
ANALYSIS ON MILK SAMPLES

Presented by Anisha Mittal

CONTENTS

OUR MAIN
TOPICS TODAY

AIM

DATA SET

EDA OF TRAITS

ANALYSIS

CLUSTERING

CLASSIFICATION

CONCLUSION



AIM

Data Visualisation and exploration for the protein and technological traits of milk samples.

Clustering the cow breeds based on the MIR spectra of their milk samples.

Classifying milk samples having heat stability less than 10mins based on their MIR spectra.

DATA SET

ABOUT THE DATA

Initial Data set

582 variables and 431 observations.

After removing NA columns

574 variables and 430 observations.

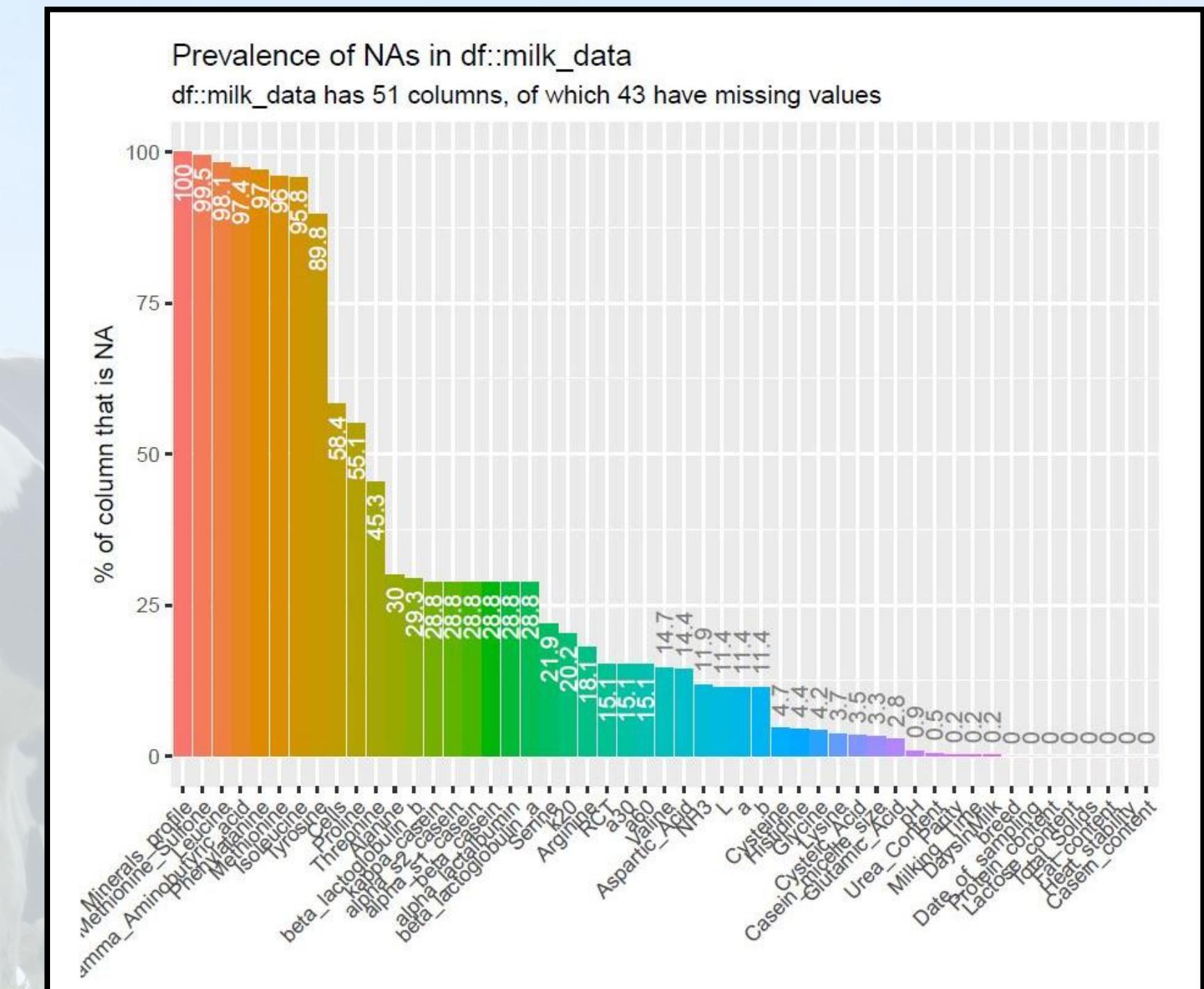
After removing outliers

574 variables and 326 observations.

Final Data set

1:43 trait variables

44:531 MIR spectra variables



EXPLORATORY DATA ANALYSIS

CASEIN

Kappa Casein
Beta Casein
Alpha s1 Casein
Alpha s2 Casein

WHEY

Alpha Lactalbumin
Beta Lactalbumin A
Beta Lactalbumin B

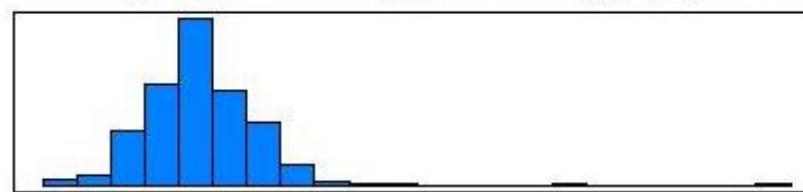
TECHNOLOGICAL

Casein Micelle size
Heat Stability
pH
RCT
A30
A60
K20

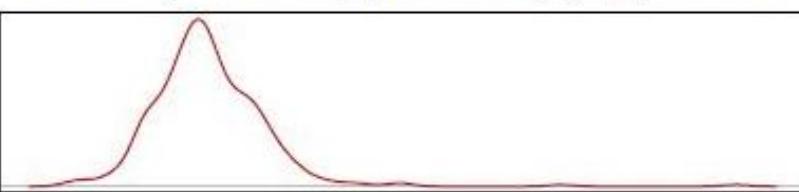
CASEIN TRAITS

EXPLORATORY DATA ANALYSIS

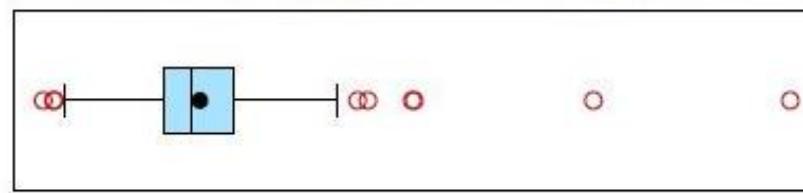
Histogram of milk_data\$kappa_casein



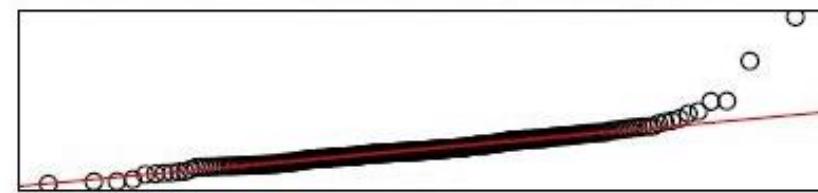
Density of milk_data\$kappa_casein



Boxplot of milk_data\$kappa_casein



Q-Q Plot of milk_data\$kappa_casein

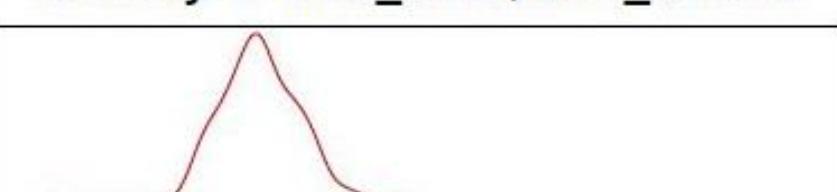


EXPLORATORY DATA ANALYSIS

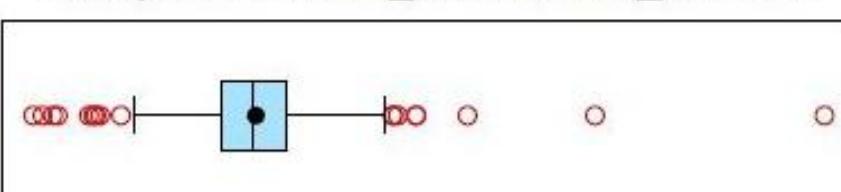
Histogram of milk_data\$beta_casein



Density of milk_data\$beta_casein



Boxplot of milk_data\$beta_casein

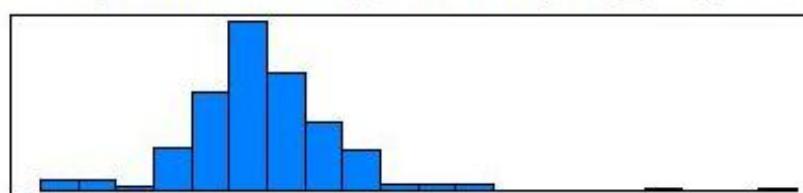


Q-Q Plot of milk_data\$beta_casein

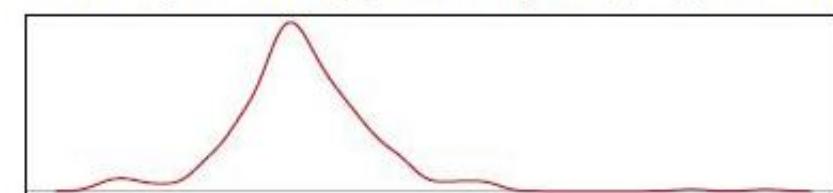


EXPLORATORY DATA ANALYSIS

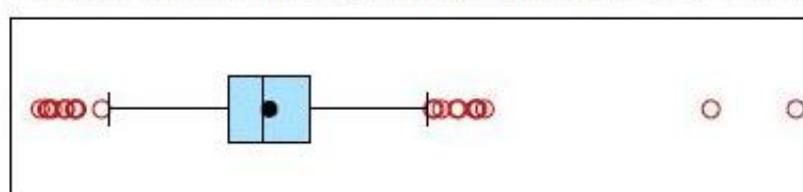
Histogram of milk_data\$alpha_s2_casein



Density of milk_data\$alpha_s2_casein



Boxplot of milk_data\$alpha_s2_casein

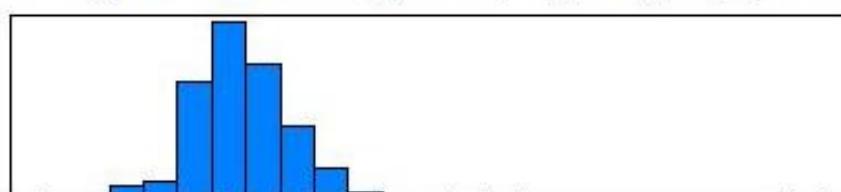


Q-Q Plot of milk_data\$alpha_s2_casein

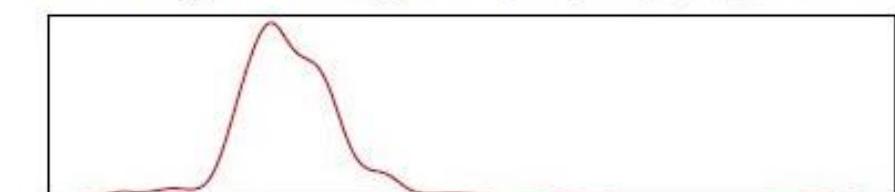


EXPLORATORY DATA ANALYSIS

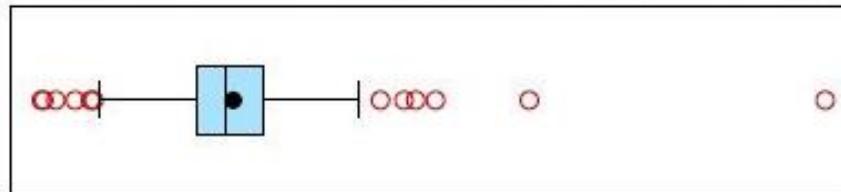
Histogram of milk_data\$alpha_s1_casein



Density of milk_data\$alpha_s1_casein



Boxplot of milk_data\$alpha_s1_casein



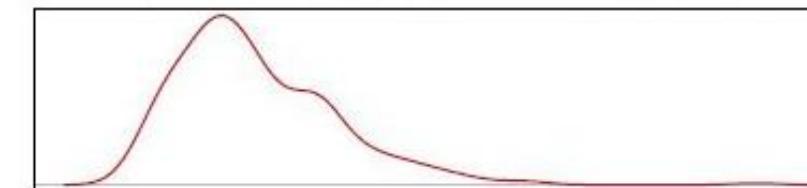
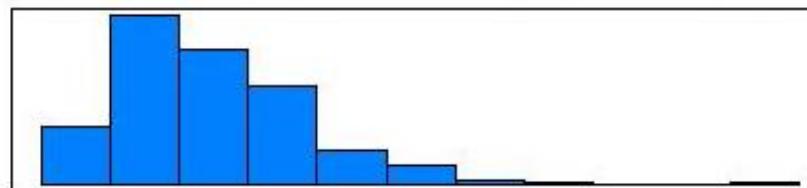
Q-Q Plot of milk_data\$alpha_s1_casein



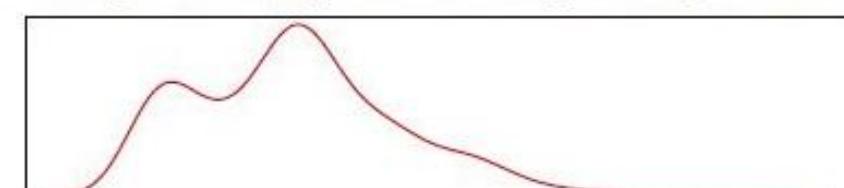
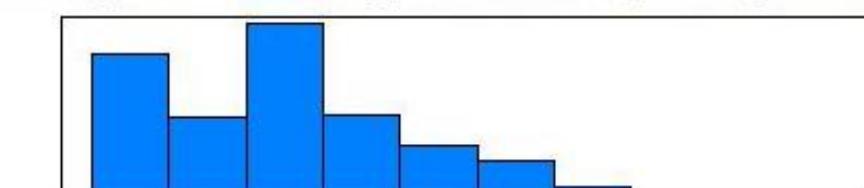
WHEY TRAITS

EXPLORATORY DATA ANALYSIS

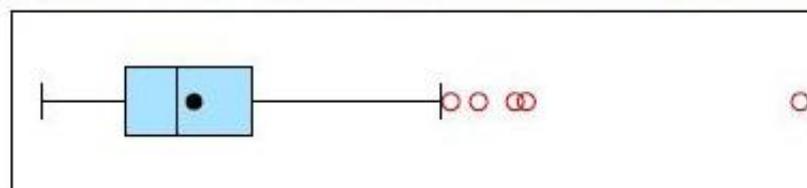
histogram of milk_data\$beta_lactoglobulin_Density of milk_data\$beta_lactoglobulin_a



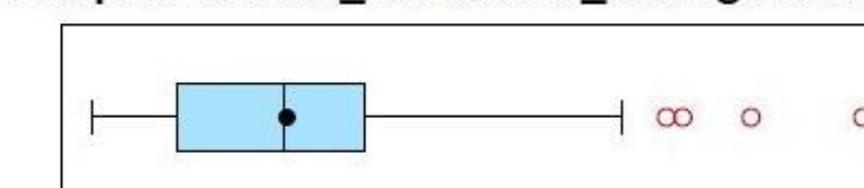
histogram of milk_data\$beta_lactoglobulin_Density of milk_data\$beta_lactoglobulin_b



Boxplot of milk_data\$beta_lactoglobulin_a



Boxplot of milk_data\$beta_lactoglobulin_b

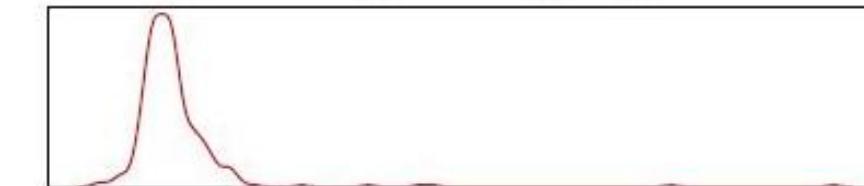
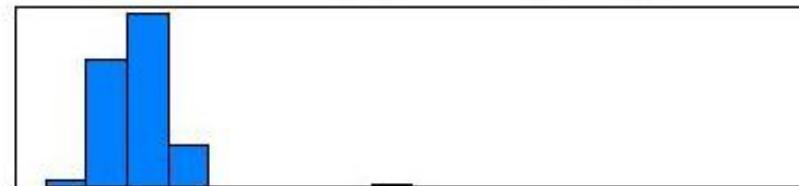


Variable Alpha
Lactalbumin shows
some variability and is
right-skewed but mean
and median are equal.

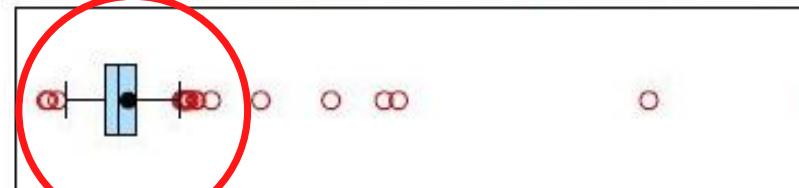


EXPLORATORY DATA ANALYSIS

Histogram of milk_data\$alpha_lactalbumin Density of milk_data\$alpha_lactalbumin



Boxplot of milk_data\$alpha_lactalbumin



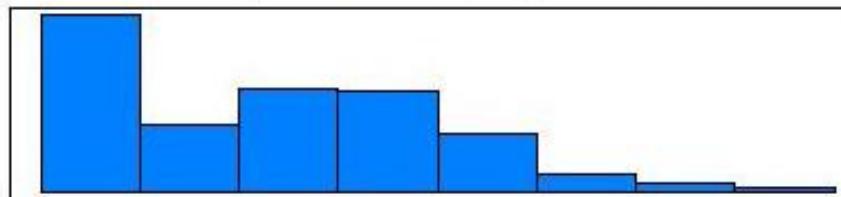
Q-Q Plot of milk_data\$alpha_lactalbumin



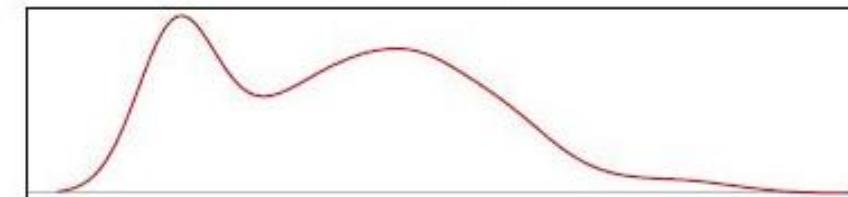
TECHNOLOGICAL TRAITS

EXPLORATORY DATA ANALYSIS

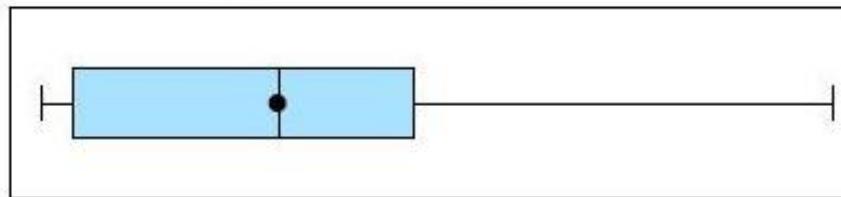
Histogram of milk_data\$a30



Density of milk_data\$a30



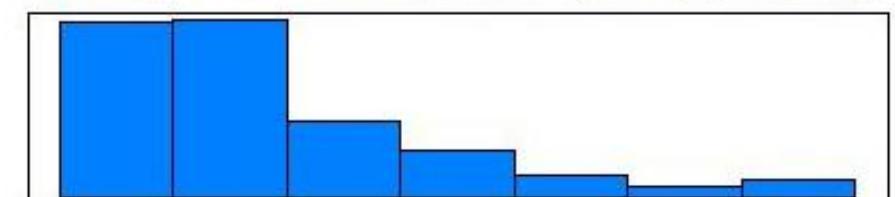
Boxplot of milk_data\$a30



Q-Q Plot of milk_data\$a30

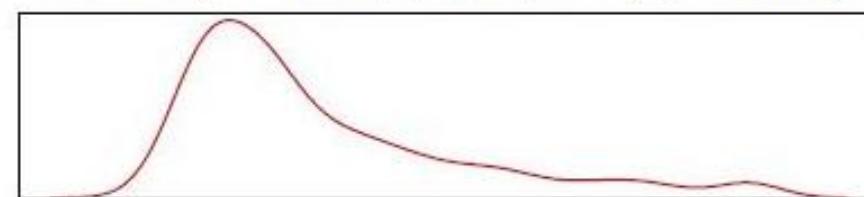


Histogram of milk_data\$Heat_stability

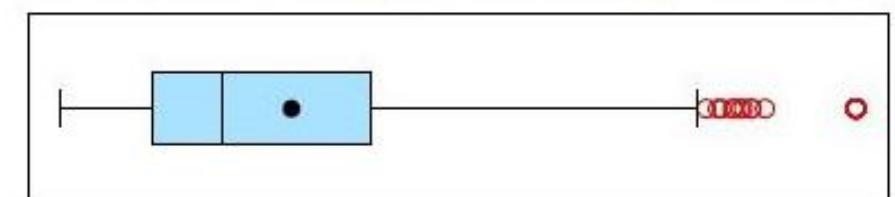


EXPLORATORY DATA ANALYSIS

Density of milk_data\$Heat_stability



Boxplot of milk_data\$Heat_stability

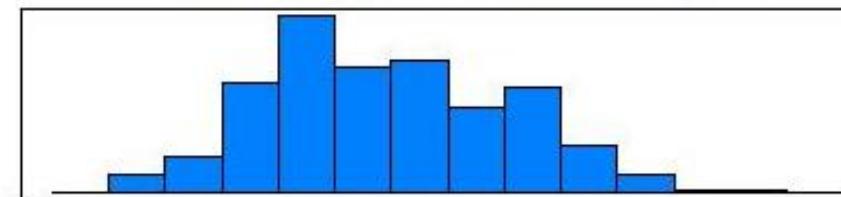


Q-Q Plot of milk_data\$Heat_stability

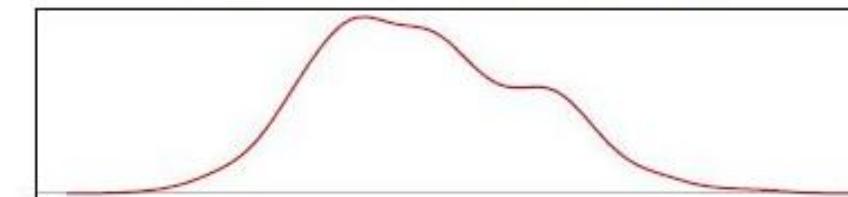


EXPLORATORY DATA ANALYSIS

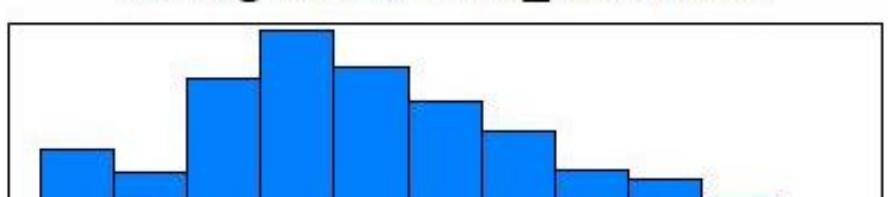
Histogram of milk_data\$pH



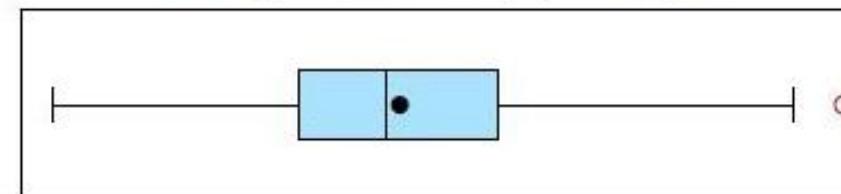
Density of milk_data\$pH



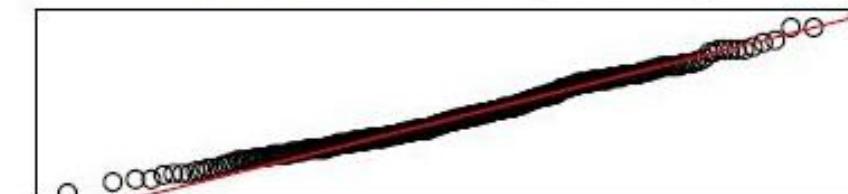
Histogram of milk_data\$RCT



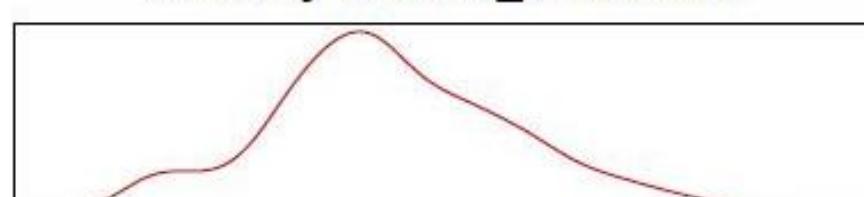
Boxplot of milk_data\$pH



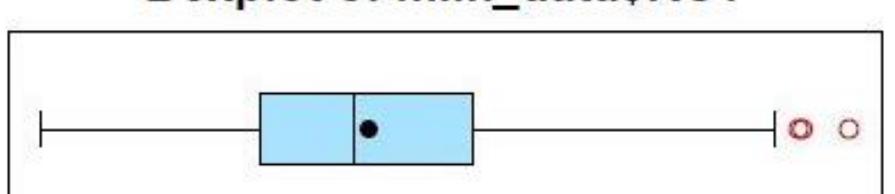
Q-Q Plot of milk_data\$pH



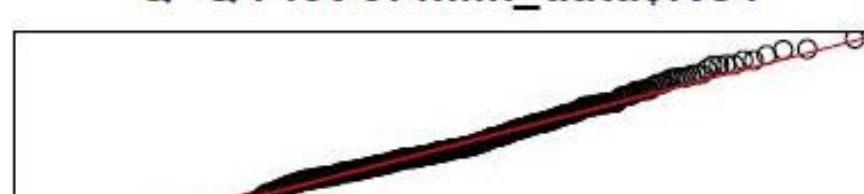
Density of milk_data\$RCT



Boxplot of milk_data\$RCT



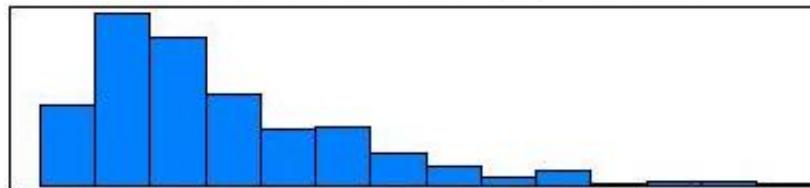
Q-Q Plot of milk_data\$RCT



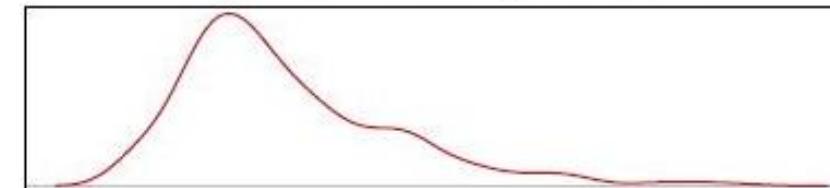
TECHNOLOGICAL TRAITS

EXPLORATORY DATA ANALYSIS

Histogram of milk_data\$k20



Density of milk_data\$k20



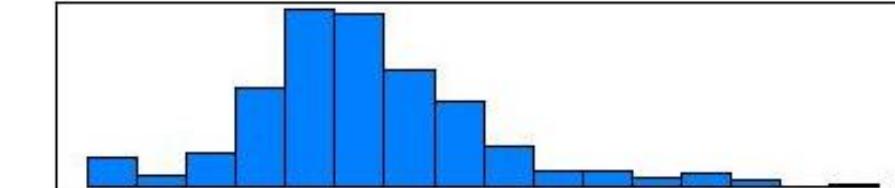
Boxplot of milk_data\$k20



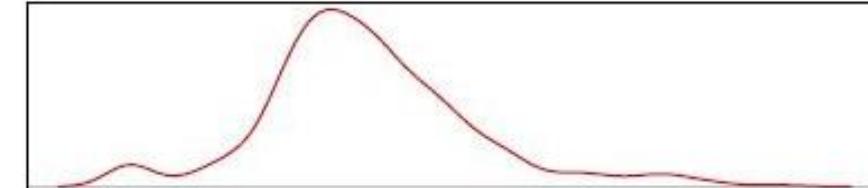
Q-Q Plot of milk_data\$k20



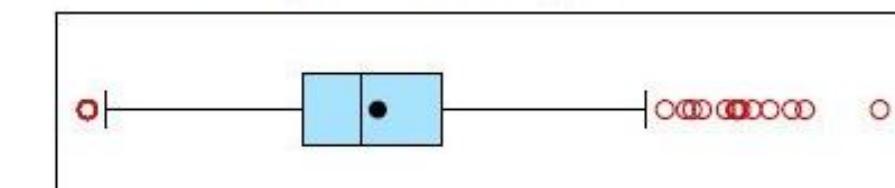
Histogram of milk_data\$a60



Density of milk_data\$a60



Boxplot of milk_data\$a60

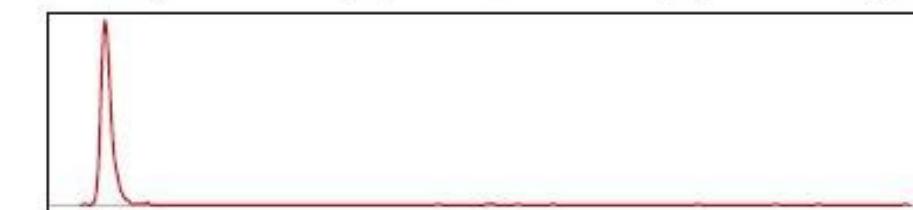


Q-Q Plot of milk_data\$a60

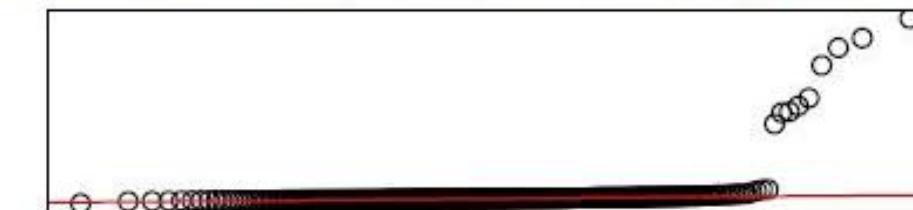
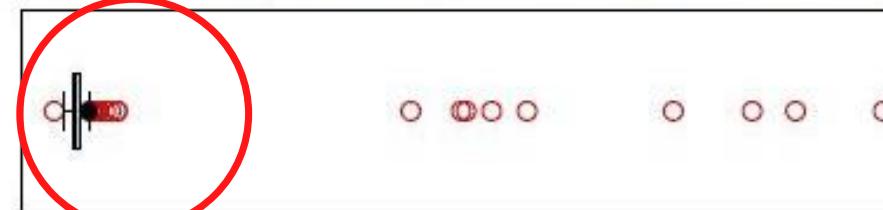


EXPLORATORY DATA ANALYSIS

Histogram of milk_data\$Casein_micelle_size



Boxplot of milk_data\$Casein_micelle_size



Variable Casein Micelle size shows a lot of variability and is right-skewed.



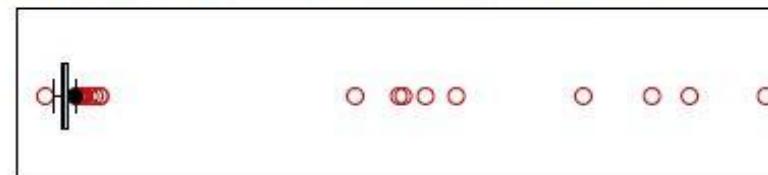
Casein Micelle Size

EXPLORATORY DATA ANALYSIS

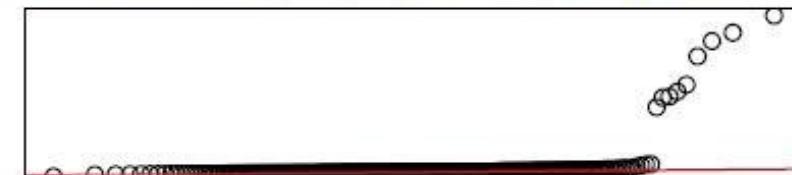
histogram of milk_data\$Casein_micelle_size
Density of milk_data\$Casein_micelle_size



Boxplot of milk_data\$Casein_micelle_size



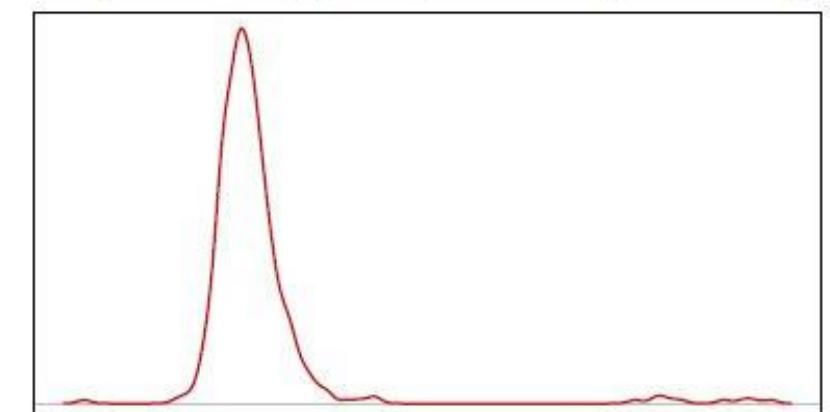
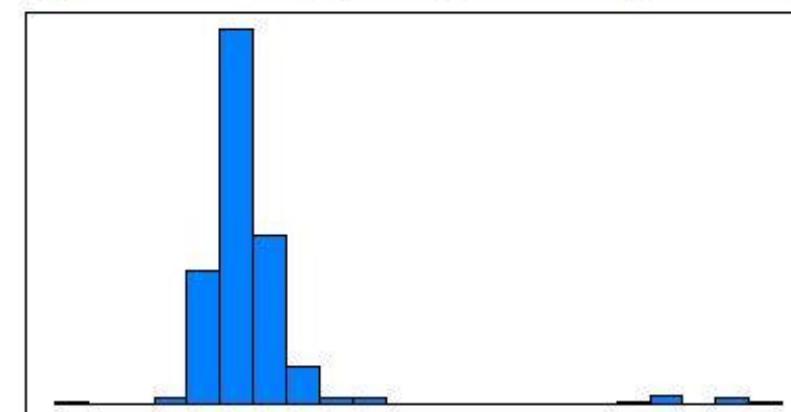
Q-Q Plot of milk_data\$Casein_micelle_size



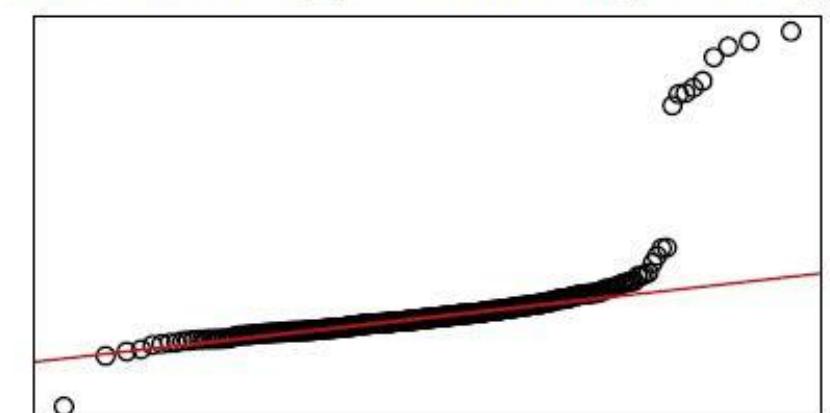
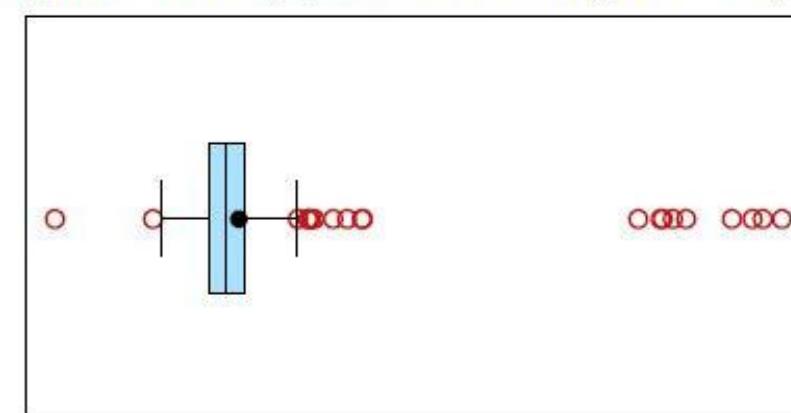
EDA of variable Casein
Micelle size after log
Transformation.

EXPLORATORY DATA ANALYSIS

histogram of milk_data\$Casein_micelle_size
Density of milk_data\$Casein_micelle_size

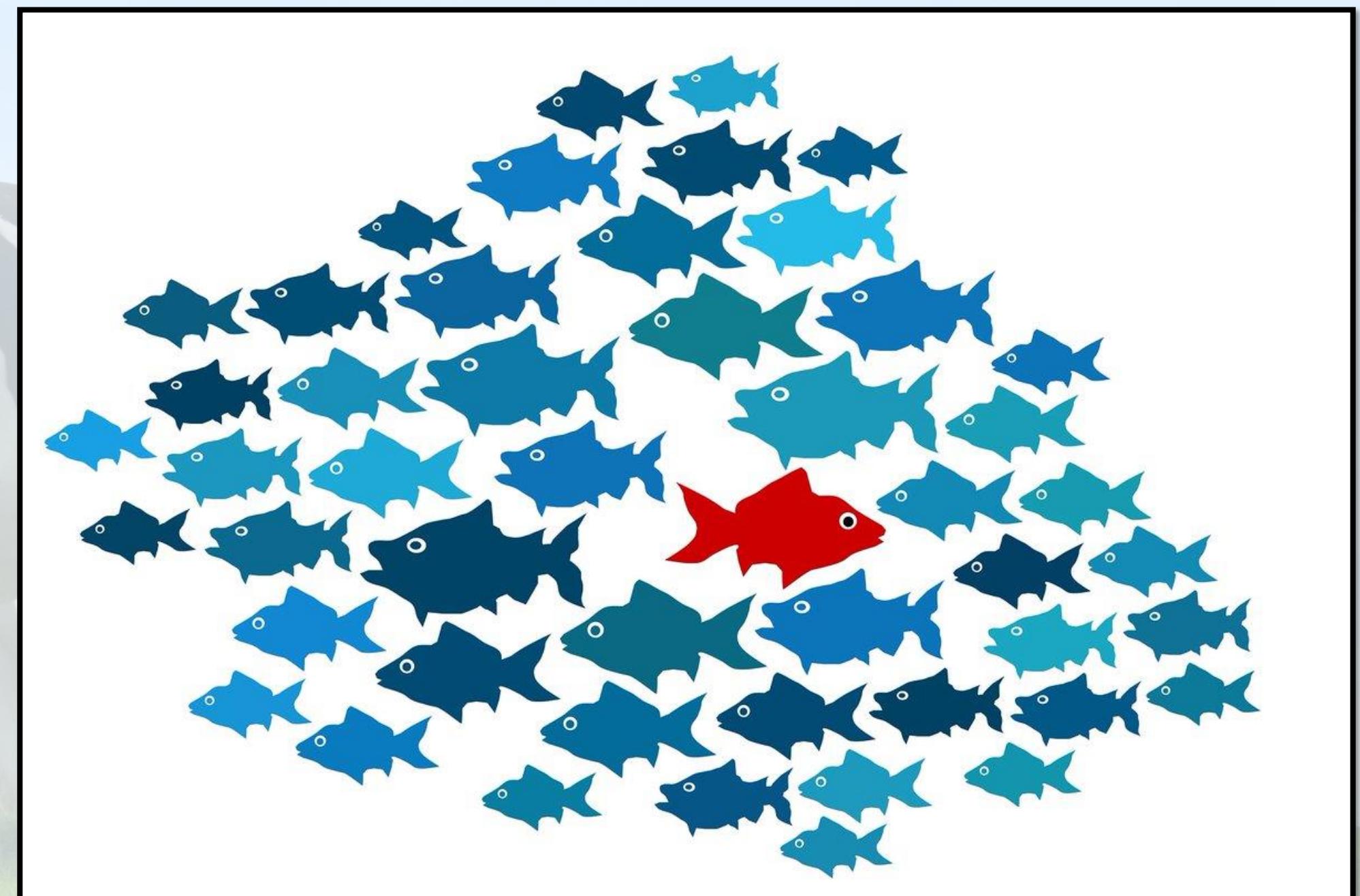


Boxplot of milk_data\$Casein_micelle_size

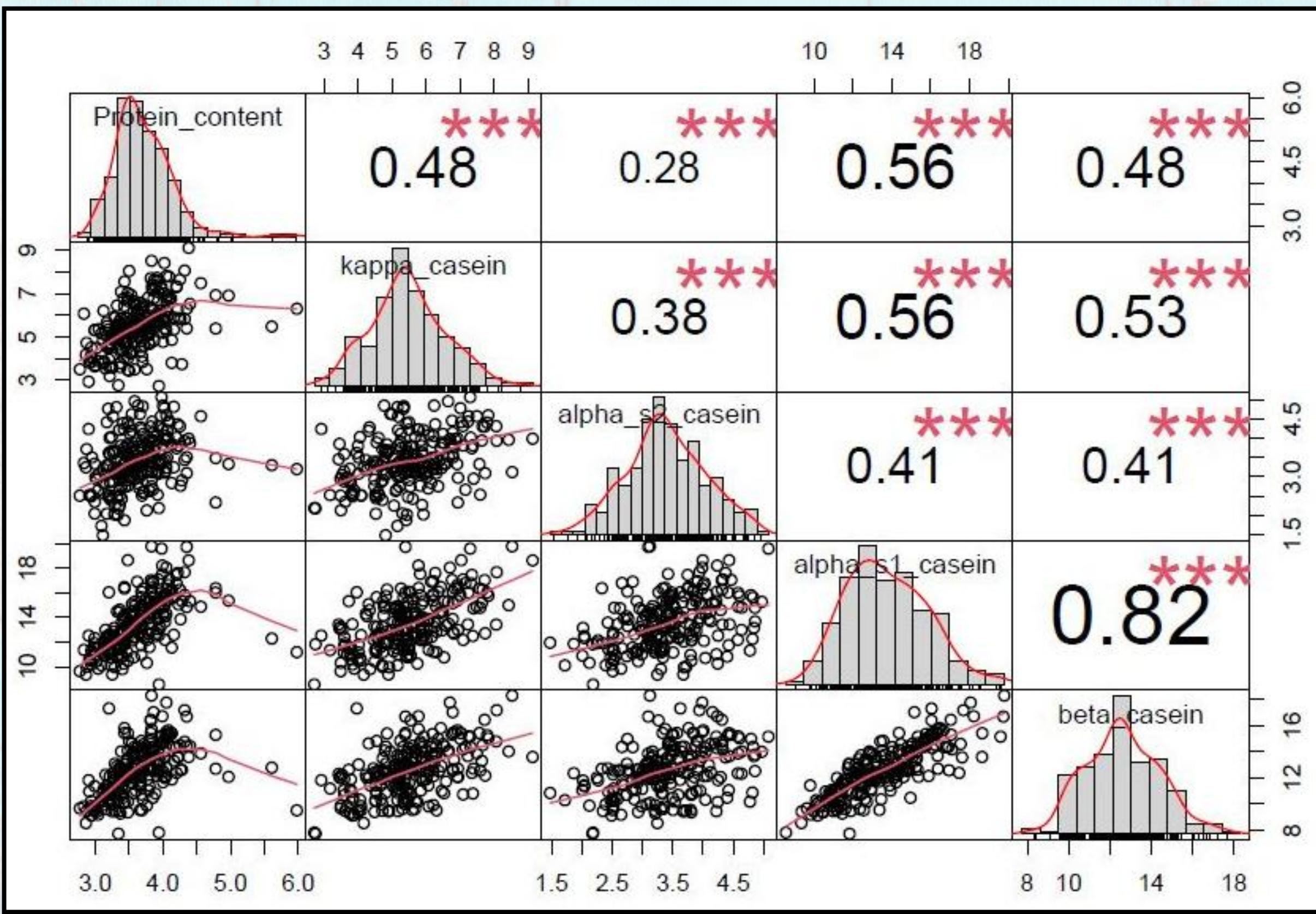


OUTLIER REMOVAL

Outliers are unusual values in a dataset, and they can distort statistical analyses and violate their assumptions.



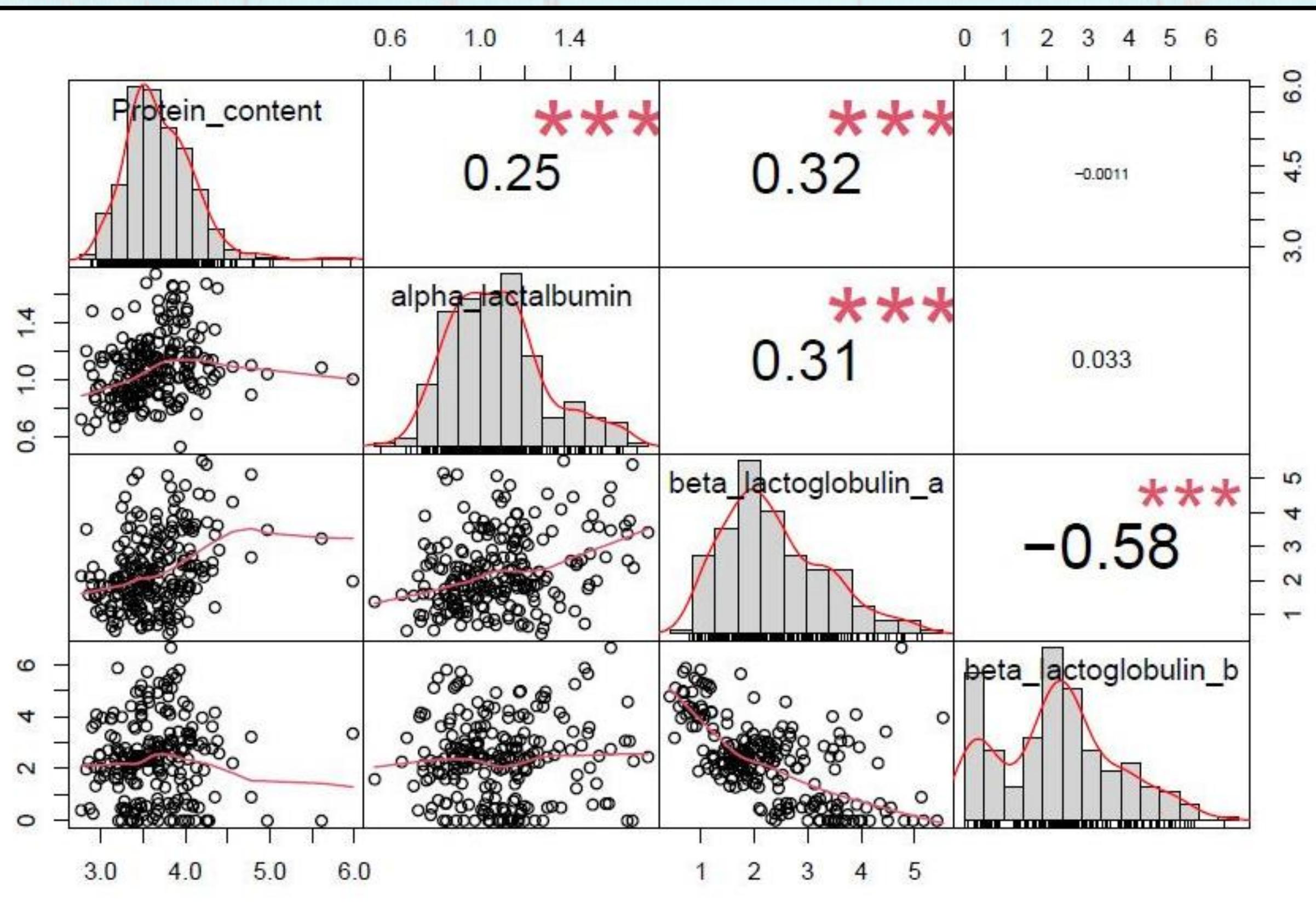
CORRELATION CHART



CASEIN PROTEIN

We can see that all the Casein protein traits are highly positively correlated to each other.

CORRELATION CHART

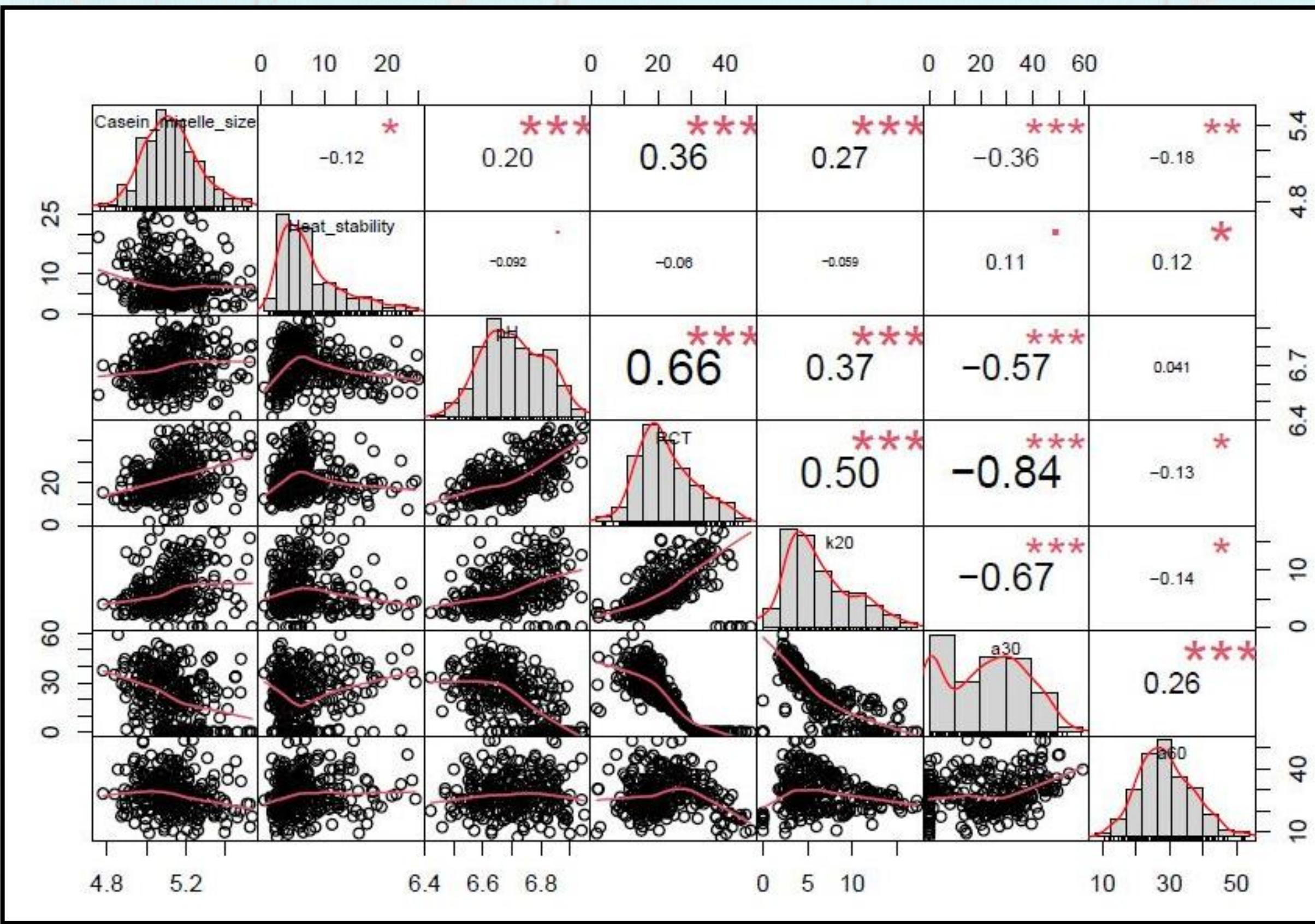


WHEY PROTEIN

Beta-lactalbumin B and Alpha-lactalbumin traits are not significantly correlated.

Beta-lactalbumin B is negatively correlated to protein content in a milk sample.

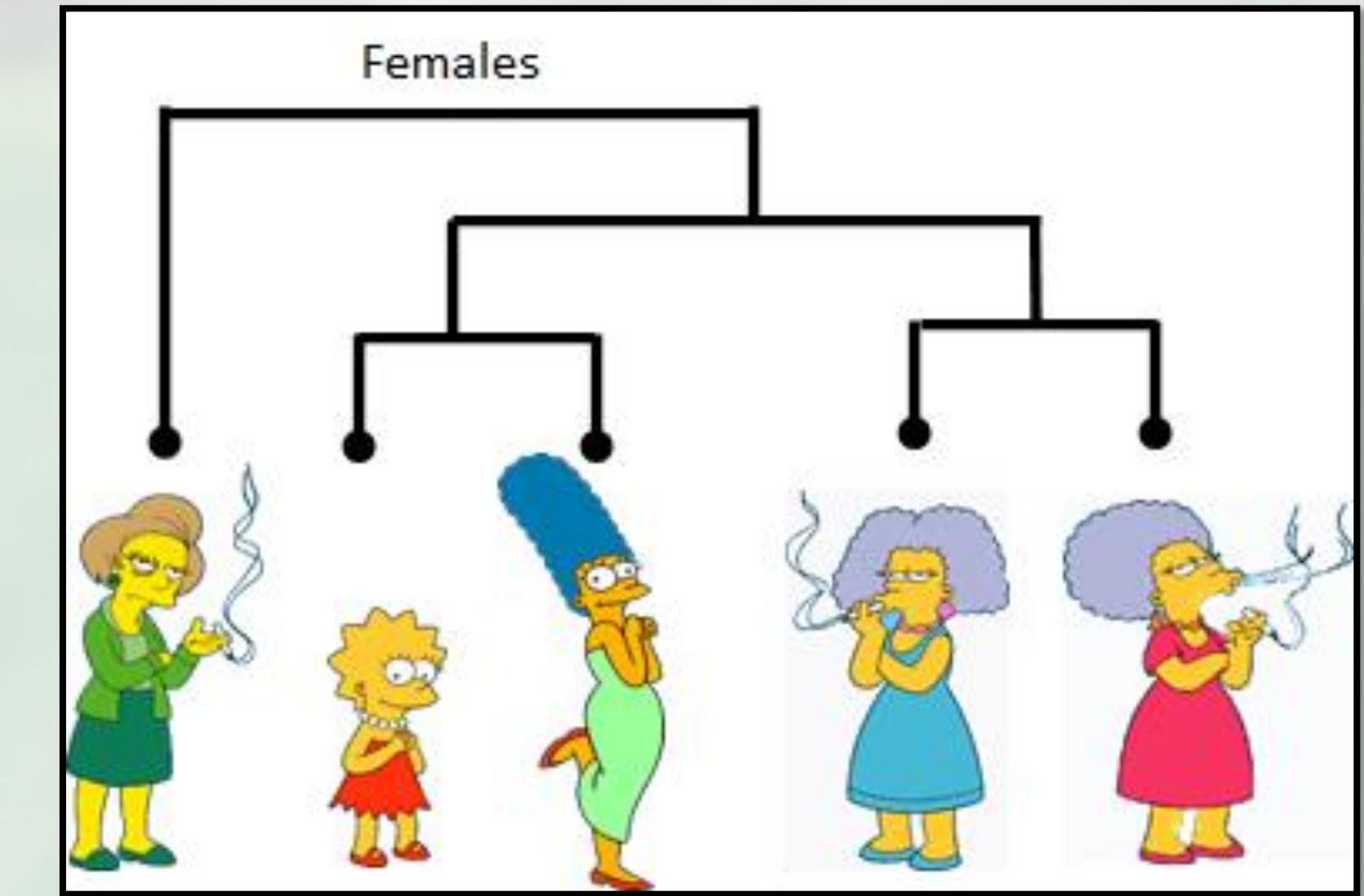
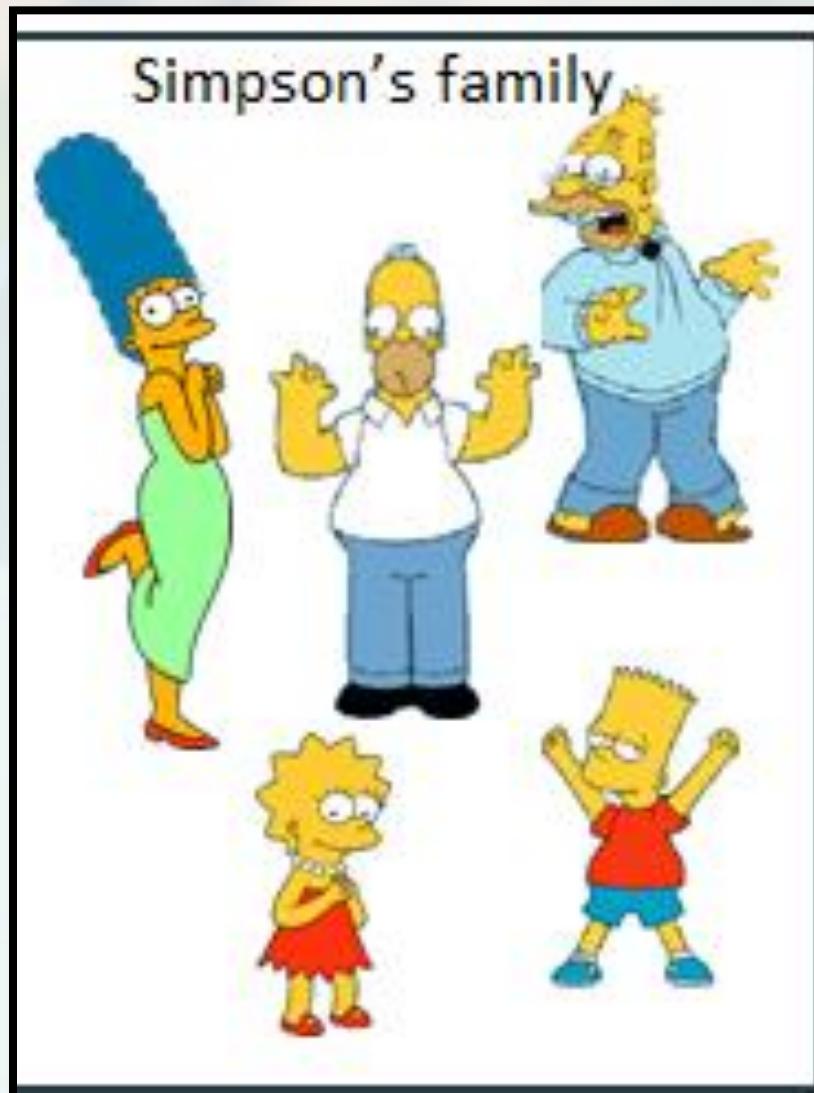
CORRELATION CHART



TECHNOLOGICAL

We can see that there are very few technological traits that are correlated to each other.

CLUSTERING

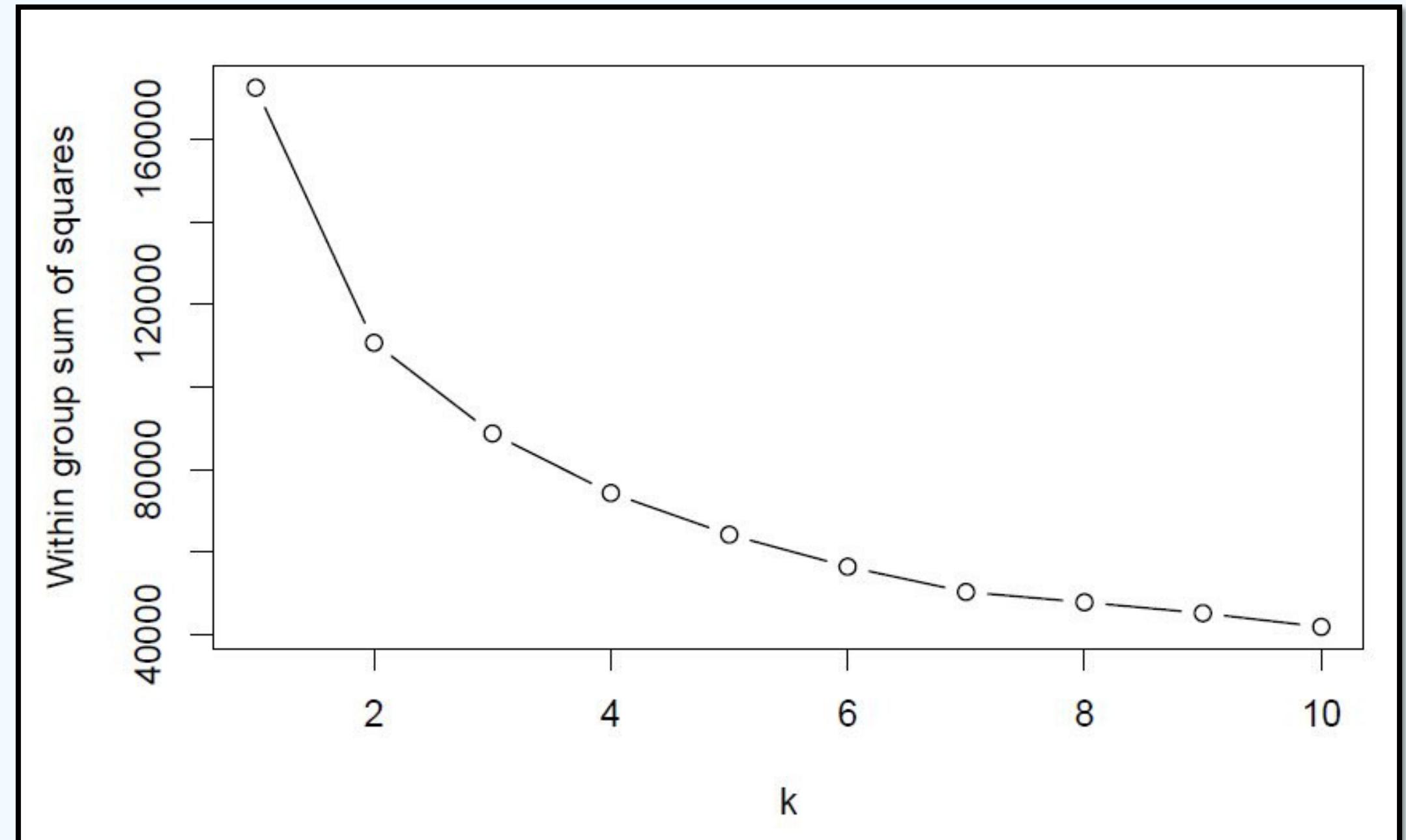


PARTITIONING

HIERARCHICAL

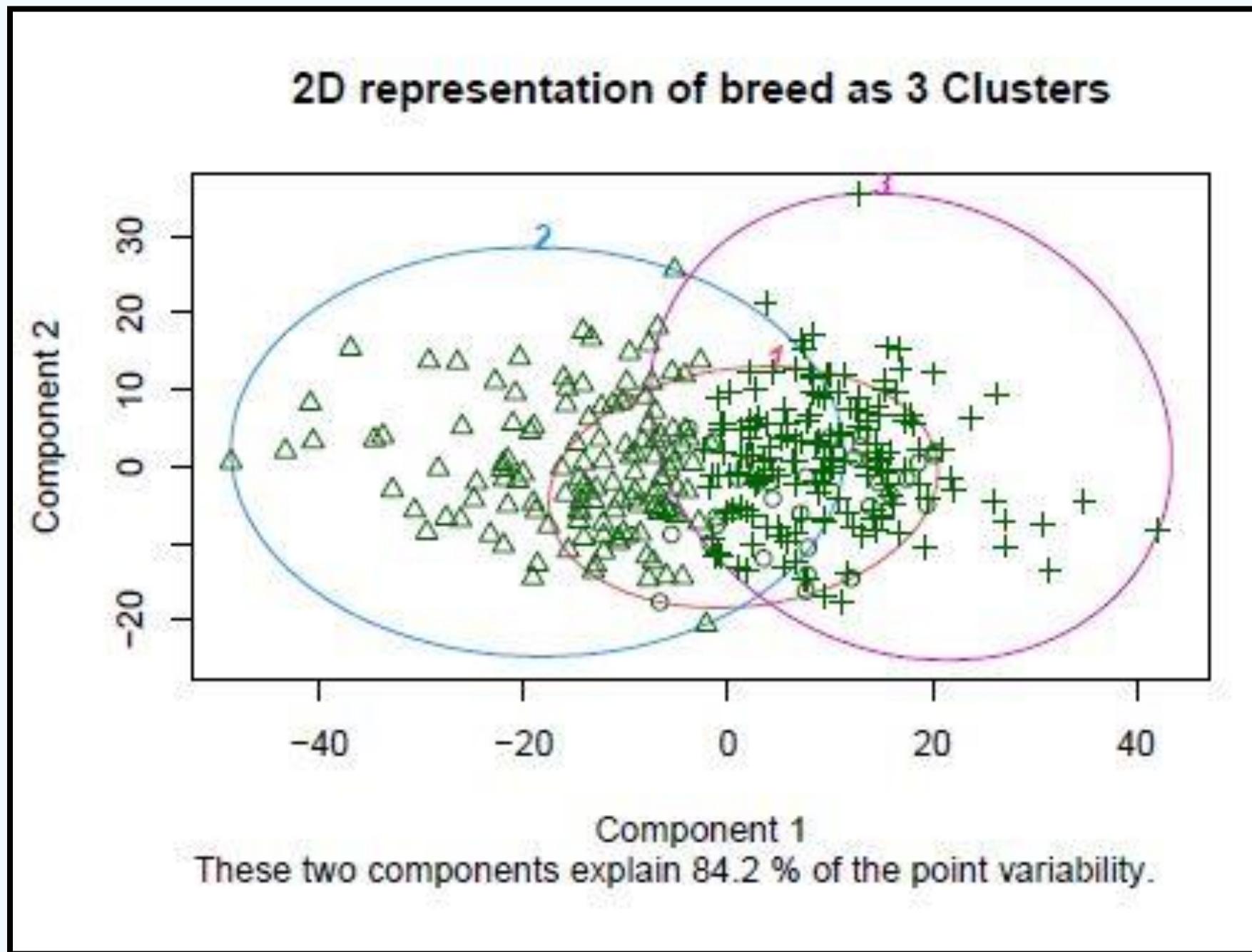
WHAT K TO SELECT? WITHIN SUM OF SQUARE

Plot a range of k values against within sum of squares and select the K with the least sum of squares value.

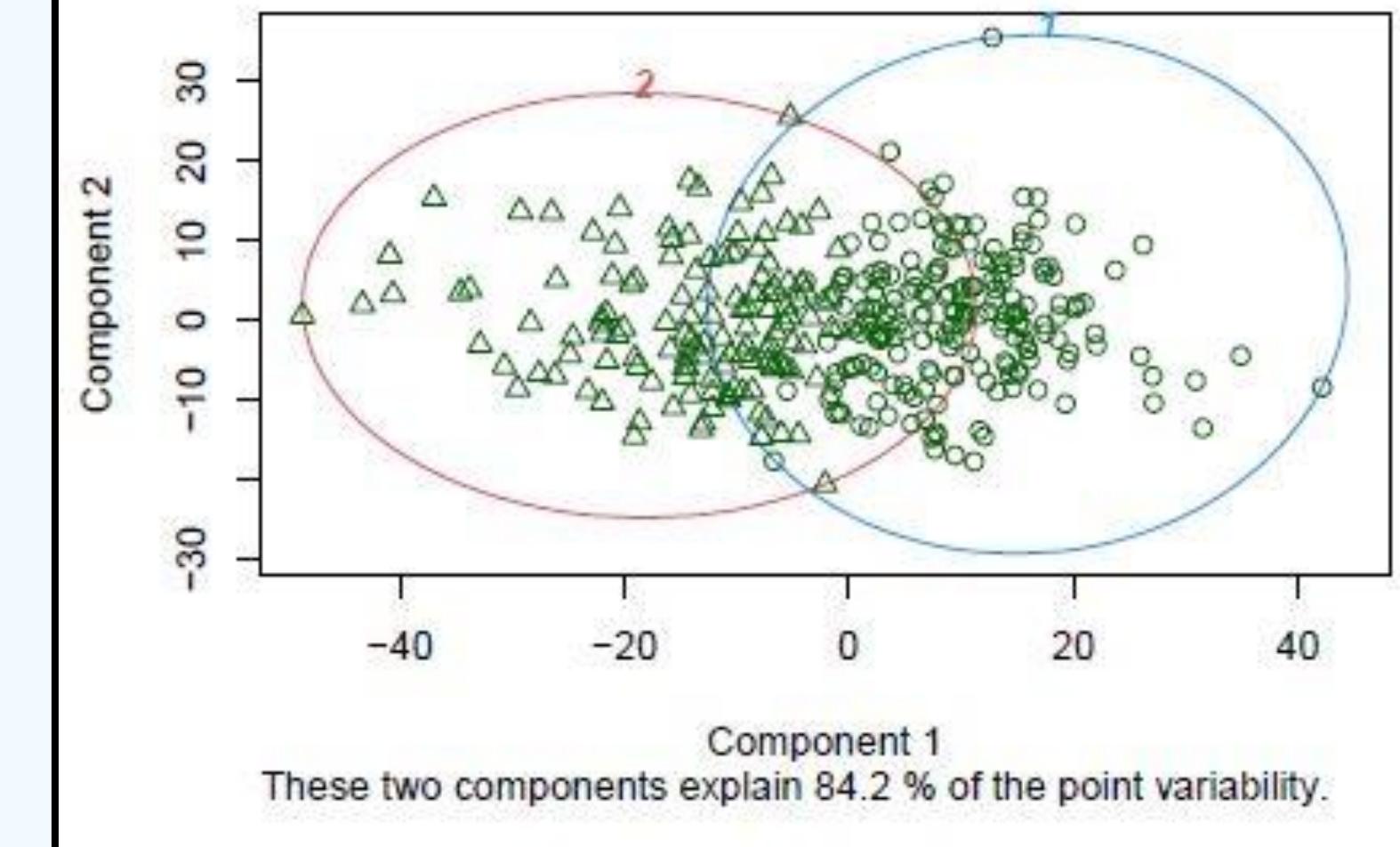


K=2

Cluster 1 = 195 observations
Cluster 2 = 131 observations



2D representation of breed as 2 Clusters



K=3

Cluster 1 = 36 observations
Cluster 2 = 126 observations
Cluster 3 = 164 observations

VALIDATION

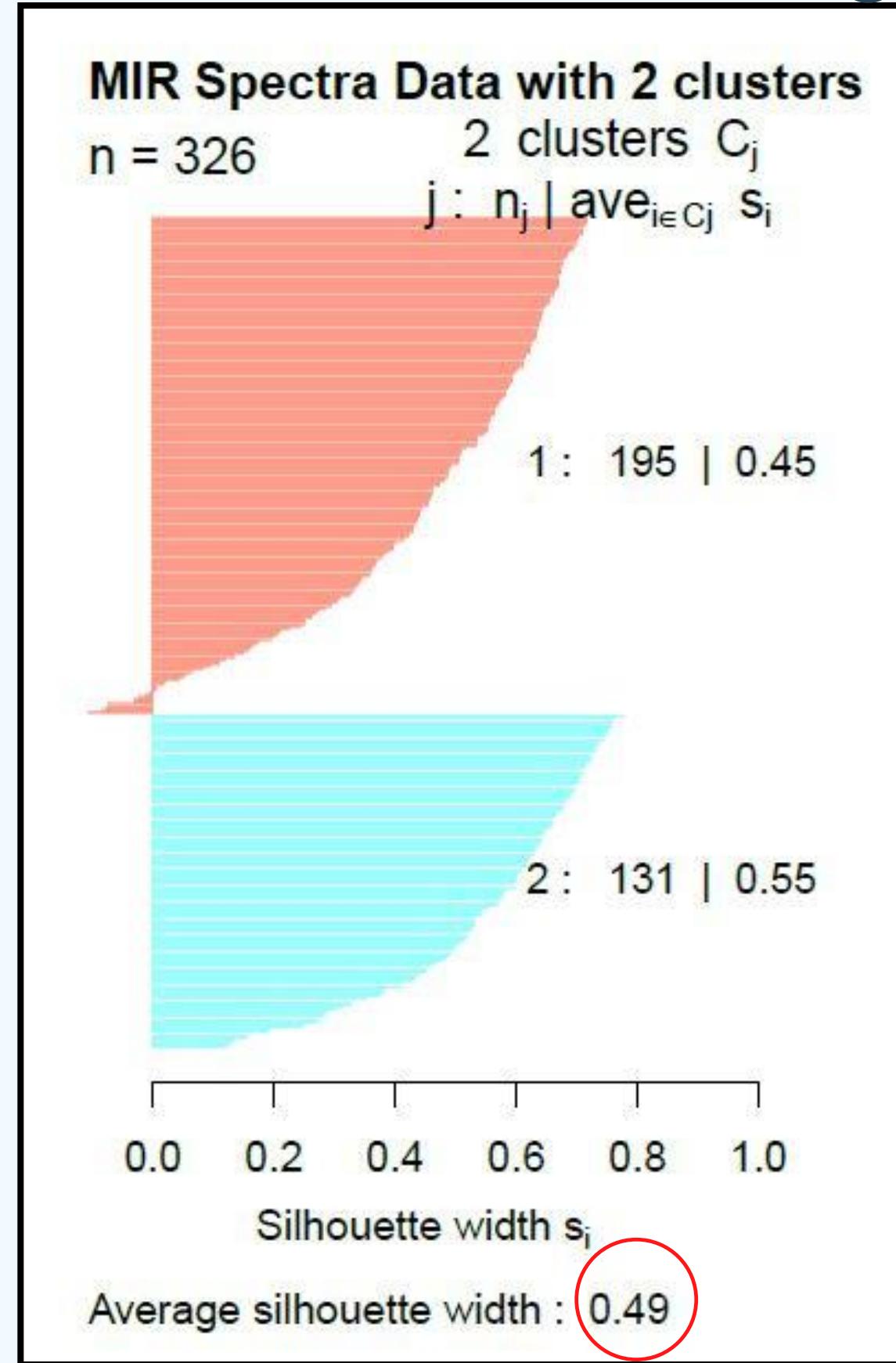
INTERNAL

Useful to check whether the appropriate K is used or not. eg: Average Silhouette width, and Calinski-Harabasz Index

EXTERNAL

Useful to assess the quality of the given clustering in comparison to a reference grouping of the observations.
eg: Rand Index, and Adjusted Rand Index

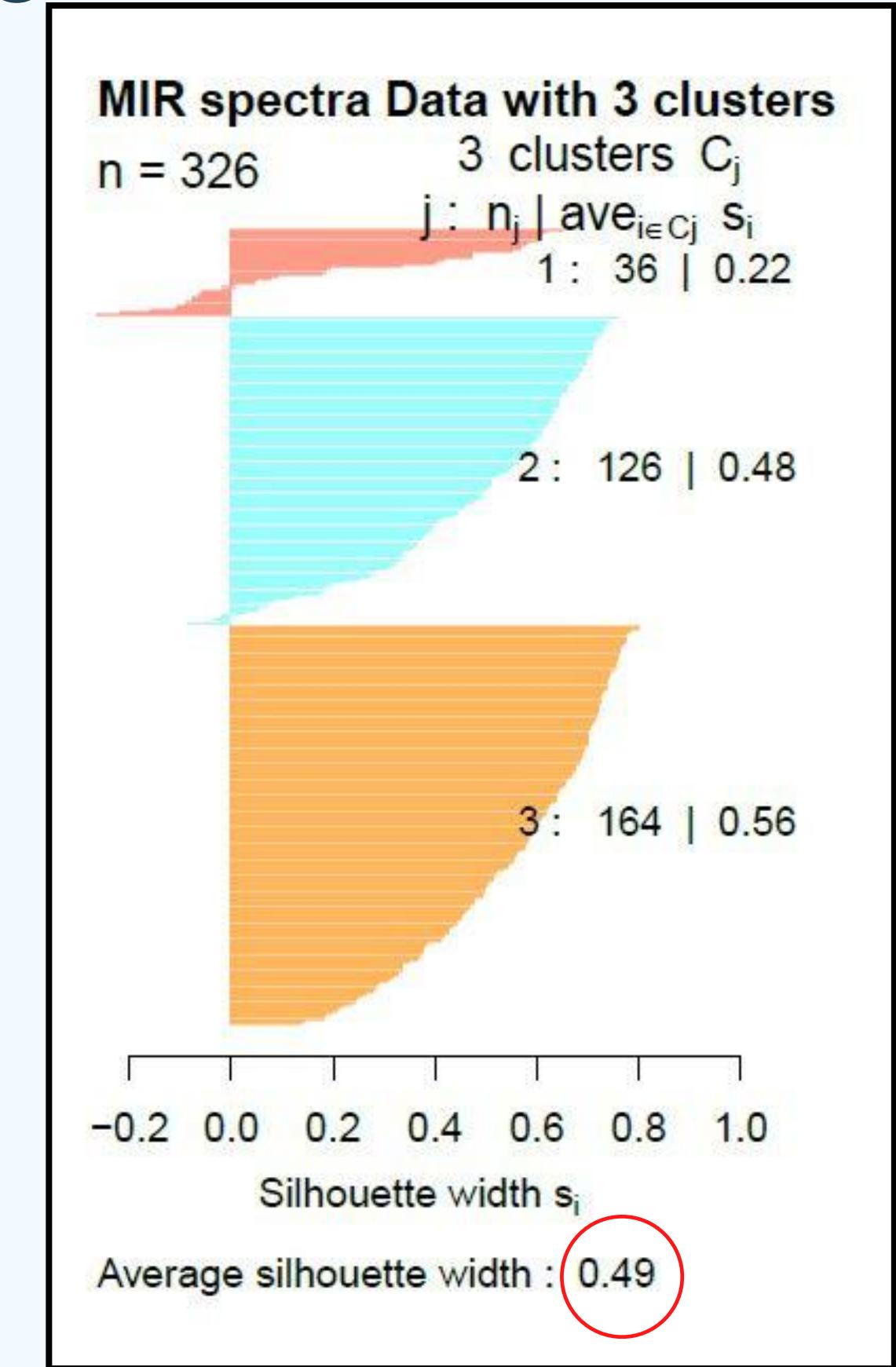
SILHOUETTE PLOTS



EUCLIDEAN

Formula :
$$\sqrt{\sum_{i=1}^p (x_{ij} - x_{ji})^2}$$

Calculates the square root
of sum of squared
differences between two
observations one by one.

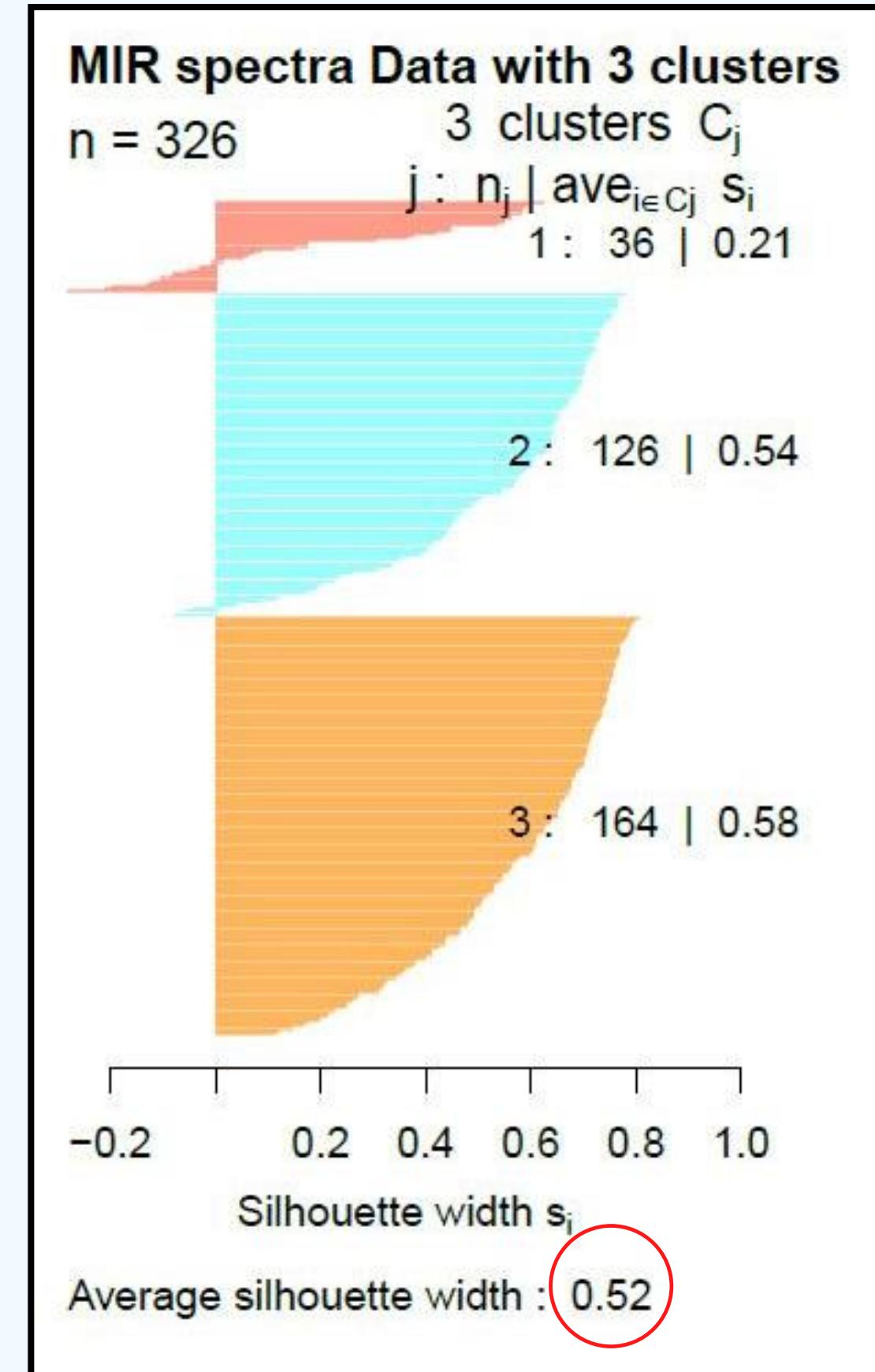
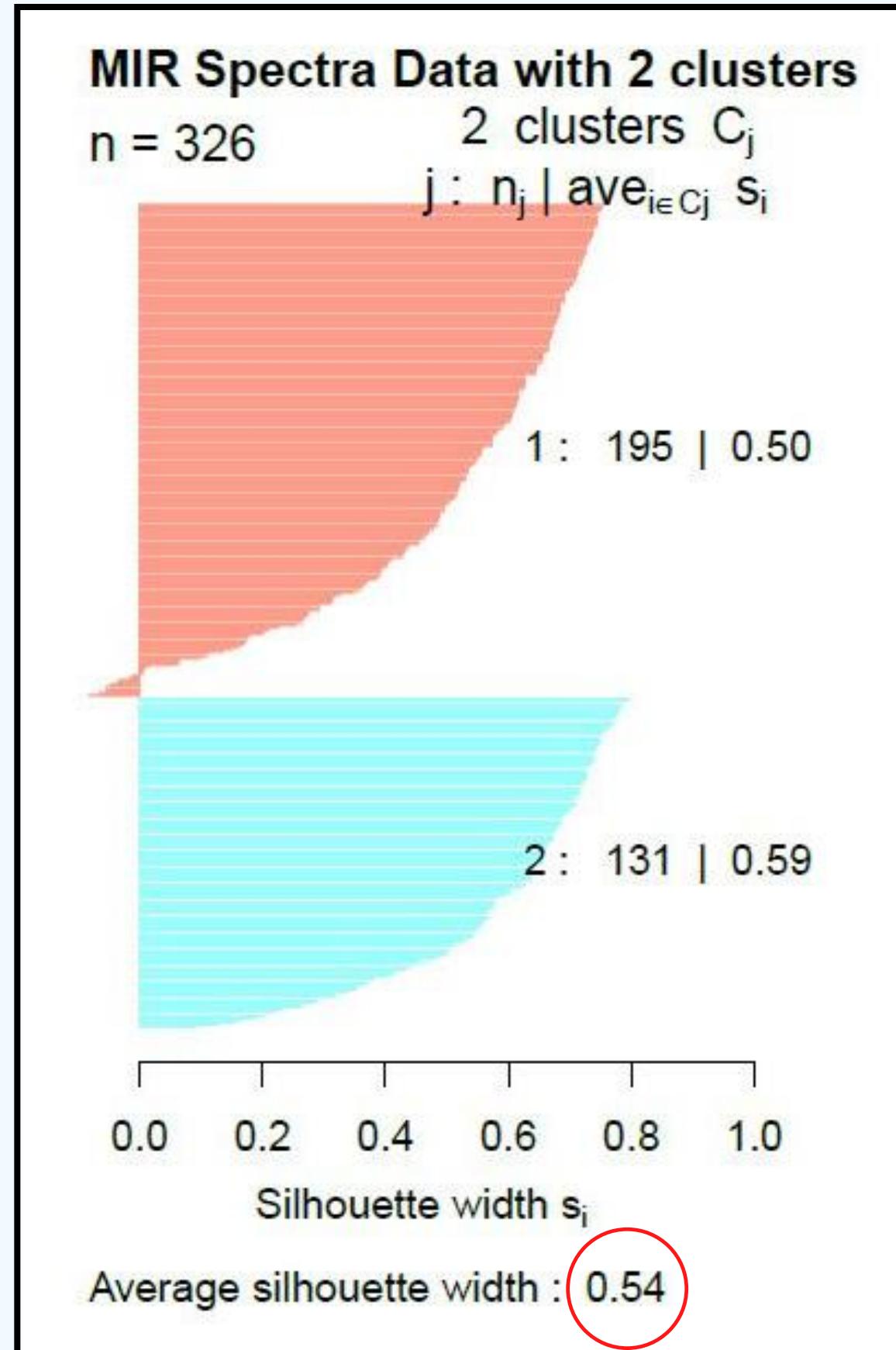


SILHOUETTE PLOTS

MANHATTAN

Formula: $\sum_{i=1}^p |x_{ij} - x_{ji}|$

Calculates the sum of absolute differences between two observations one by one.



CLASS AGREEMENT

A higher percentage means a better fit.

	Rand Index	Adjusted Rand Index
$k = 2$	52%	5%
$k = 3$	51%	1%

RESULT

No. of cows of different breeds in similar clusters.

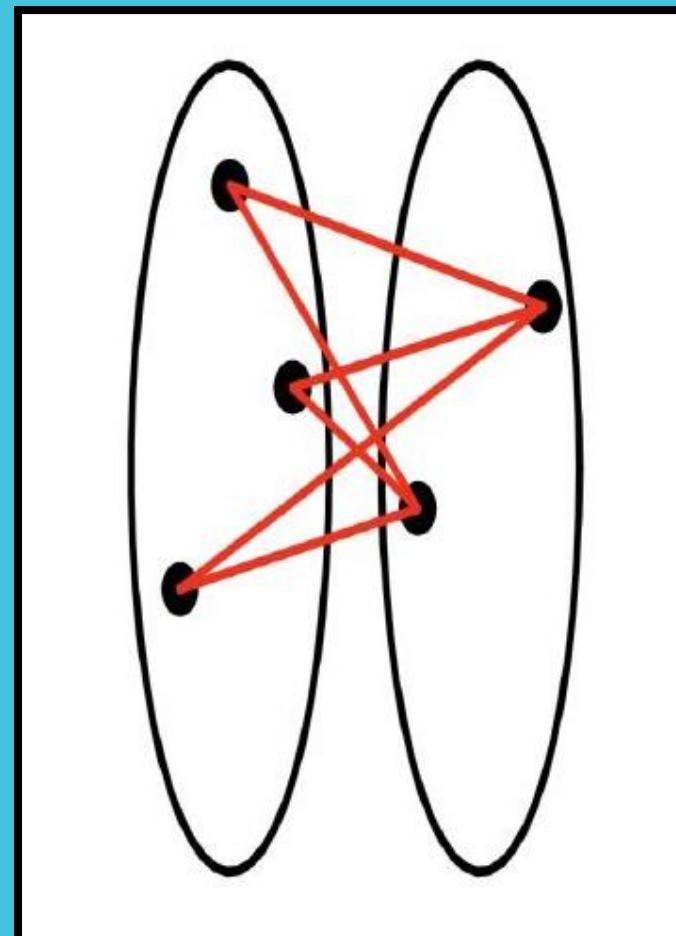
	FRX-	FRX-	HOX	HOL FRI	HOX	HOX-	JE	JEX-	MO	NR
1	0	9	1	144	10	1	18	4	1	6
2	1	4	0	76	1	2	28	14	0	6

1 observation did not have a breed name and was put in cluster 1.

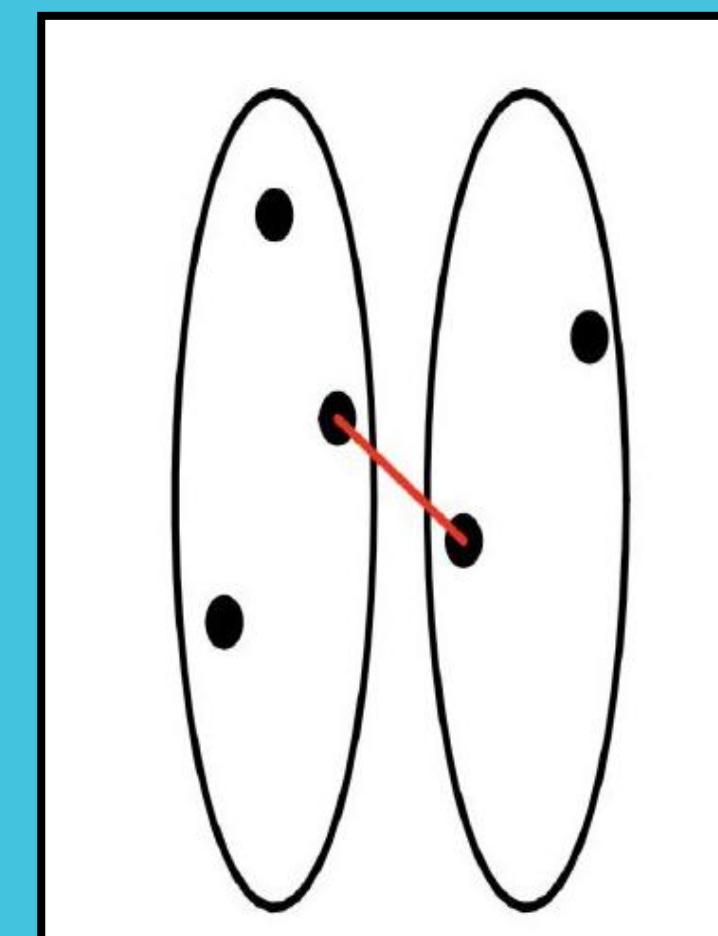
HIERARCHICAL CLUSTERING

LINKAGE

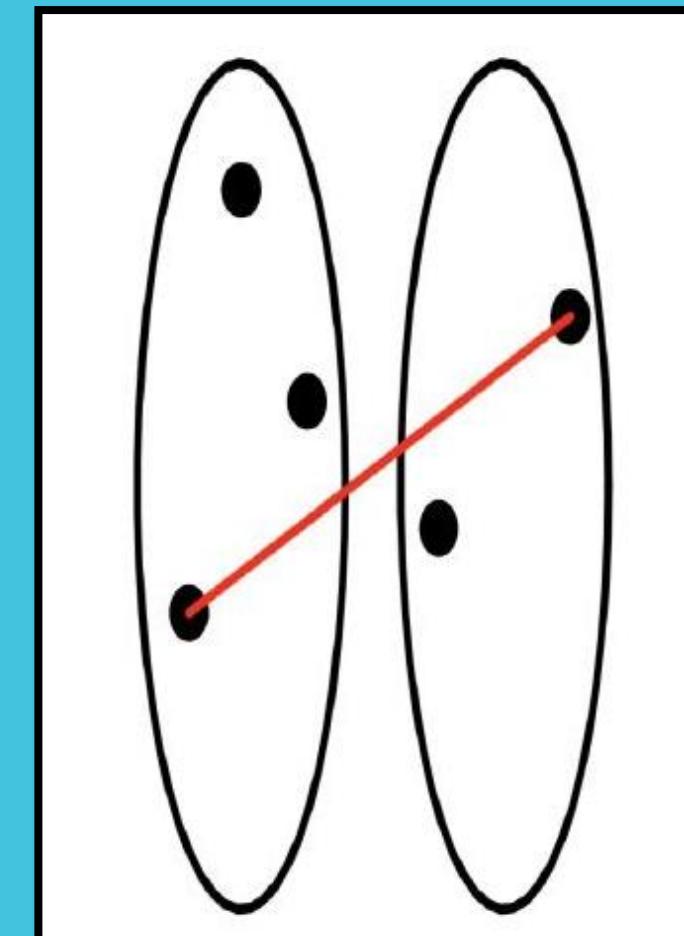
MEASURE OF DISSIMILARITY



Average



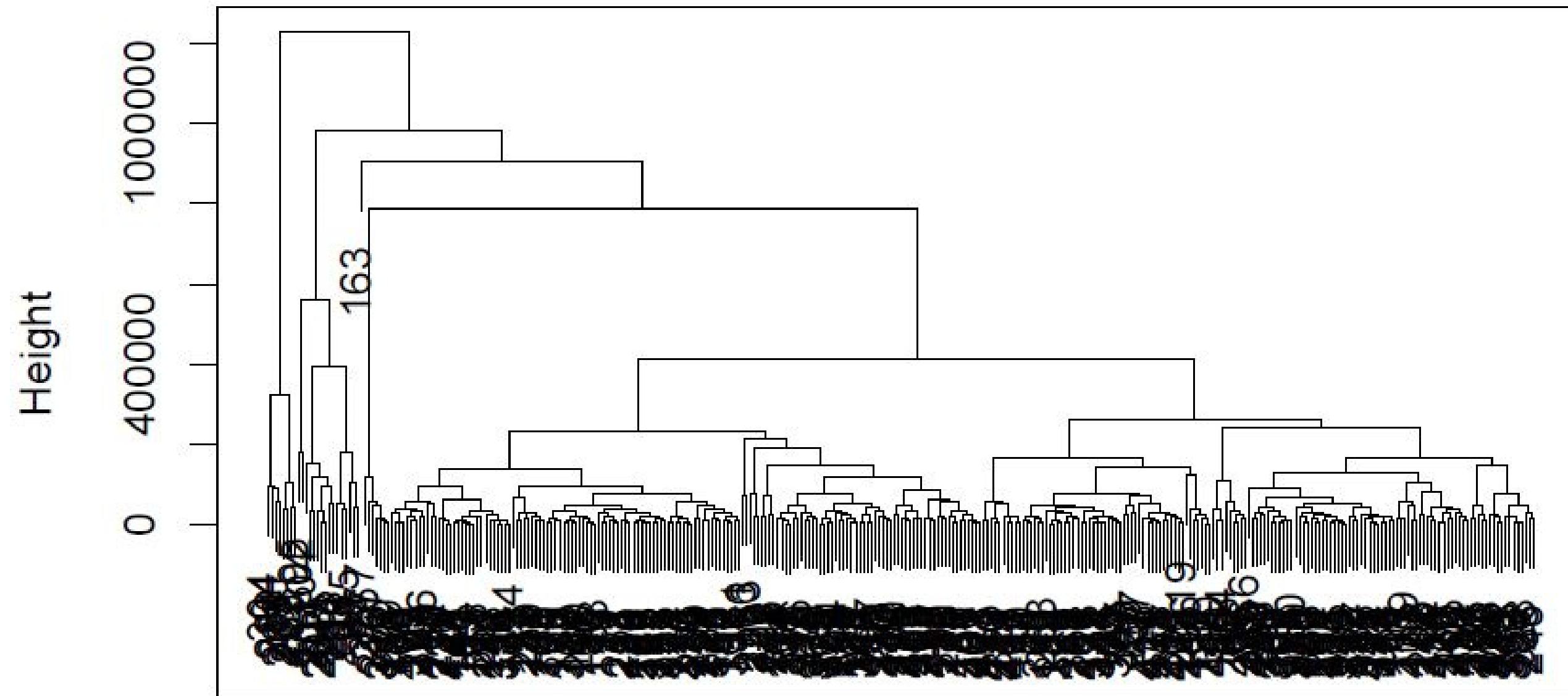
Single



Complete

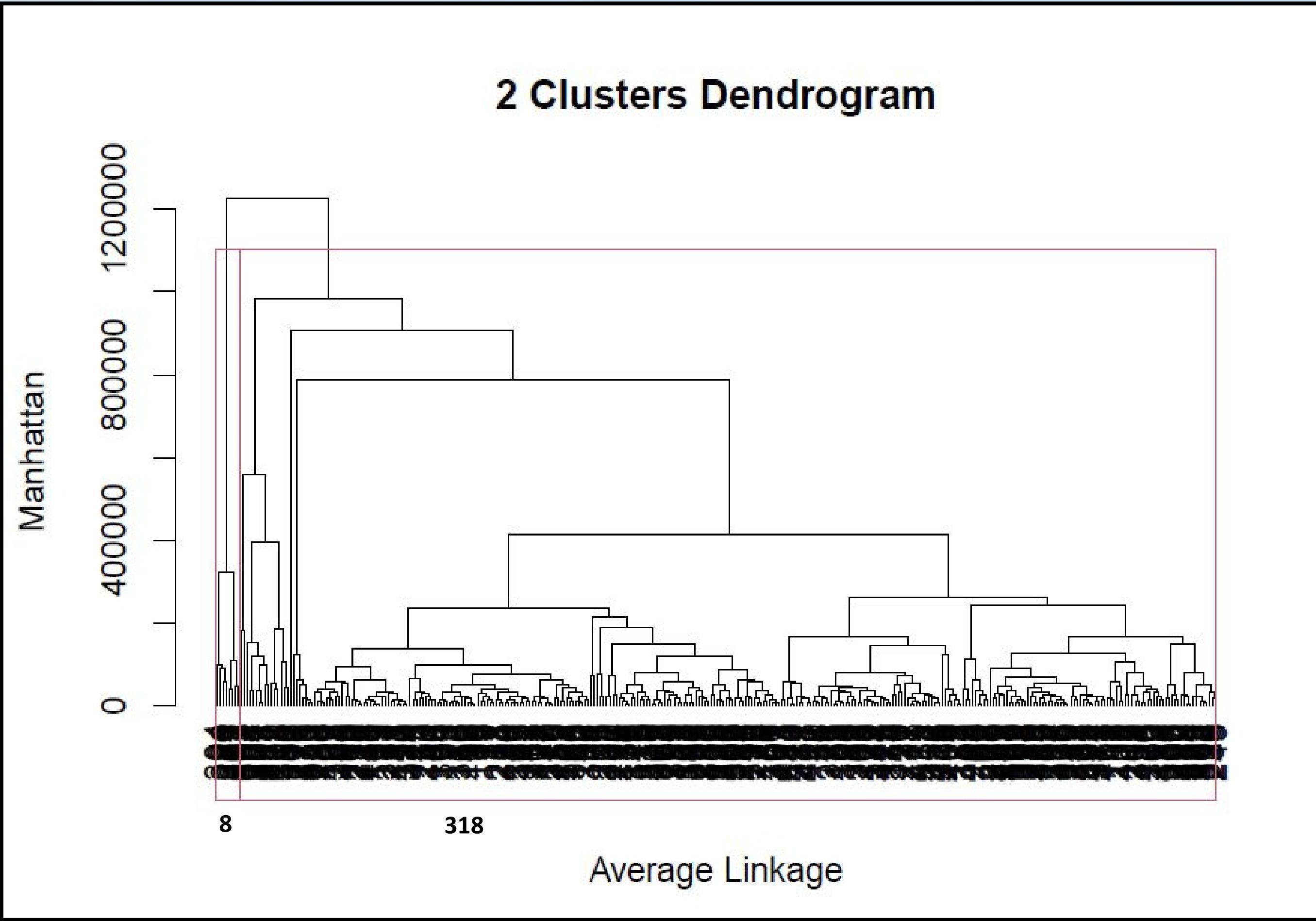
AVERAGE LINKAGE

Cluster Dendrogram (manhattan)



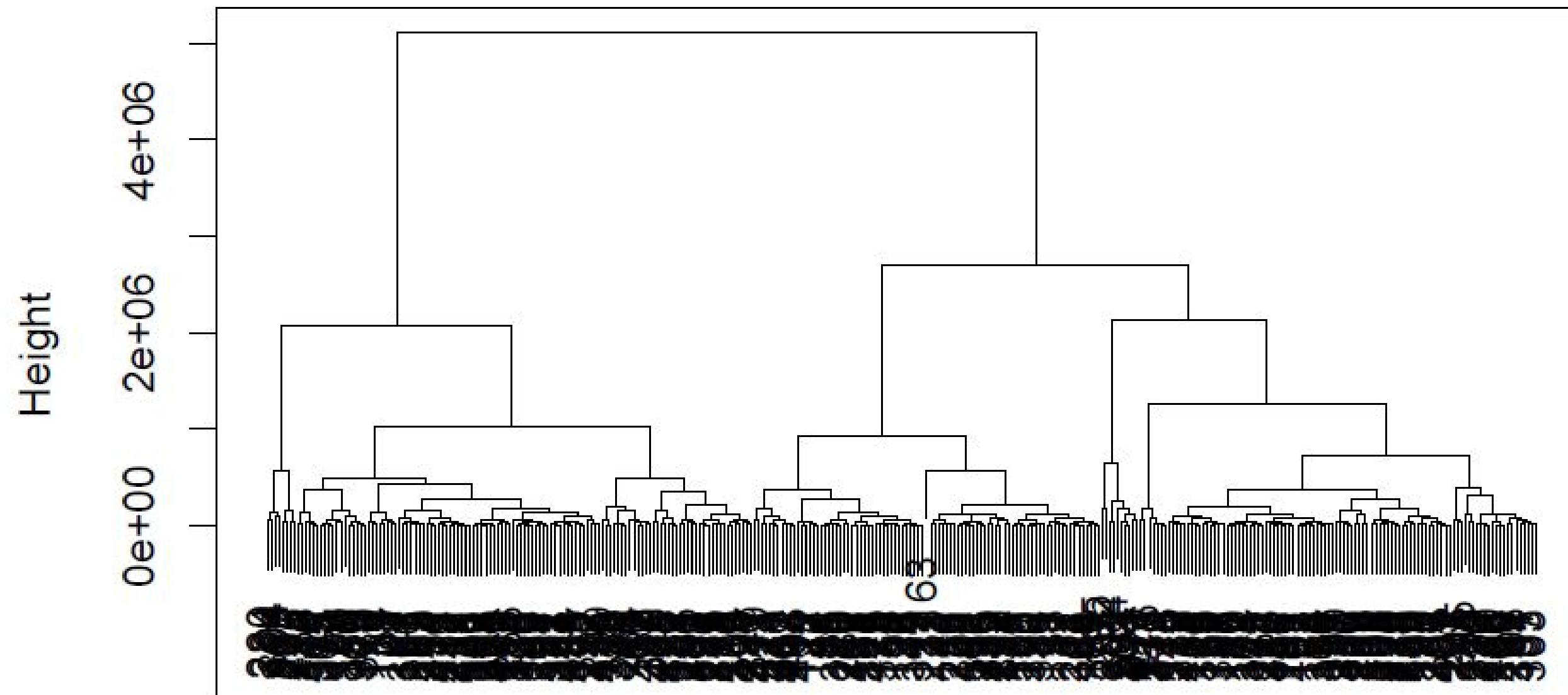
AVERAGE LINKAGE

2 Clusters Dendrogram



COMPLETE LINKAGE

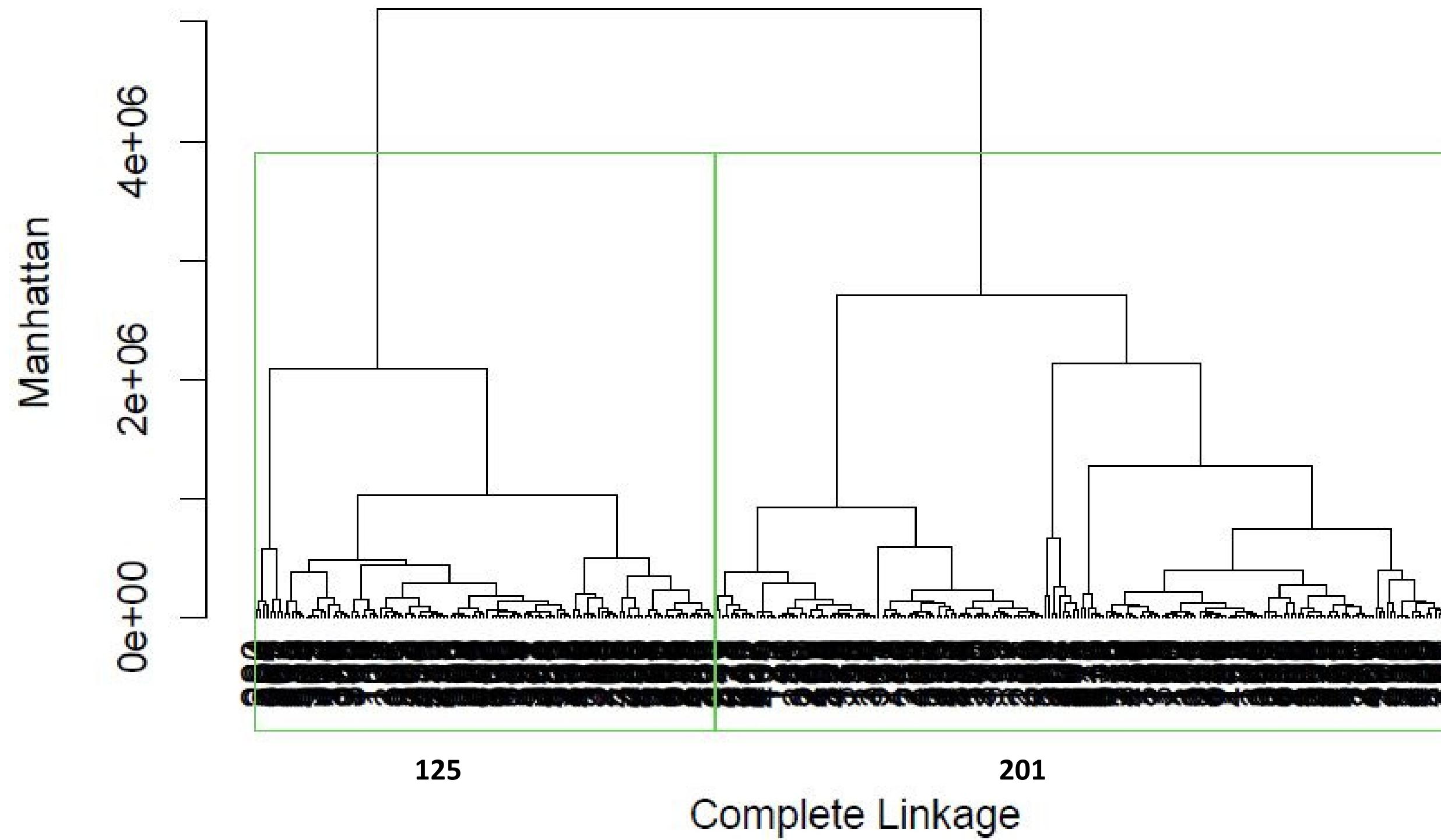
Cluster Dendrogram (manhattan)



d2
hclust (*, "complete")

COMPLETE LINKAGE

2 Clusters Dendrogram



CLASSIFICATION

K-NEAREST NEIGHBORS

Non - Parametric classifier

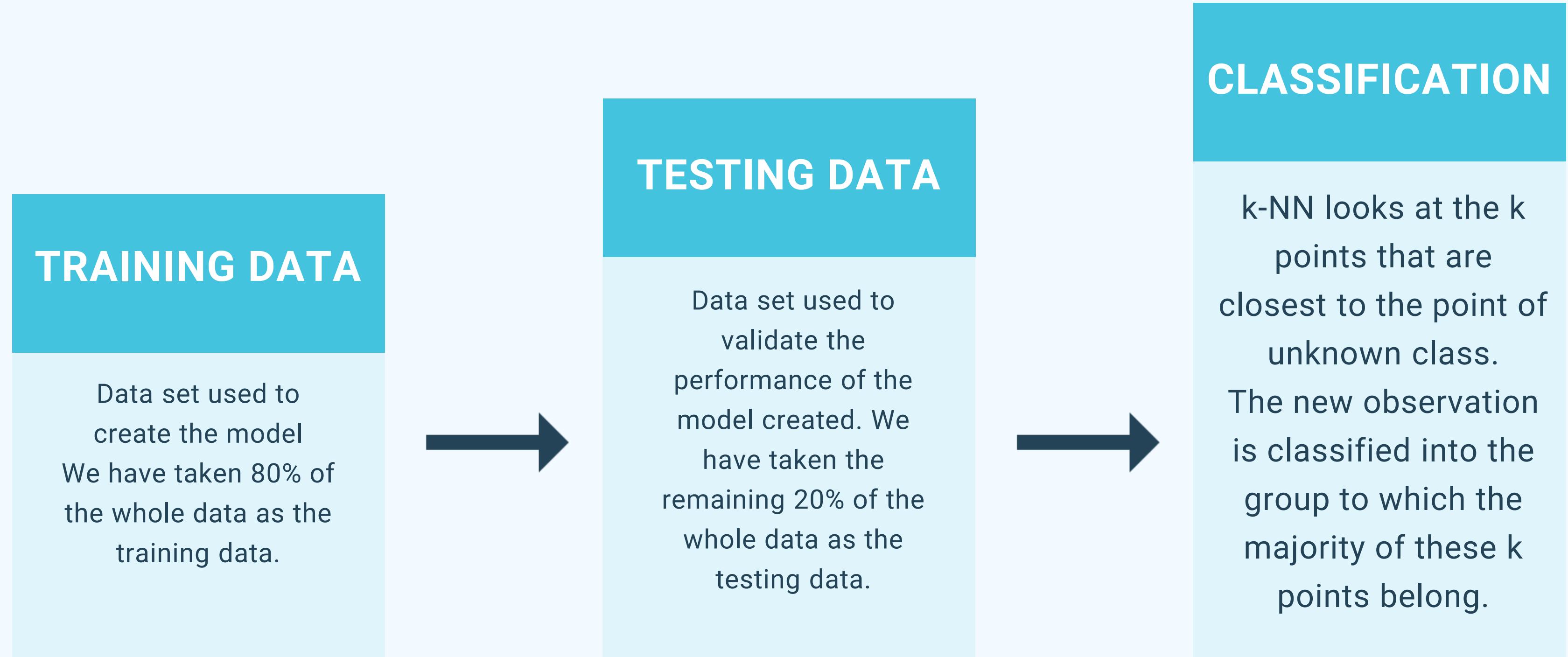
LINEAR DISCRIMINANT ANALYSIS

Parametric classifier

QUADRATIC DISCRIMINANT ANALYSIS

Parametric classifier

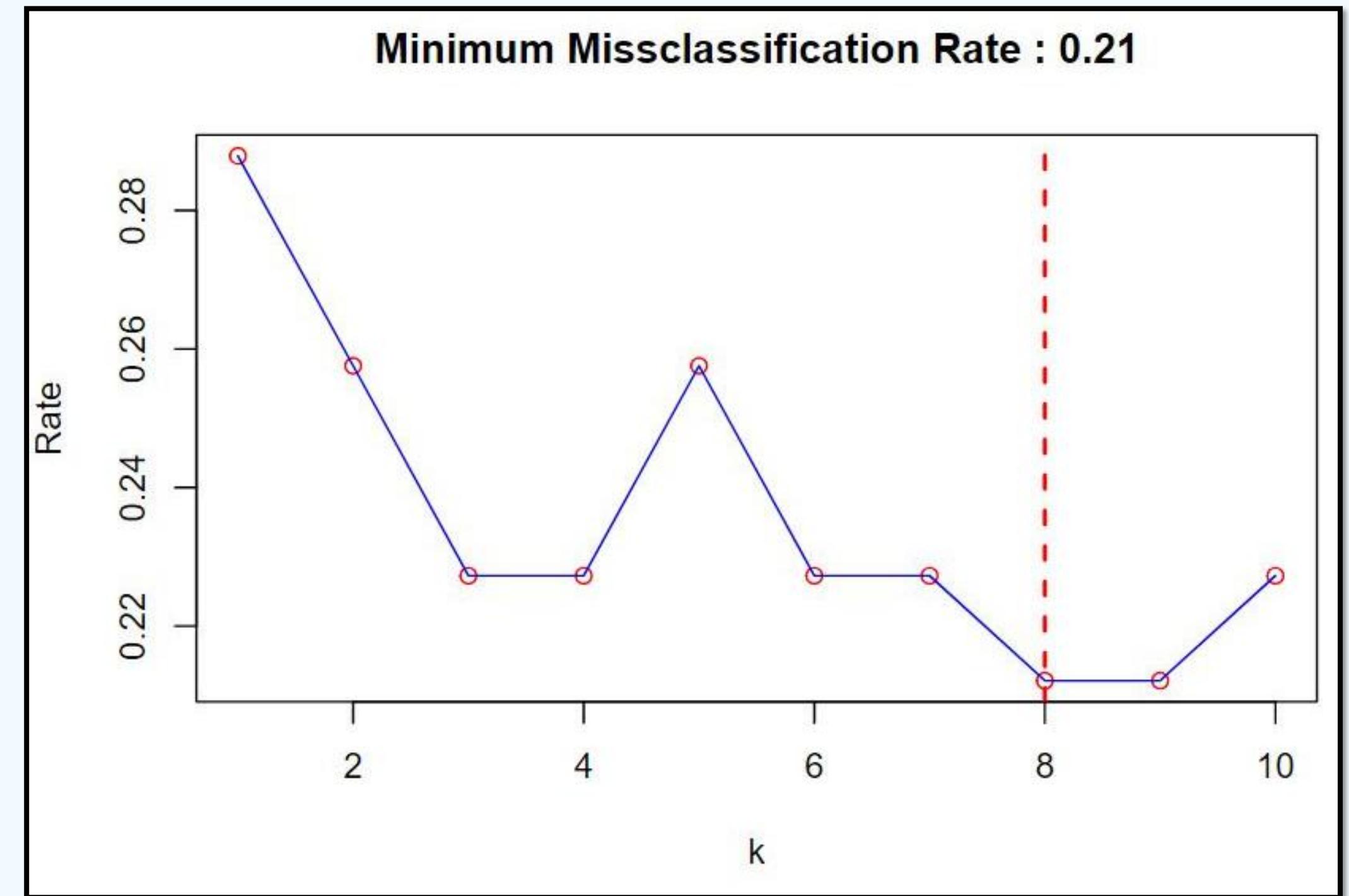
K-NN CLASSIFICATION



WHAT K TO SELECT?

MISCLASSIFICATION RATE

Plot a range of k values against
the misclassification rate and
select the K with the least rate.



RESULT

0 represents: heat stability > 10 mins

1 represents: heat stability < 10 mins

Result	0	1
0	4	4
1	10	48

Accuracy rate = 78.79%

CONCLUSION

All the casein protein traits are positively correlated to protein content. Beta Lactalbumin B is negatively correlated to protein content but is not too significant. Very few technological traits are correlated to each other.

Cow breeds can be divided into two clusters based on the MIR spectra of the milk samples. The two clusters are most populated by Hol fri and JE breed.

We can classify approximately 75% of the observations have a heat stability of less than 10 mins based on the MIR spectra of their milk samples.

**THANK YOU
QUESTIONS?**

