"><li>• Low optimal cluster count of 2</li><li>• Uneven distribution</li><li>• Highly distinct cluster traits (transactional)</li><li>• Distinct cluster traits (demographic)</li><li>• Little difference seen between reduct-2 and reduct-3</li></ul>	<ul style="list-style-type: none"><li>• Low optimal cluster count of 2</li><li>• Uneven distribution</li><li>• Highly distinct cluster traits (transactional)</li><li>• Distinct cluster traits (demographic)</li><li>• Reduct-3 had more similar cluster traits than reduct-2</li></ul>	<ul style="list-style-type: none"><li>• High optimal cluster counts of 19</li><li>• Uneven distribution</li><li>• Similar cluster traits (transactional) for <math>k = 8</math></li><li>• Similar cluster traits (demographic)</li><li>• Little difference seen between reduct-2 and reduct-3</li></ul>

# DISCUSSION



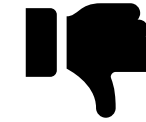
## DISTINCTIVE REDUCT DATASET

Rough Set's feature reduction generated more distinctive cluster attributes



## SLOWER SOFT CLUSTERING

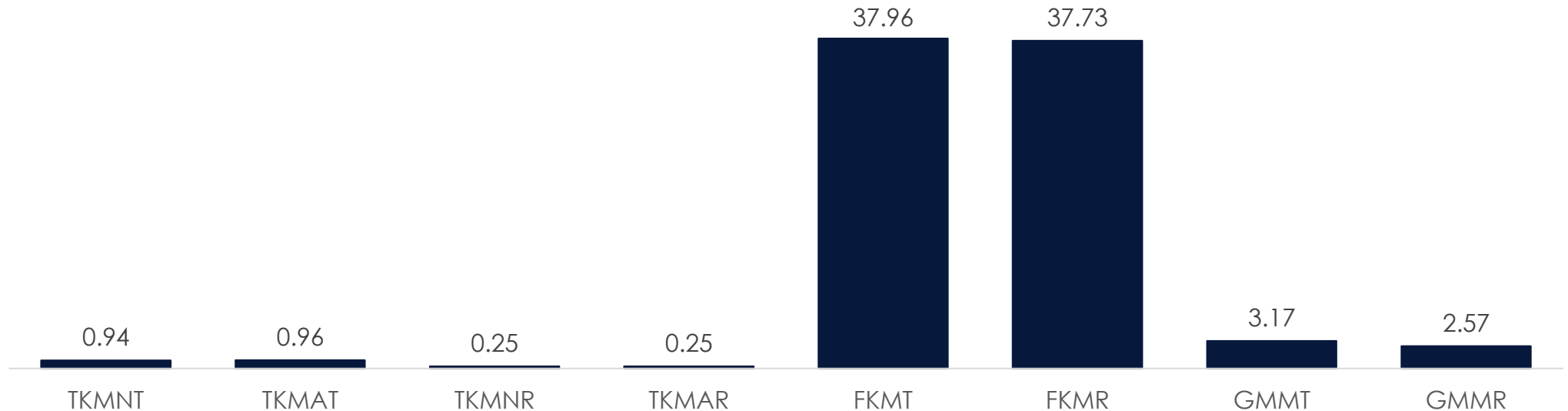
Slower processing speed observed for FKM and GMM than TKM



## HIGHER $k \neq$ BETTER & CLEAR CLUSTERS

Higher cluster counts ( $k > 2$ ) does not necessarily show distinctive clusters

Average Execution Time (Secs) of R Output with Cluster  $k=2$  with 20 iterations



\*TKM and FKM time includes silhouette calculations

# CONCLUSION & CONTRIBUTIONS

## USEFUL ADDITIONS TO A DATA ANALYST'S TOOLBOX



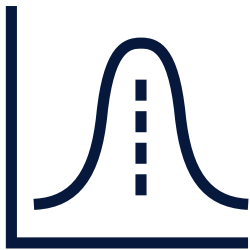
### TRADITIONAL K-MEANS

Clustering Comparison



### FUZZY K-MEANS

Clustering Comparison



### GAUSSIAN MIXTURE MODEL

Clustering Comparison



### ROUGH SET

Feature Reduction

## FUTURE WORK

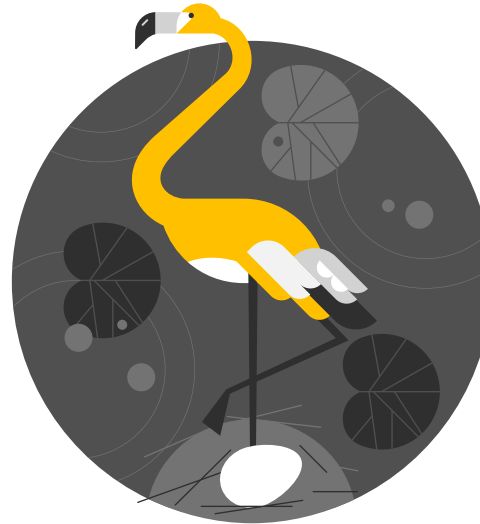
- Exploring other Rough Set reduct method
  - Capstone used global discernibility
  - Available methods include local discernibility, quick reduct.rst and quick reduct.frst
- Explore use of Rough Set reduct on both *pre* and *post* clustering
  - Capstone used reduct *pre* clustering
  - Other journals used reduct *post* clustering to distil significant attributes for each cluster

# MILESTONES & REFLECTIONS



## NEVER UNDER-ESTIMATE CAPSTONE RIGOUR

Taking additional courses would've sufficed as MITB course requirements, but capstone opportunity proved invaluable



## JOURNEY BEGINS WITH THE FIRST STEP

Initially fuzzy; Gained clarity at every step; Where clarity was at its best, scale of work seemed daunting



## SUSTAINABLE SAS & R INTEGRATION

R packages complements existing SAS environment. Challenge is to make it sustainable, amidst compatibility concerns

# THANK YOU

---

## FOR THE OPPORTUNITY!

