

Creative Data Mining

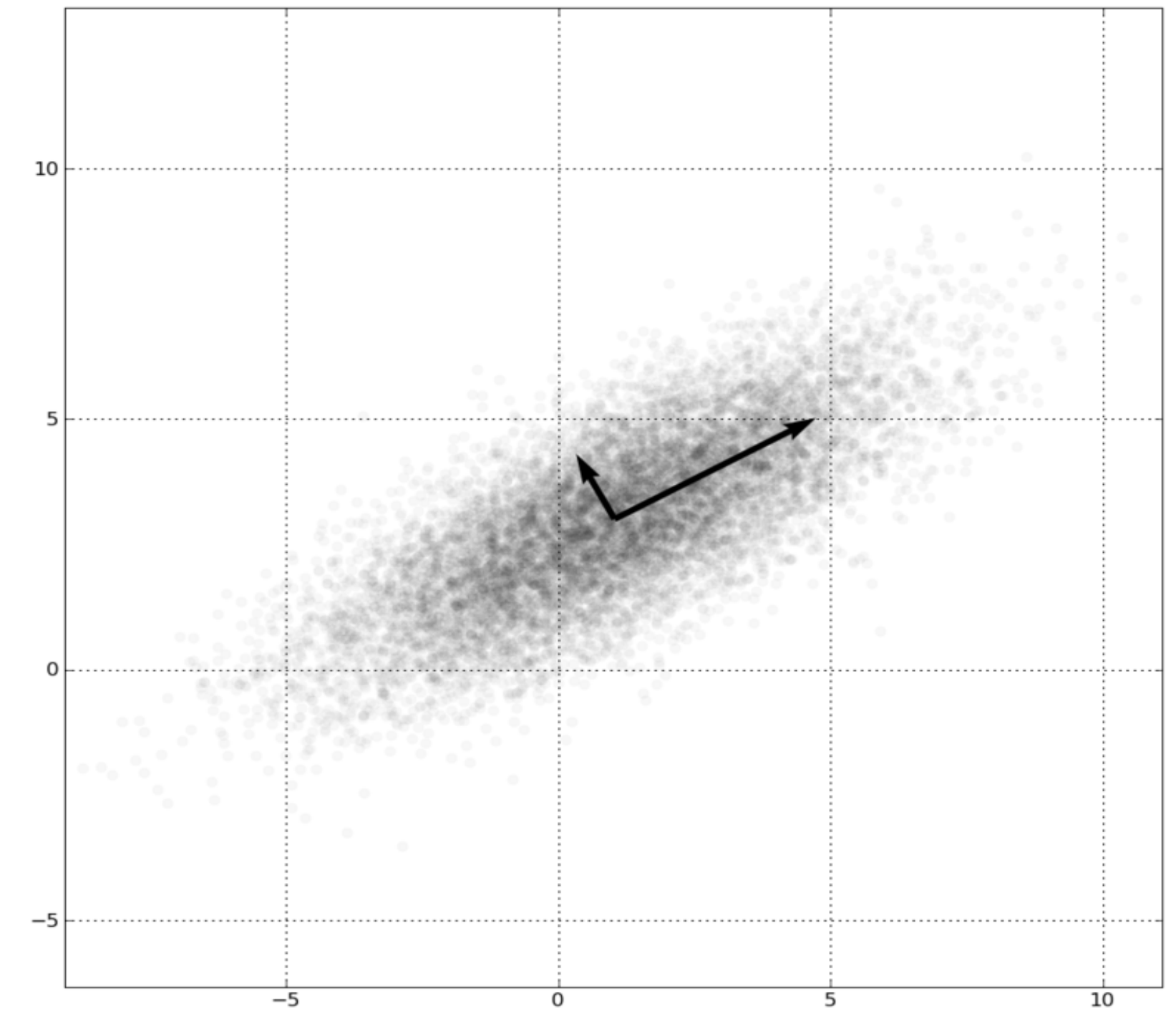
Lecture 03: Intro to RStudio and Clustering
7 March 2016

Danielle Griego, griegod@ethz.ch

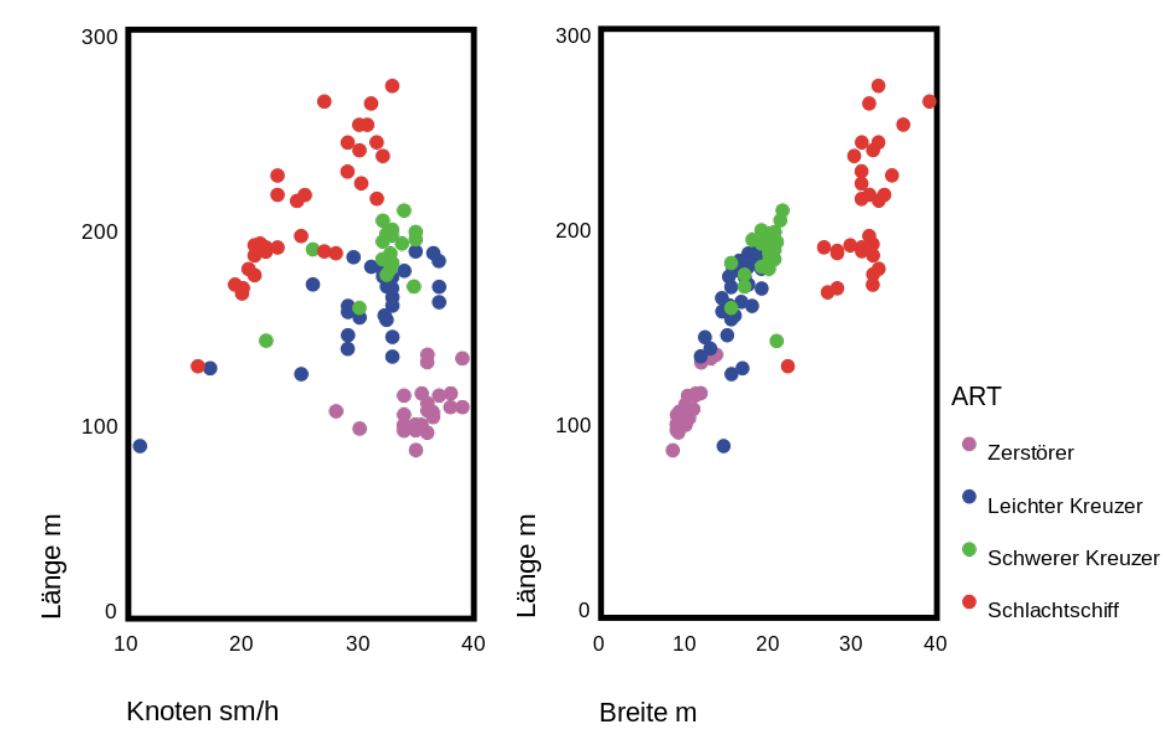
Matthias Standfest, standfest@arch.ethz.ch

THE MATH UNDERNEATH: INTRODUCTION TO CLUSTERING

PRIMARY COMPONENT
ANALYSIS (PCA): MANY
APPROACHES, WE USE
SINGULAR VALUE
DECOMPOSITION (**SVD**)

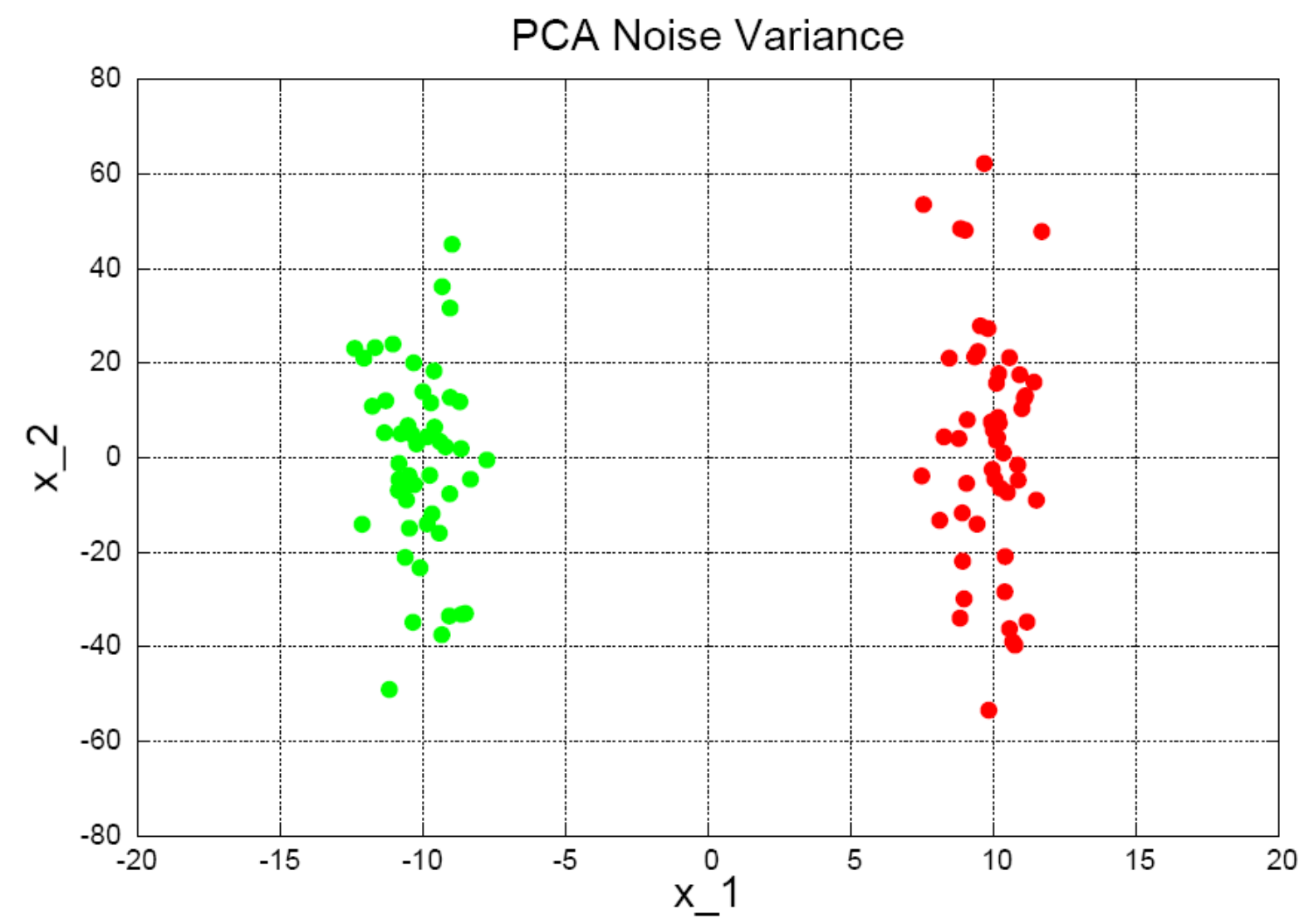
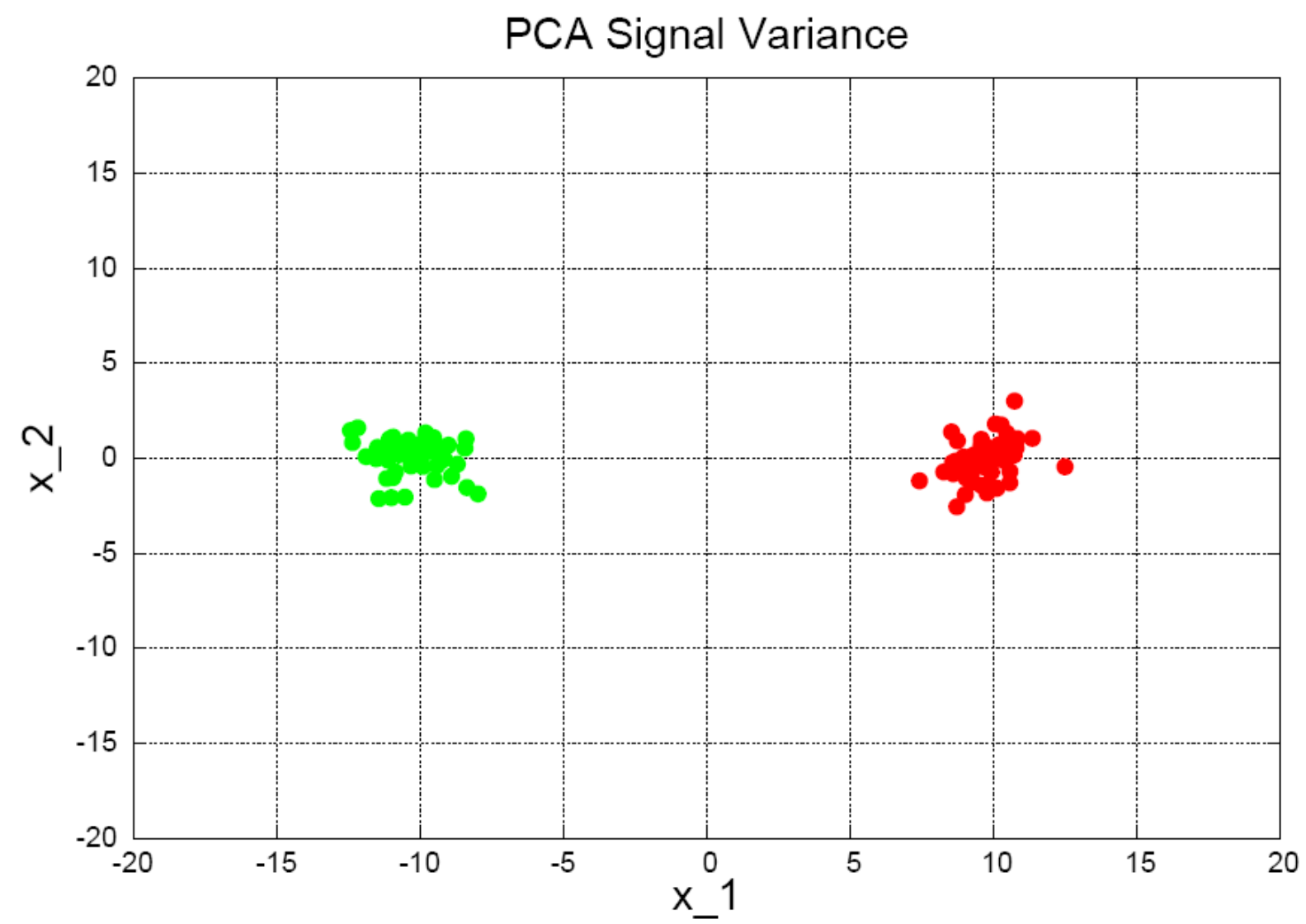


FIRST COMPONENT OF WAR SHIPS:

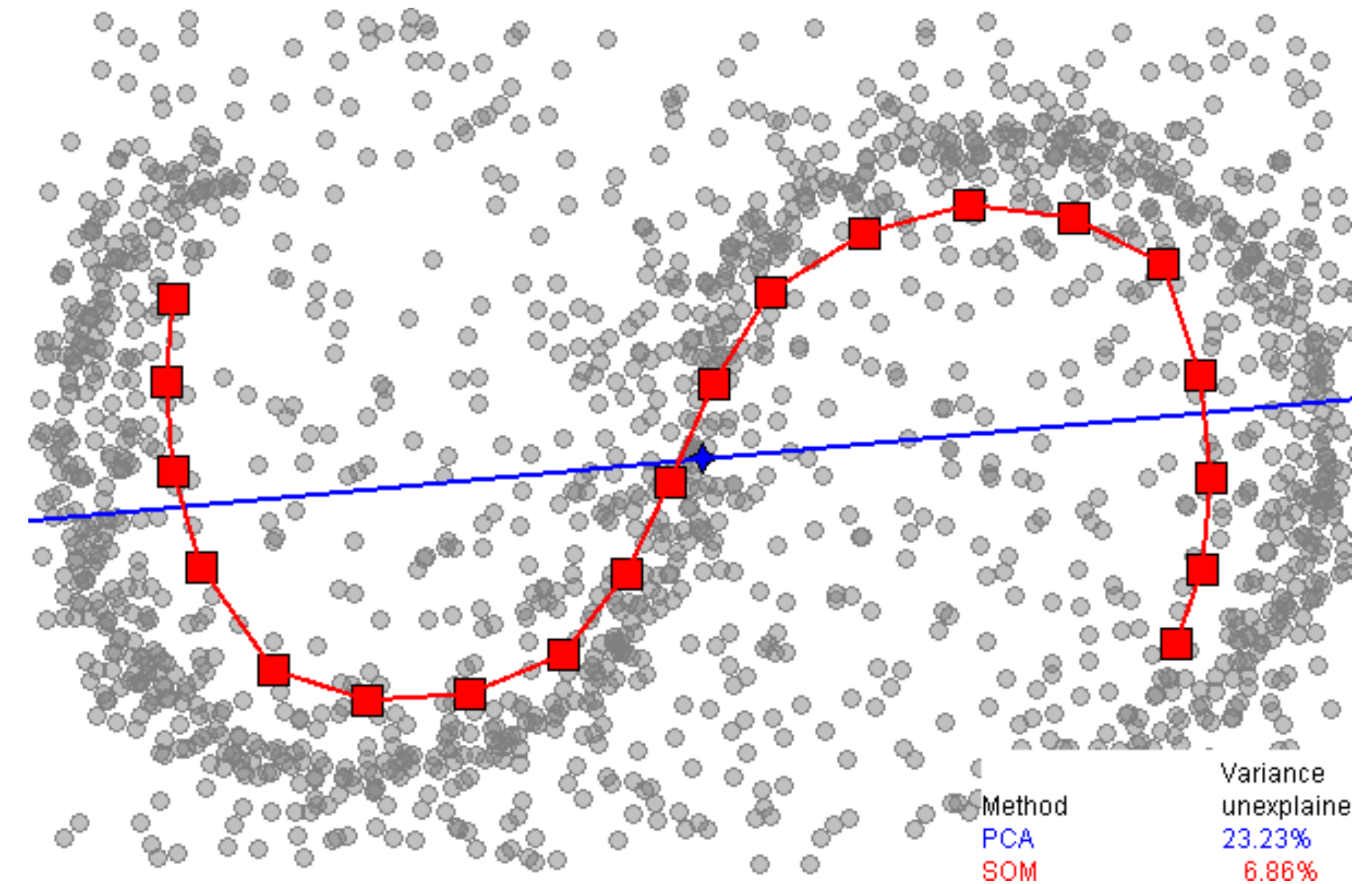


| FACTOR | A | B | C |
|--------|--------|-------|--------|
| LENGTH | 0,862 | 0,481 | -0,159 |
| WIDTH | 0,977 | 0,083 | 0,198 |
| SPEED | -0,679 | 0,730 | 0,082 |

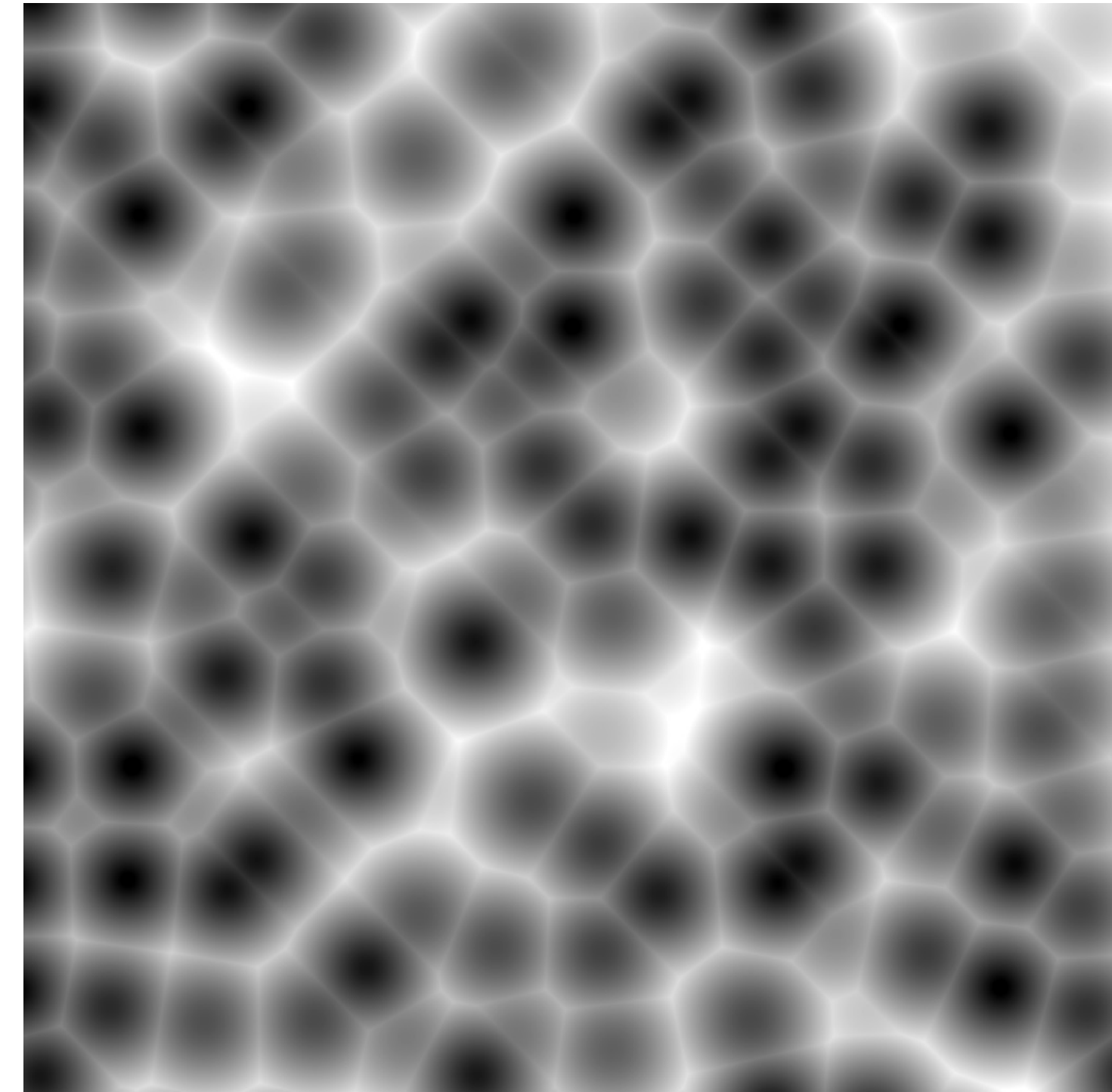
$$Y_A=0.862*LENGTH + 0.977 *WIDTH-0.679*SPEED$$

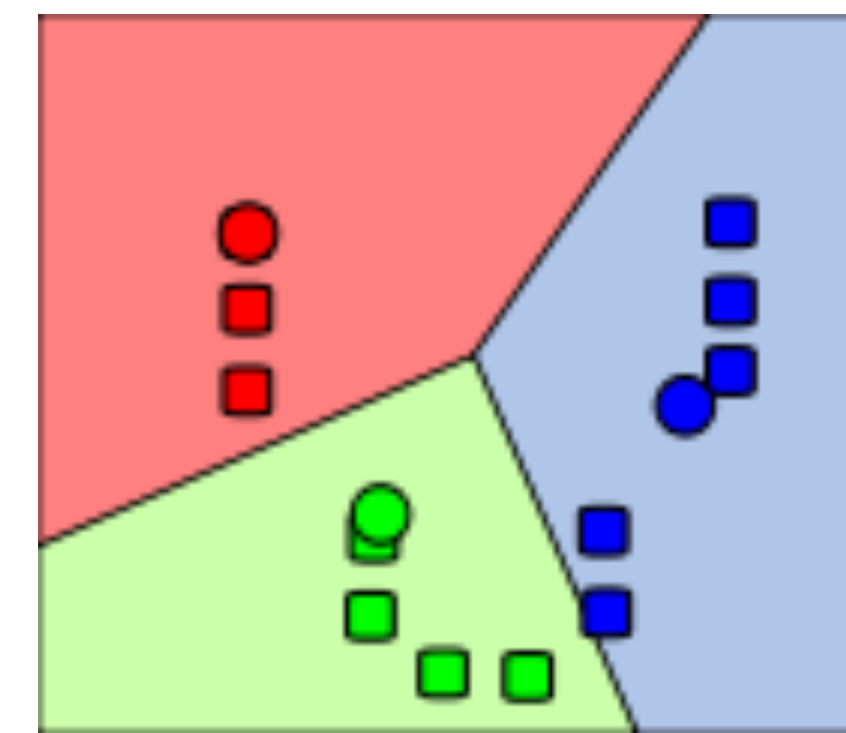
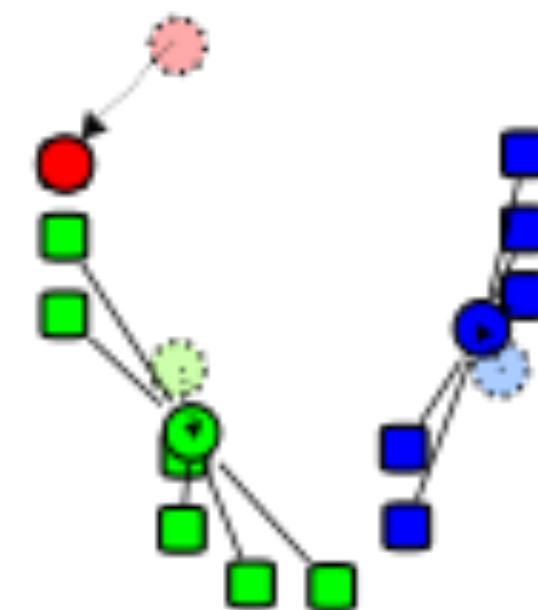
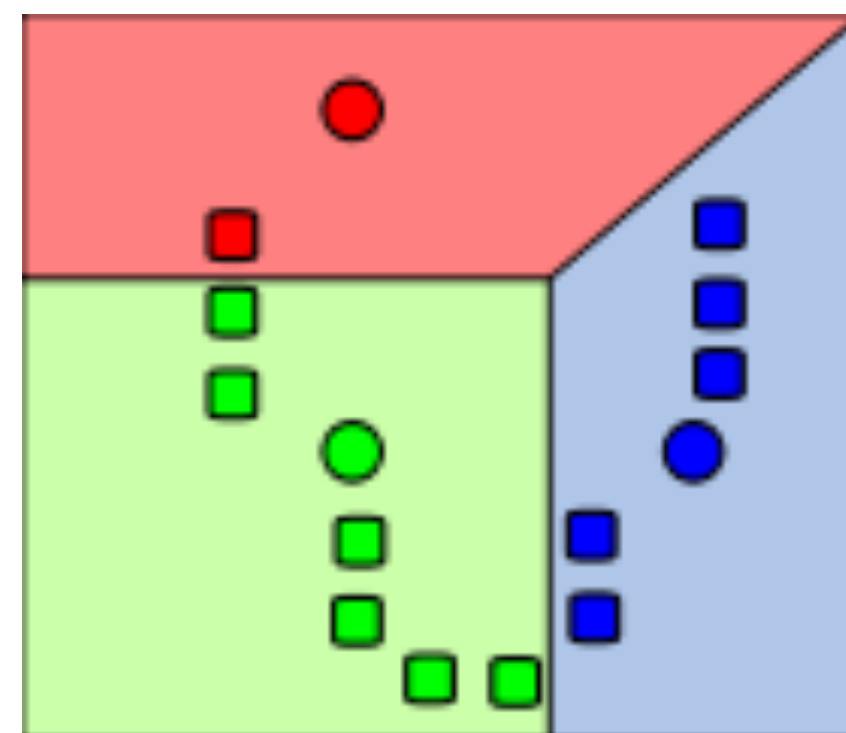
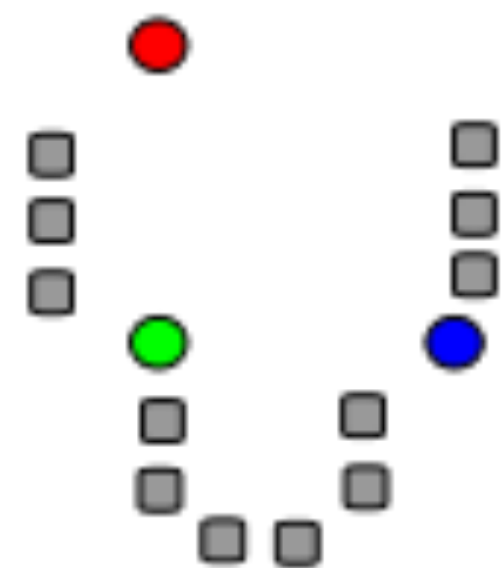


FOR FITTING
STRAIGHT LINES
WITH HIGHEST
POSSIBLE VARIANCE.
PROBLEM: LINEAR.



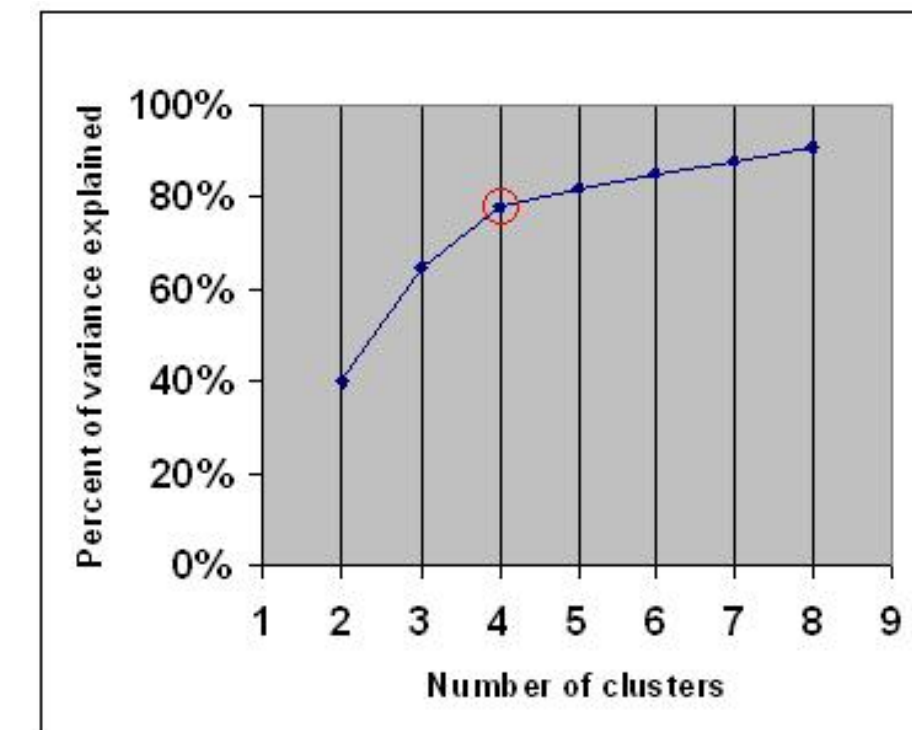
KMEANS:
VORONOI CELLS
FOR EACH
CLUSTER





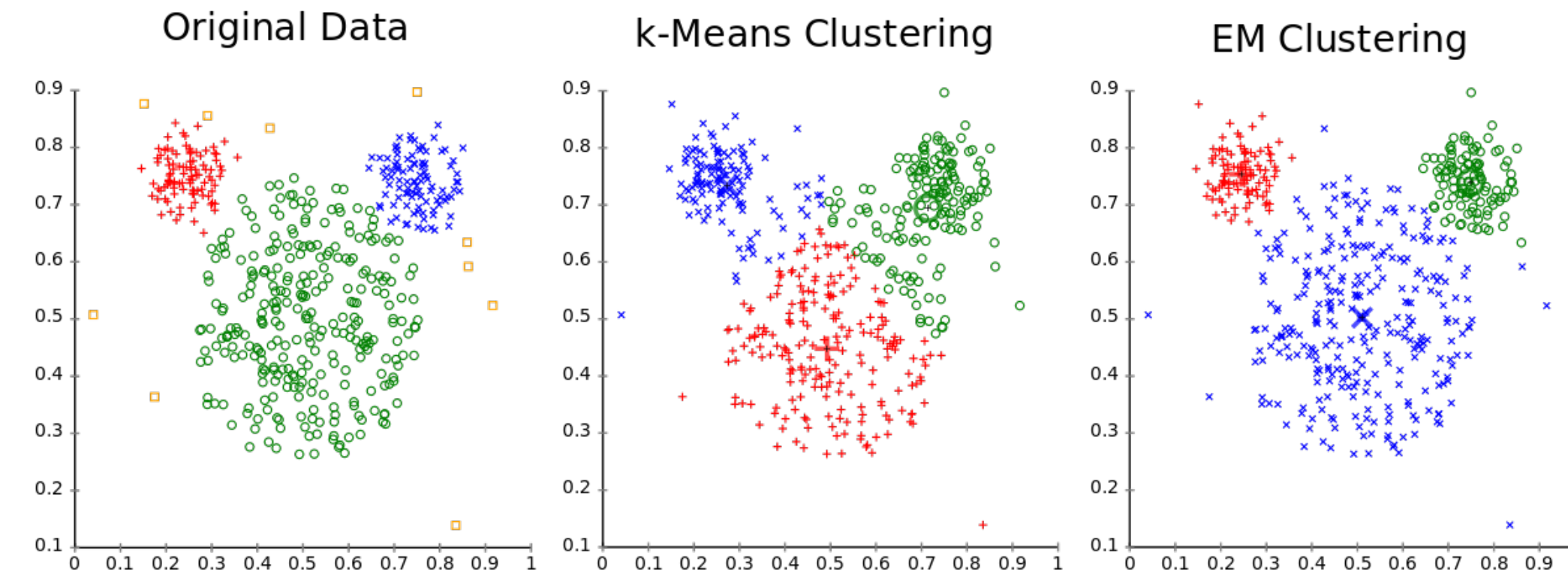
KMEANS:
NUMBER OF
CLUSTERS NEEDS

TO BE ESTIMATED.
ELBOW CURVE.



KMEANS PROBLEM: SAME SIZED CLUSTERS

Different cluster analysis results on "mouse" data set:



FIRST K-MEANS CLUSTERING

lecture_2_05.R

```
1 #
2 # CREATIVE DATA MINING - FS 2016
3 # LECTURE 2 - R101
4 # Matthias Standfest
5 # Danielle Griego
6
7 ## NOW, LET'S JUMP RIGHT IN AND MAKE YOUR FIRST K-MEANS ANALYSIS!
8
9 # 1. set your working directory
10 setwd("~/Dropbox/00_Work/01_Teaching/Creative Data Mining/001_FS15-DataMining/lecture_2/RSCRIPTS_FS16/")
11
12 # 2. get data
13 iris2 <- iris[,1:4]
14
15 # 3. Prepare data and inspect data
16 mydata <- na.omit(iris2) # listwise deletion of missing
17 summary(mydata) # look at the data
18 plot(mydata$Sepal.Length)
19
20 help(scale) # scale is generic function whose default method centers and/or scales the columns of a numeric matrix
21 mydata_s <- scale(mydata) # standardize variable
22 summary(mydata_s)
23 #plot(mydata_s$Sepal.Length) # note that this will not work, error will say that "$ operator is invalid for atomic vectors"
24 plot(mydata_s[,1]) # so need to look at lecture_2_04.R for other column references
25
26 # 4. Determine number of clusters
27 maxCluster <- round(sqrt(nrow(mydata_s)/2)*2) # max cluster estimation
28 set.seed(1) #set seed of random variable so we have the same results each time (to be able to compare)
29 library(NbClust) #import a library
30 result<-NbClust(mydata, diss=NULL, distance = "euclidean", min.nc=2, max.nc=maxCluster,
31   | | | | | method = "complete", index = "kl")
32 (CLUSTERESTIMATION <- result$Best.nc[1])
33
34 # 5. K-Means cluster analysis
35 set.seed(1)
36 fit <- kmeans(mydata, CLUSTERESTIMATION)
37
38 # 6.1. make plot using package cluster
39 library(cluster)
40 clusplot(mydata,fit$cluster)
41
42 # 6.2. make plot using package ade4
43 library(ade4)
44 pca <-prcomp(mydata, scale.=T, retx=T) # principal components analysis
45 plot.mydata <- cbind(pca$x[,1], pca$x[,2]) # first and second PC
46 s.class(plot.mydata, factor(fit$cluster))
47
```


DIFFERENT VISUALISATION

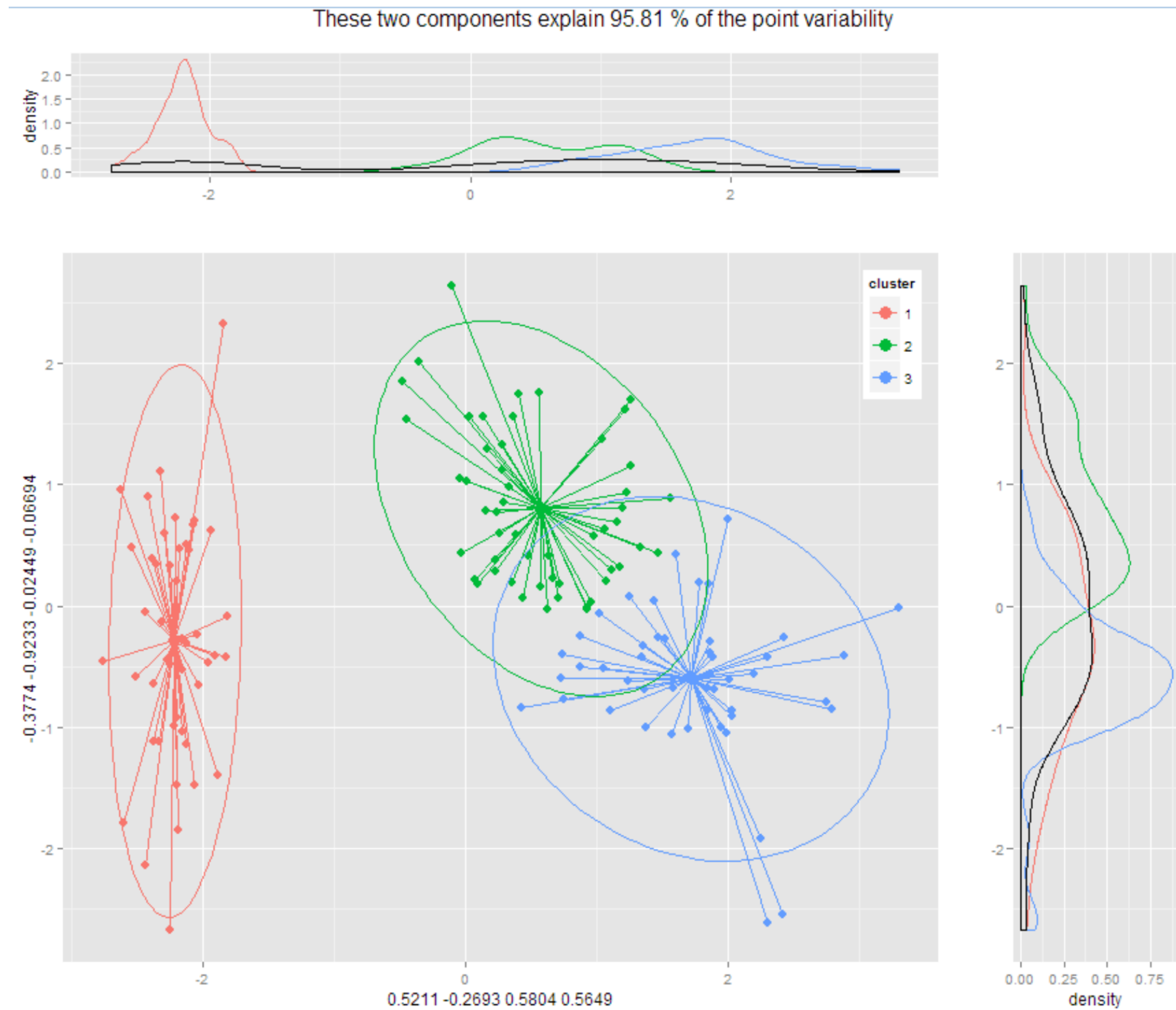
lecture_2_06.R

```
1 #
2 # CREATIVE DATA MINING - FS 2016
3 # LECTURE 2 - R101
4 # Matthias Standfest
5 # Danielle Griego
6
7 ## NOW LET'S IMPROVE THE VISUALIZATION FROM THE PREVIOUS ANALYSIS
8
9 # 1. set your working directory
10 setwd("~/Dropbox/00_Work/01_Teaching/Creative Data Mining/001_FS15-DataMining/lecture_2/RSCRIPTS_FS16/")
11
12 # 2. get data
13 iris2 <- iris[,1:4]
14
15 # 3. prepare data
16 mydata <- na.omit(iris2) # listwise deletion of missing
17 summary(mydata)
18 mydata <- scale(mydata) # standardize variables
19 summary(mydata)
20
21 # 4. Determine number of clusters
22 maxCluster <- round(sqrt(nrow(mydata)/2)*2) # max cluster estimation
23 set.seed(1) #set seed of random variable so we have the same results each time (to be able to compare)
24 library(NbClust) #import a library
25 result<-NbClust(mydata, diss=NULL, distance = "euclidean", min.nc=2, max.nc=maxCluster,
26 | | | | | method = "complete", index = "kl")
27 (CLUSTERESTIMATION <- result$Best.nc[1])
28
29 # 5. K-Means cluster analysis
30 set.seed(1)
31 fit <- kmeans(mydata, CLUSTERESTIMATION)
32
33 # 6. plot clustering using GGLOT2
34 # Cluster Plot against 1st 2 principal components
35 # vary parameters for most readable graph
36 # ggplot solution for clusplot(mydata, fit$cluster, color=FALSE, shade=TRUE, labels=2, lines=0)
37 pca <- prcomp(mydata, scale.=T, retx=T) # principal components analysis
38 # gg: data frame of PC1 and PC2 scores with corresponding cluster
39 gg <- data.frame(cluster=factor(fit$cluster), x=pca$x[,1], y=pca$x[,2])
40 # calculate cluster centroid locations
41 centroids <- aggregate(cbind(x,y)~cluster, data=gg, mean)
42 # merge centroid locations into ggplot dataframe
43 gg <- merge(gg, centroids, by="cluster", suffixes=c("", ".centroid"))
44 # calculate 95% confidence ellipses
```

```
45 library(ellipse)
46 conf.rgn <- do.call(rbind,lapply(1:3,function(i)
47   cbind(cluster=i,ellipse(cov(gg[gg$cluster==i,2:3]),centre=as.matrix(centroids[i,2:3])))))
48 conf.rgn <- data.frame(conf.rgn)
49 conf.rgn$cluster <- factor(conf.rgn$cluster)
50 # plot cluster map
51 library(ggplot2)
52
53 cumulativeVariability <- (cumsum((pca$sdev)^2) / sum(pca$sdev^2))[2] #cumulativeVariability might return a characte
54 class(cumulativeVariability)
55 #cumulativeVariability <- as.numeric((cumsum((pca$sdev)^2) / sum(pca$sdev^2))[2]) # so might need to use as.numeric
56
57 labelx <- paste(formatC(pca$rotation[,1], width=5), collapse = ' ')
58 labely <- paste(formatC(pca$rotation[,2], width=5), collapse = ' ')
59 clusterscatter <- ggplot(gg, aes(x,y, color=cluster))+
60   geom_point(size=3) +
61   geom_point(data=centroids, size=4) +
62   geom_segment(aes(x=x.centroid, y=y.centroid, xend=x, yend=y))+
63   geom_path(data=conf.rgn)+
64   ylab(labely)+
65   xlab(labelx)+
66   theme(legend.position=c(1,1),legend.justification=c(1,1))
67 plot(clusterscatter) # plot nices clusters with ellipses using ggplot2
68
69 # 7. add distribution clusterwise to the primary component axes
70 # http://www.r-bloggers.com/ggplot2-cheatsheet-for-visualizing-distributions/ #good link!
71 plot_top <- ggplot(gg, aes(x=x, col=cluster)) +
72   geom_density(alpha=.5) +
73   geom_density(color="black") +
74   theme(legend.position = "none", axis.title.x = element_blank())
75 plot_right <- ggplot(gg, aes(x=y, col=cluster)) +
76   geom_density(alpha=.5) +
77   geom_density(color="black") +
78   coord_flip() +
79   theme(legend.position = "none", axis.title.y = element_blank())
80 #placeholder plot - prints nothing at all
81 empty <- ggplot()+geom_point(aes(1,1), colour="white") +
82   theme(
83     plot.background = element_blank(),
84     panel.grid.major = element_blank(),
85     panel.grid.minor = element_blank(),
86     panel.border = element_blank(),
87     panel.background = element_blank(),
88     axis.title.x = element_blank(),
89     axis.title.y = element_blank(),
90     axis.text.x = element_blank(),
91     axis.text.y = element_blank(),
92     axis.ticks = element_blank()
93   )
94 #arrange the plots together, with appropriate height and width for each row and column
95 library(gridExtra)
96 library(grid)
97 chart1 <- grid.arrange(plot_top, empty, clusterscatter, plot_right, ncol=2, nrow=2, widths=c(4, 1), heights=c(1, 4),
98 | | | | | top=paste("These two components explain", formatC(cumulativeVariability*100, width=4), "% of the point variability"))
99
```

DIFFERENT VISUALISATION

lecture_2_06.R

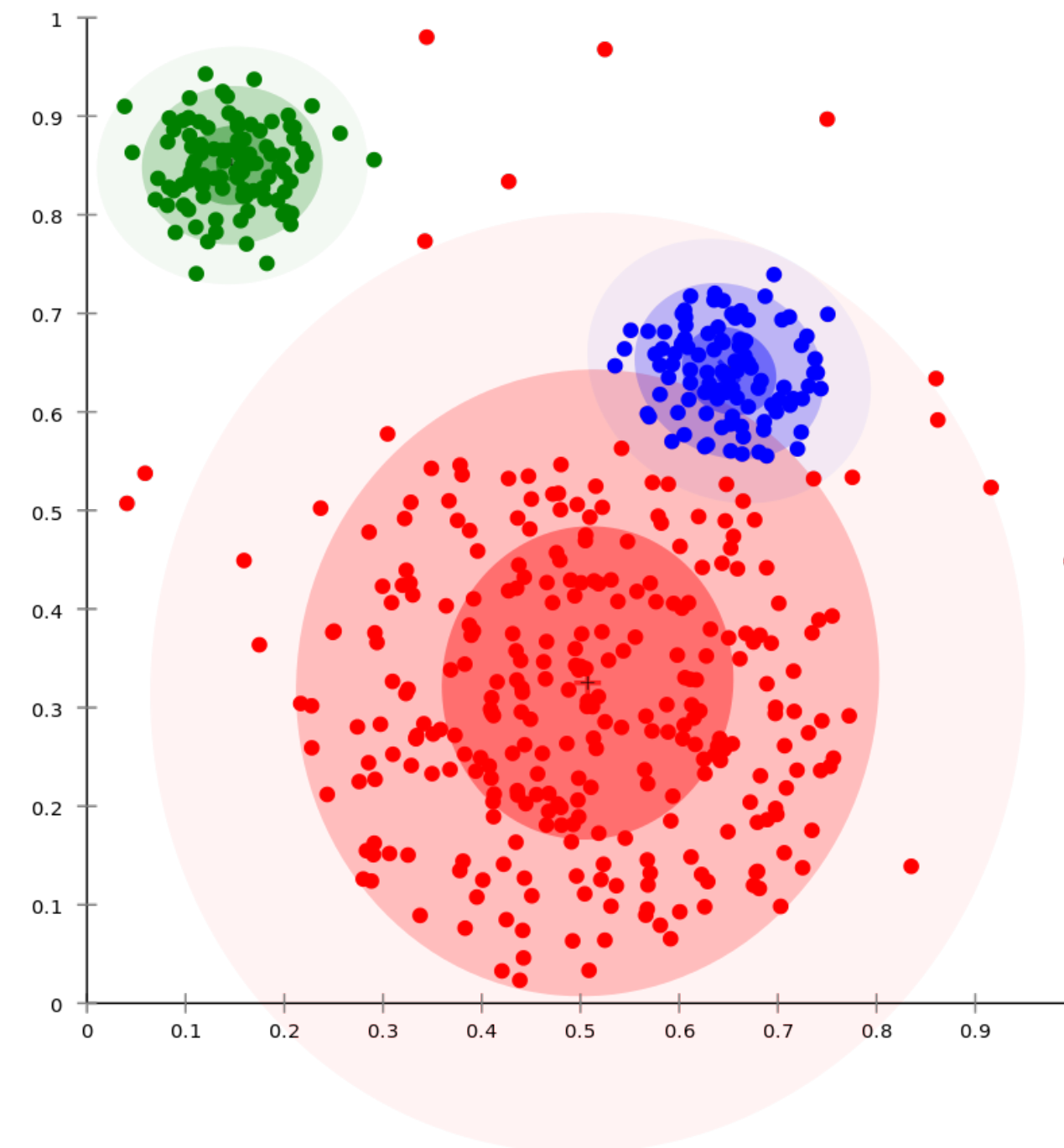


codingStyle – StyleTheCode() *#TODO(you):restyle*

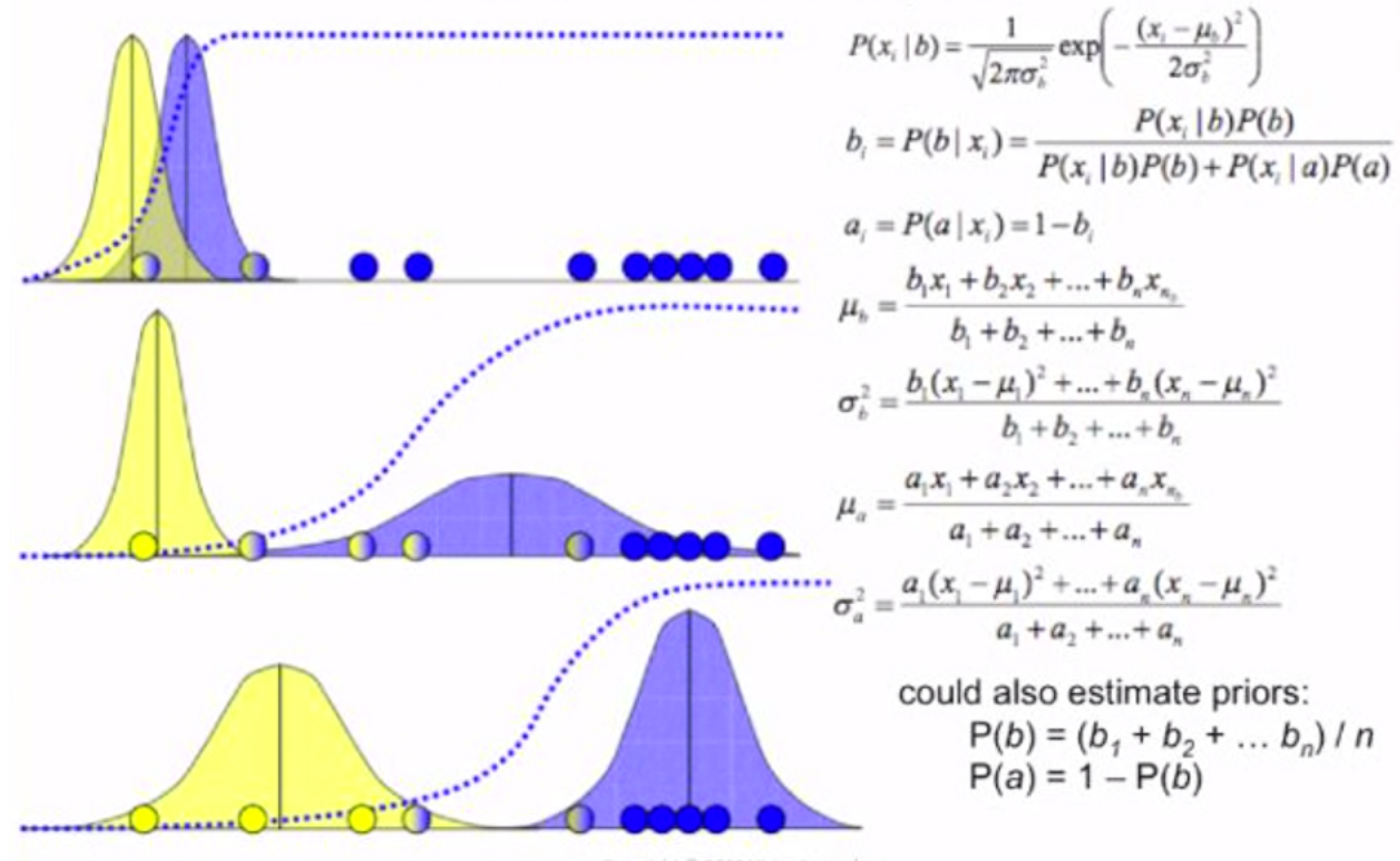
- <https://google.github.io/styleguide/Rguide.xml>
- kApple, vectorBanana, matrixCherry, scalarDurian
- GatherElderBerries(100)
- <https://stat.ethz.ch/R-manual/R-devel/library/datasets/html/00Index.html>
- str(faithful) *# mtcars, rock, trees, LifeCycleSavings*
- *#vectorized instead of for, boolean tables as filters instead of if clause*

OTHER CLUSTERING TECHNIQUES

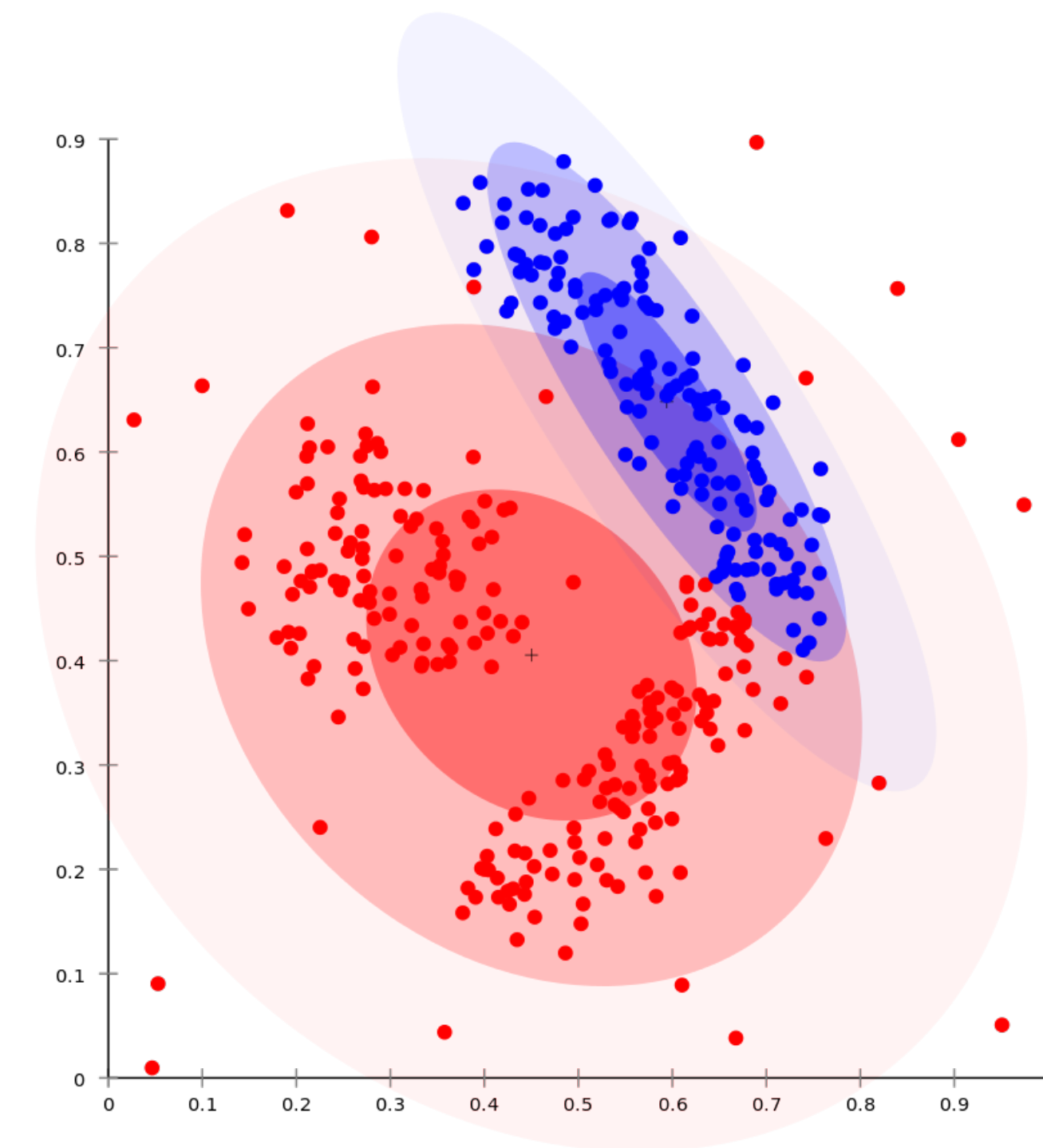
EXPECTATION MAXIMIZATION ALGORITHM WORKS WITH DISTRIBUTION



EM: 1-d example

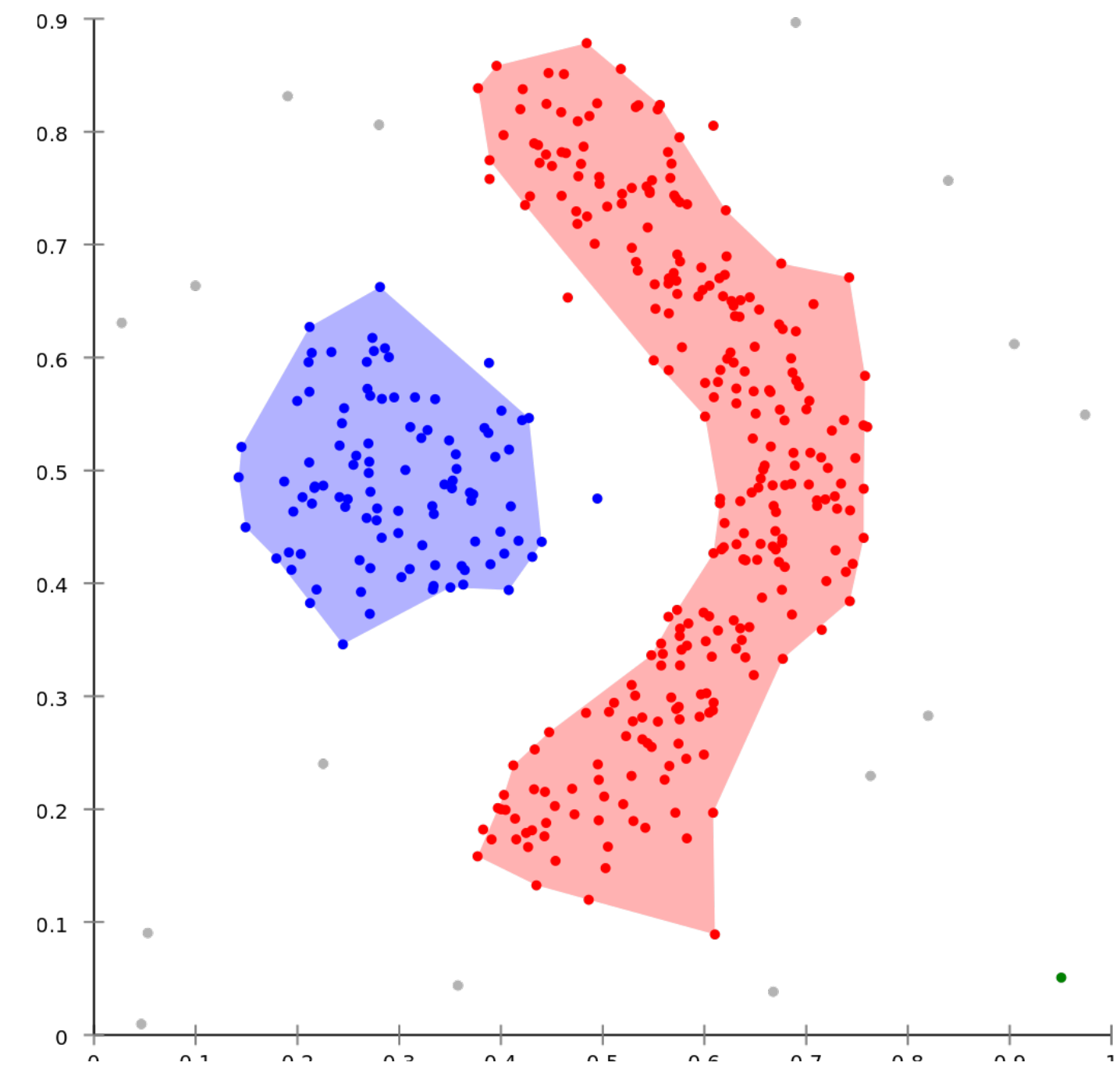


EM ALGO.: PROBLEM DENSITY CLUSTERS

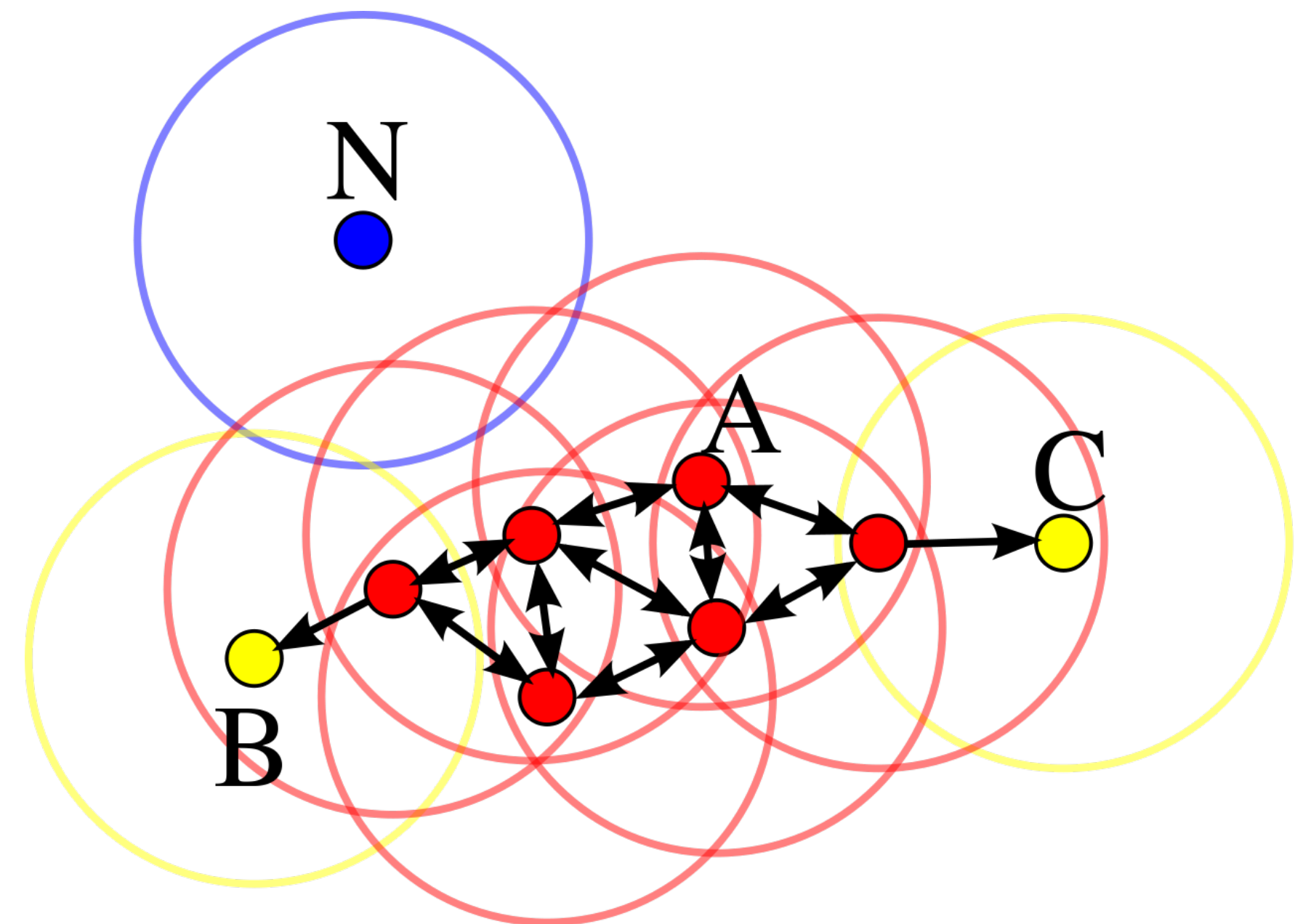


DBSCAN

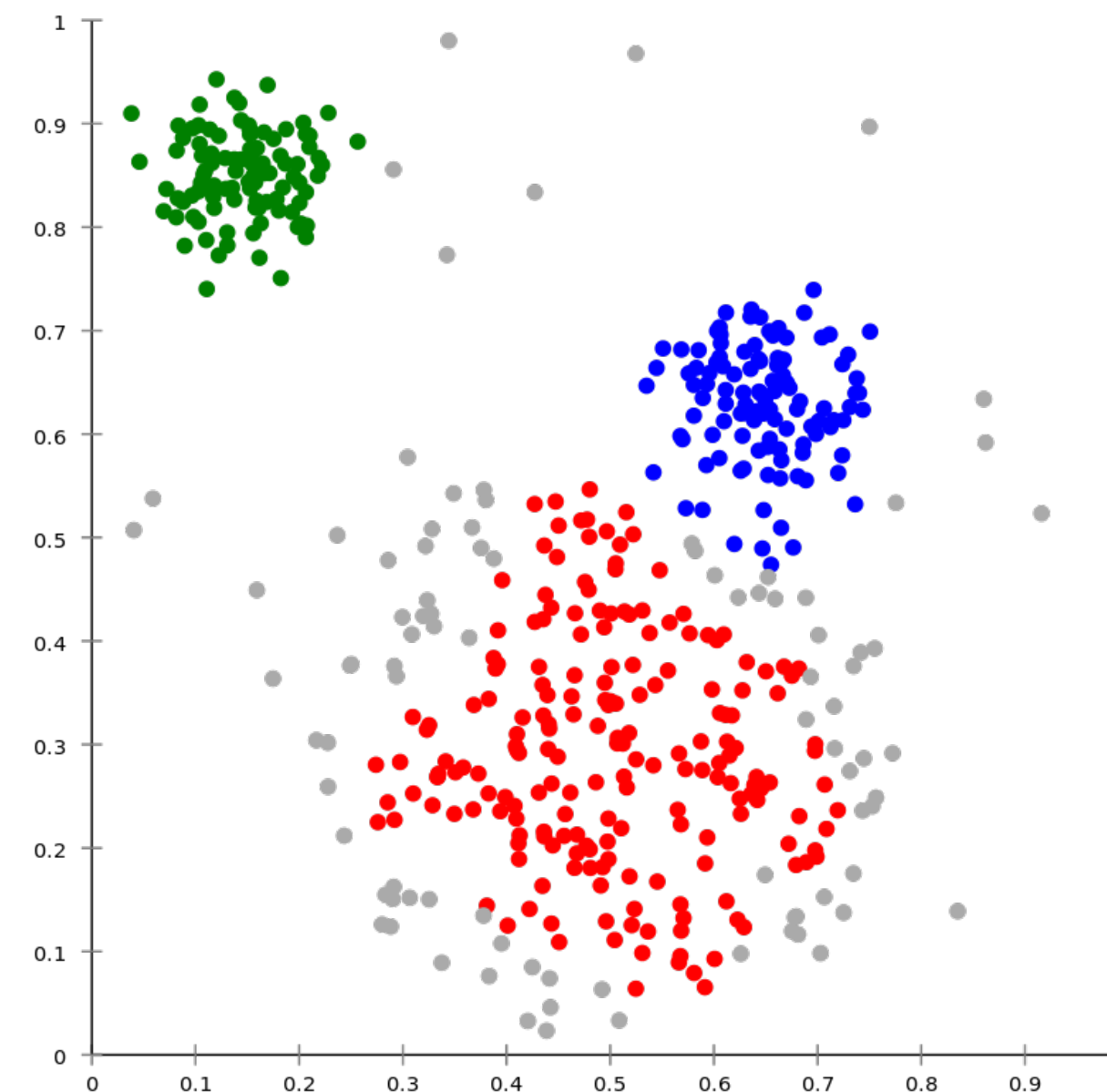
DENSITY-BASED SPATIAL CLUSTERING OF APPLICATIONS WITH NOISE



IN THIS DIAGRAM, $\text{MINPTS} = 3$. A AND THE OTHER RED POINTS ARE CORE POINTS, BECAUSE AT LEAST THREE POINTS SURROUND IT IN AN R RADIUS. BECAUSE THEY ARE ALL REACHABLE FROM ONE ANOTHER, THEY FORM A SINGLE CLUSTER. POINTS B AND C ARE NOT CORE POINTS, BUT ARE REACHABLE FROM A (VIA OTHER CORE POINTS) AND THUS BELONG TO THE CLUSTER AS WELL. POINT N IS A NOISE POINT THAT IS NEITHER A CORE POINT NOR DENSITY-REACHABLE.

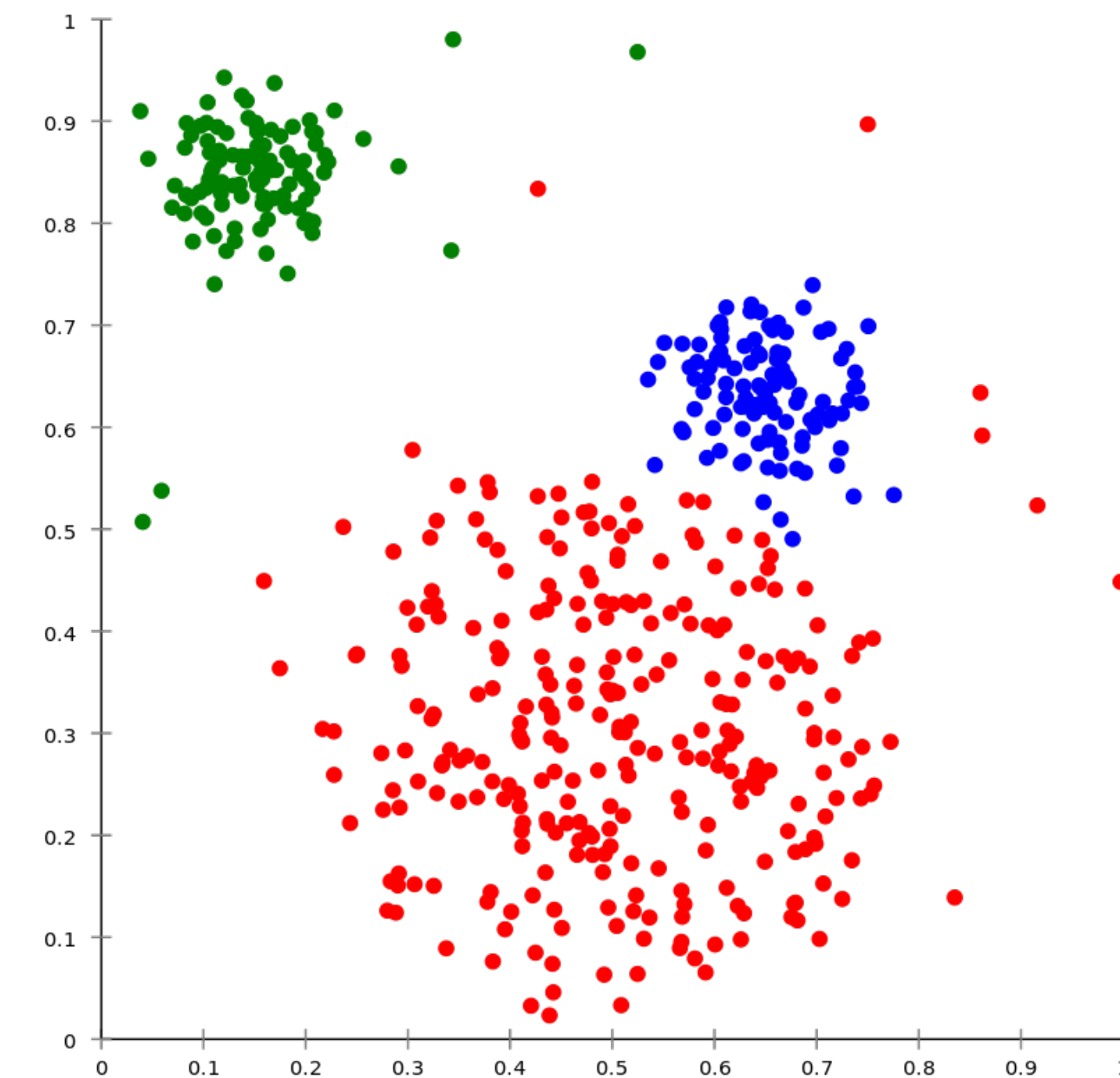


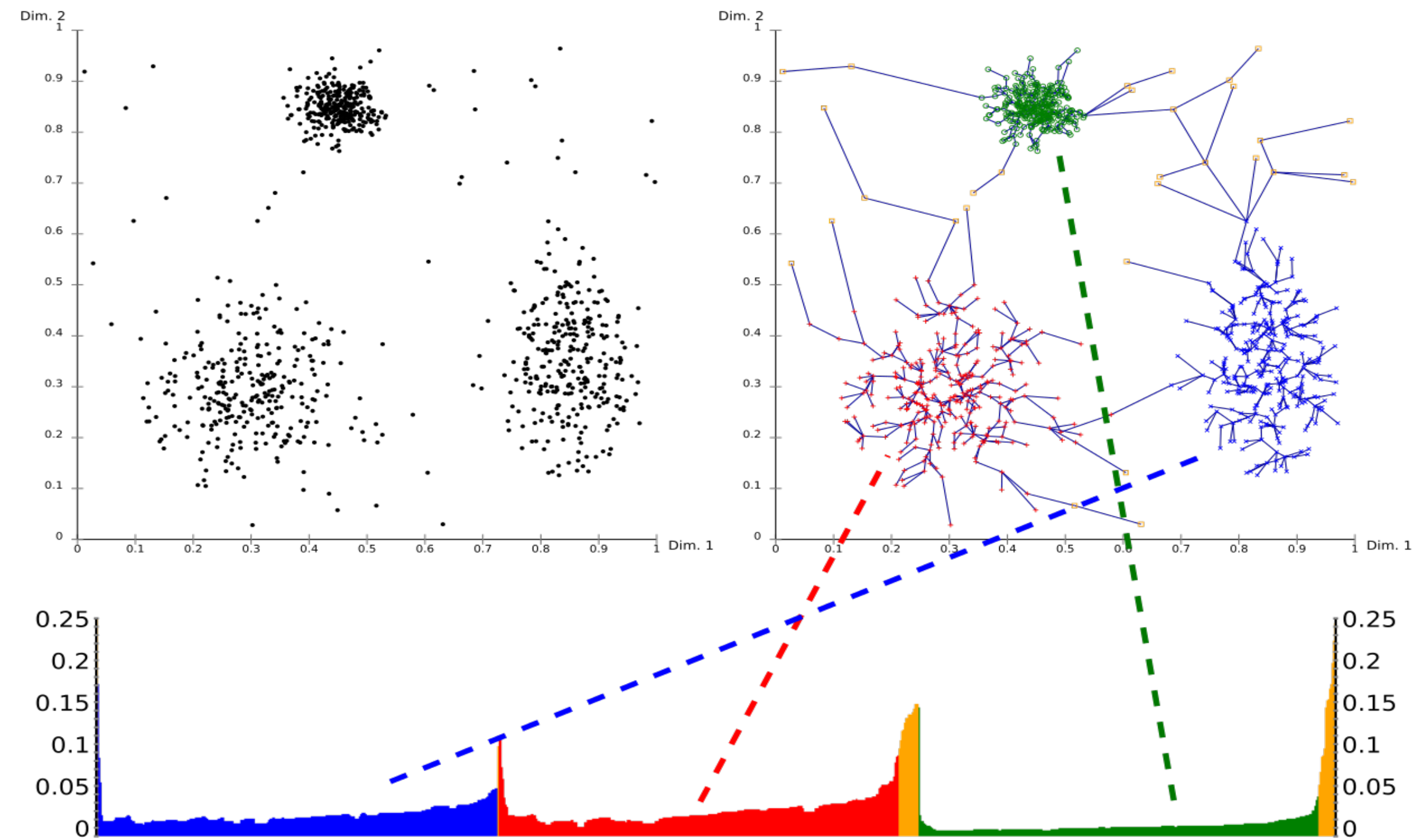
DBSCAN PROBLEMS WITH VARYING DENSITY



OPTICS

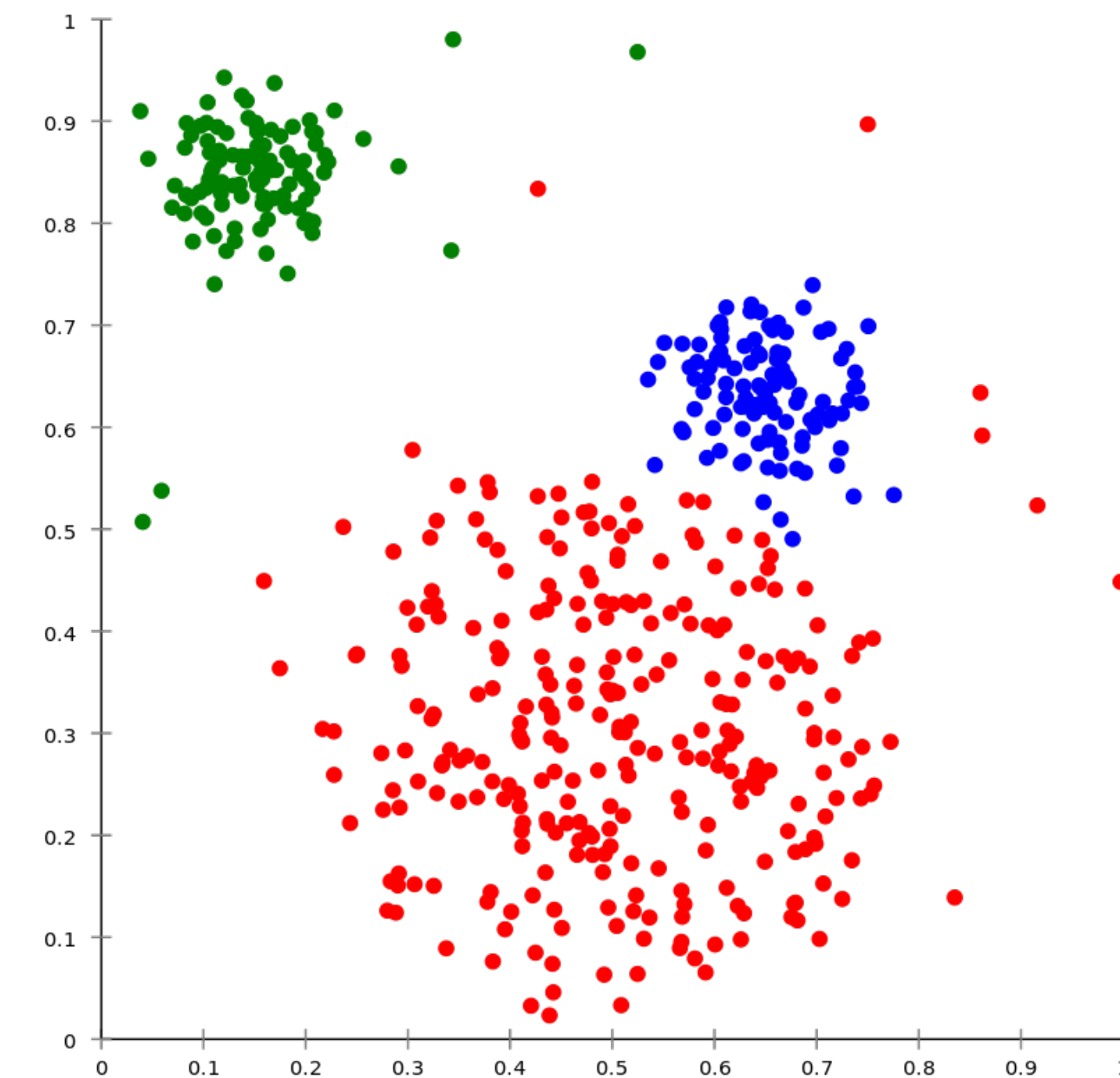
ORDERING POINTS TO IDENTIFY THE CLUSTERING STRUCTURE



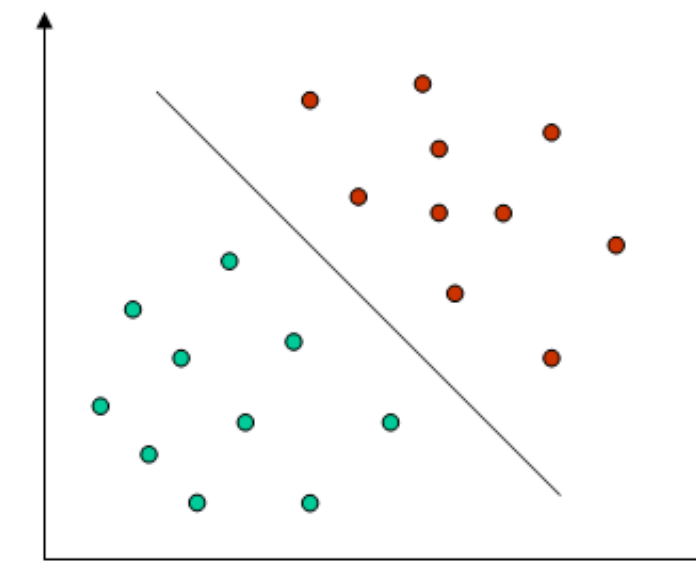


OPTICS

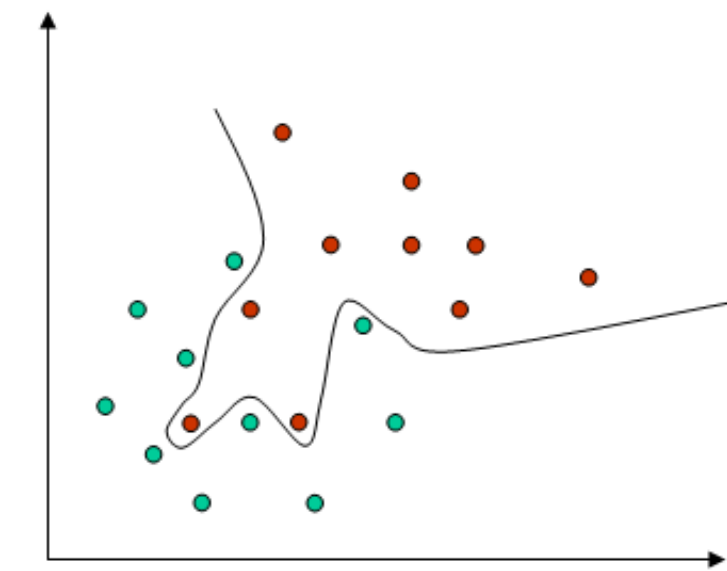
PROBLEM WITH LINEARITY OF PARTITION



SUPPORT VECTOR MACHINE (SVM) MAXIMIZE BORDERS

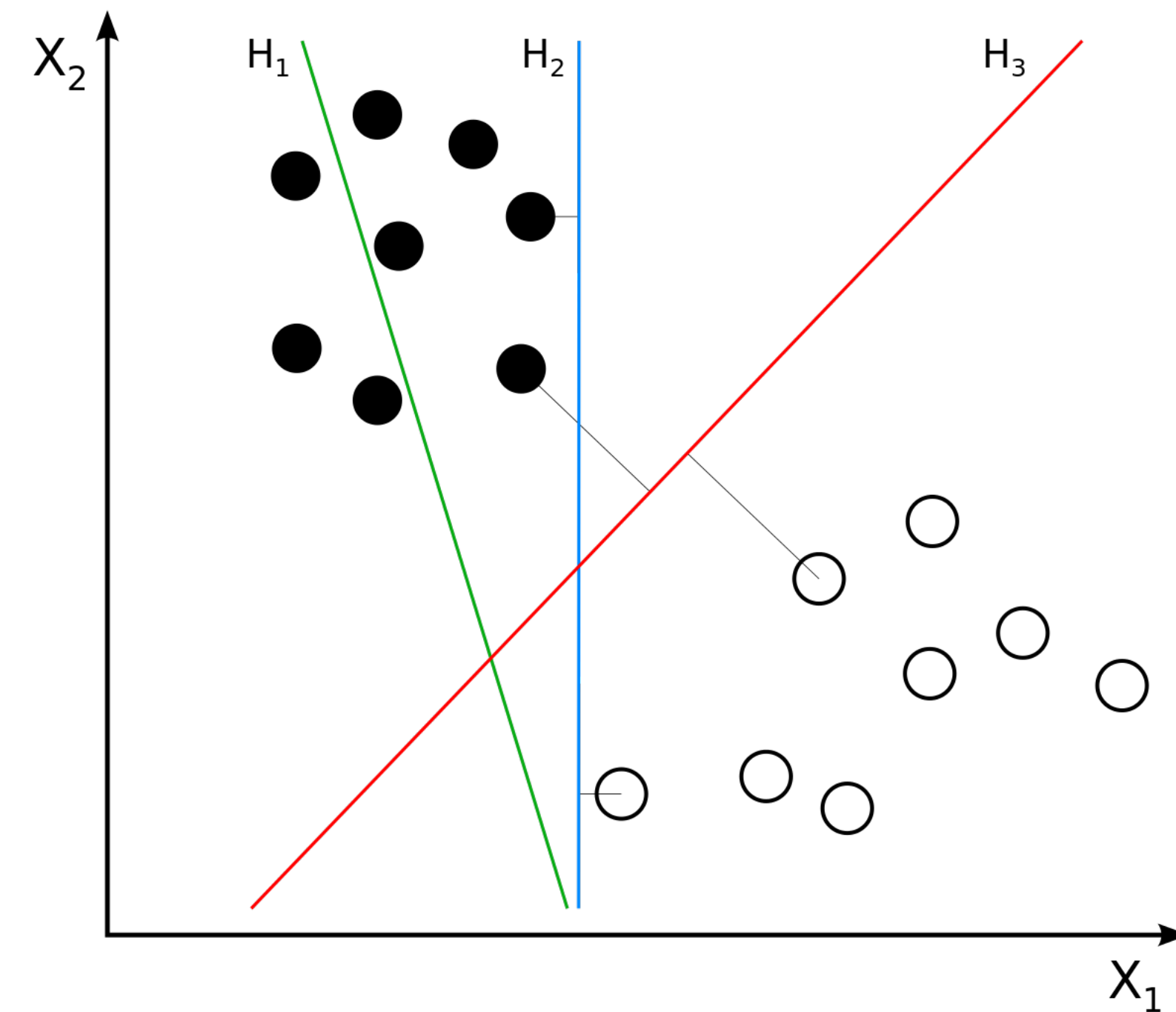


linear trennbar

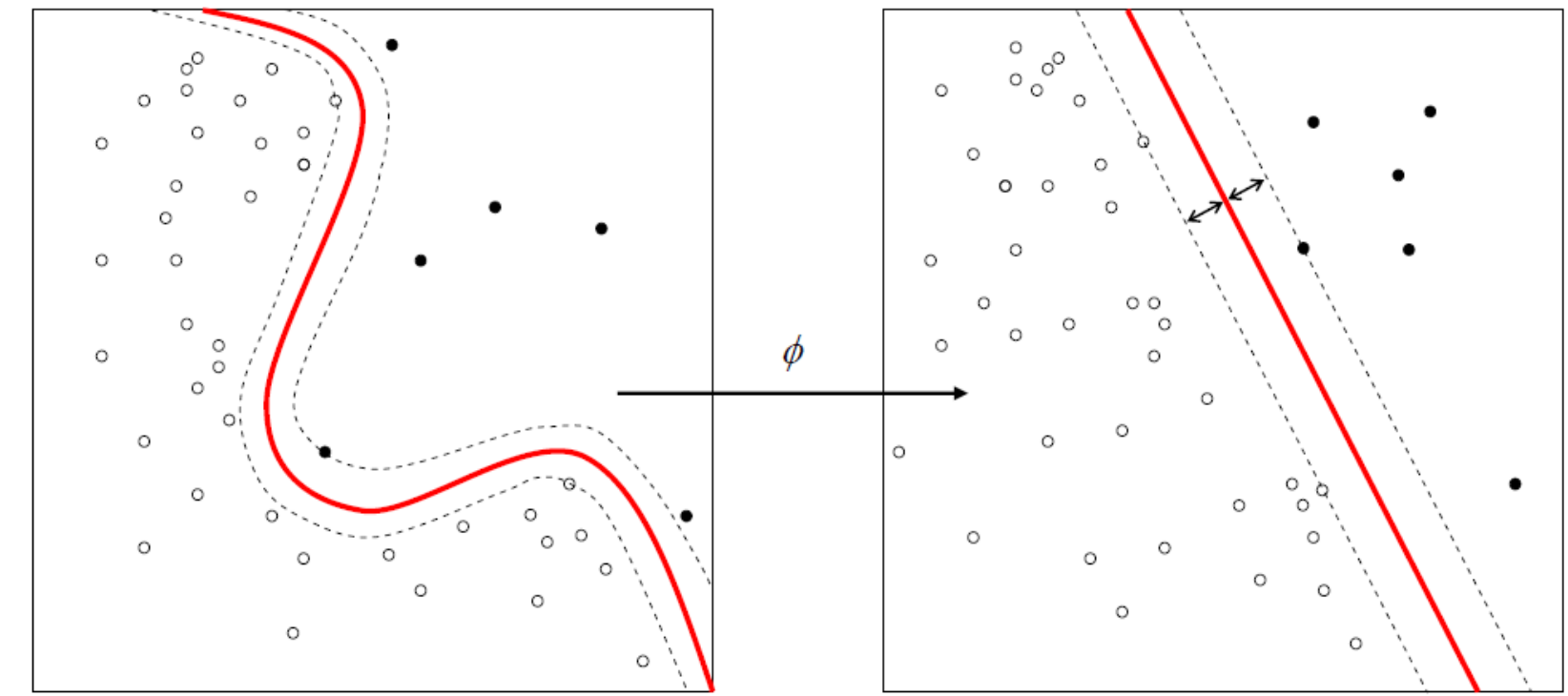


nicht linear trennbar

MAXIMIZE DISTANCE TO SEPARATING HYPERPLANES



KERNEL TRICK MAP TO HIGHER DIMENSIONS AND BACK



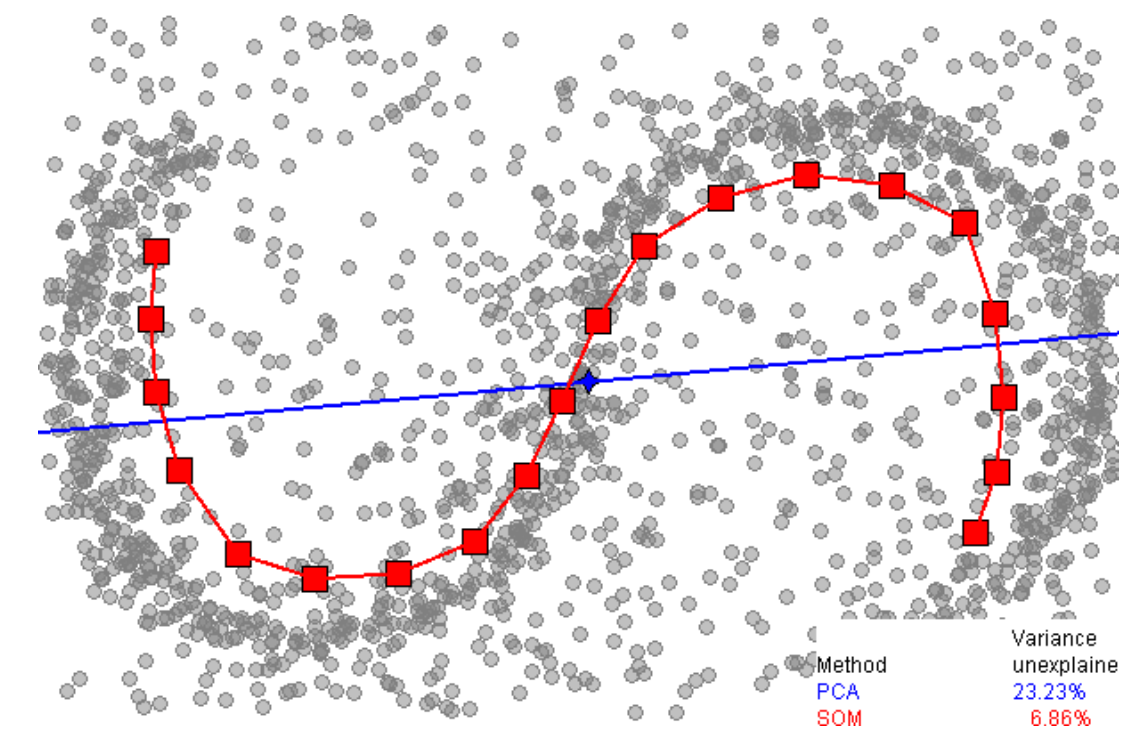
SVM

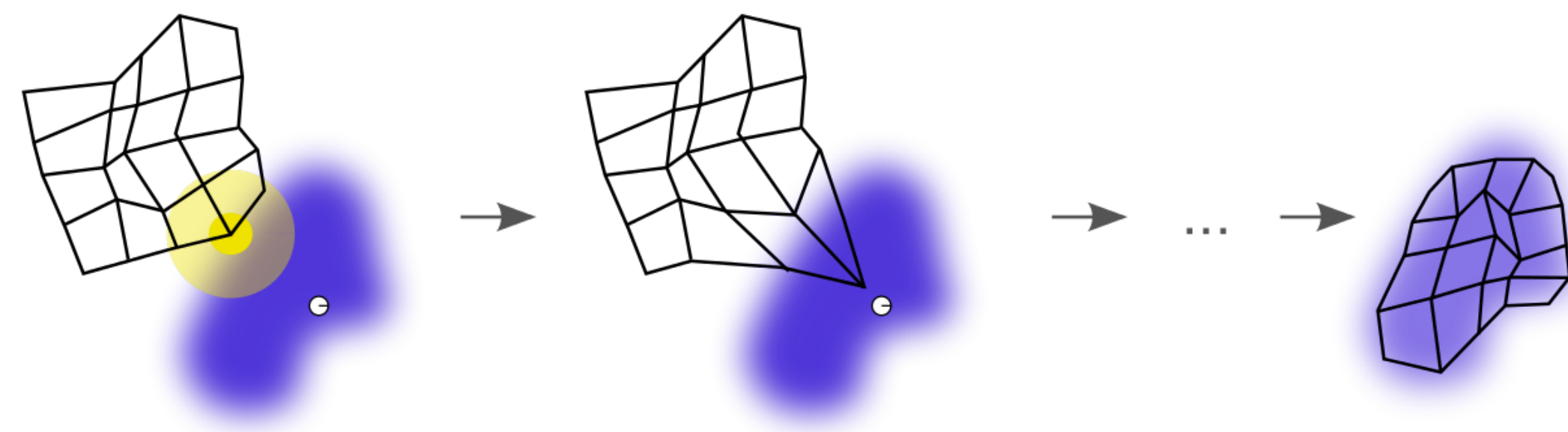
„PROBLEM“: SUPERVISED APPROACH

FIT A STRUCTURE OR FIT THE FEATURES

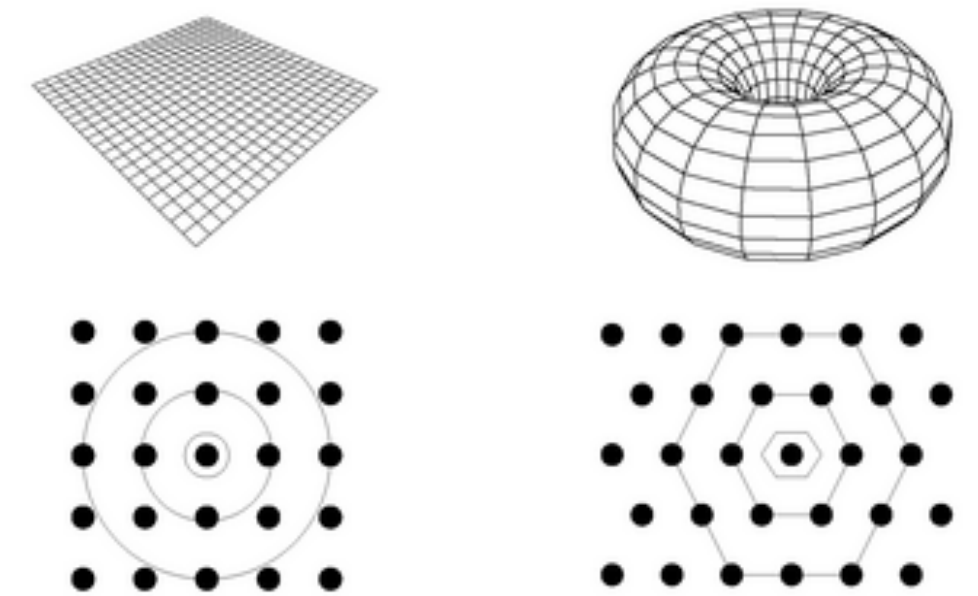


SELF-ORGANIZING MAP (SOM) NEURAL NETWORK APPROACH

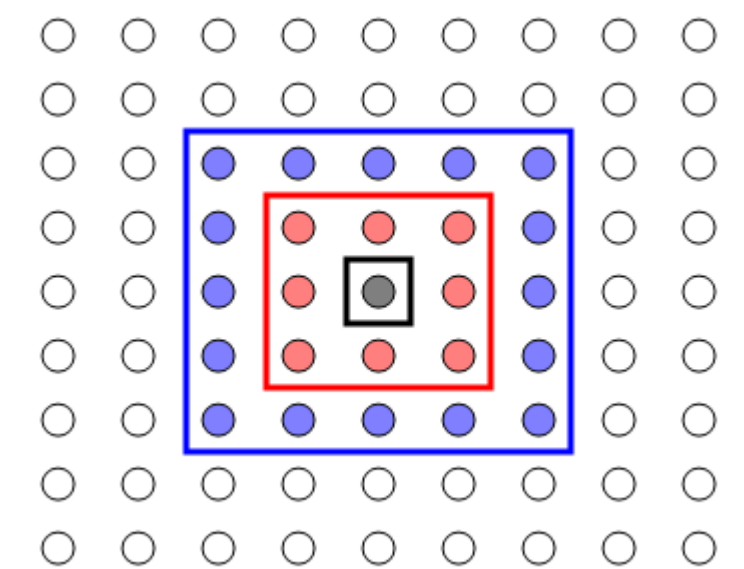
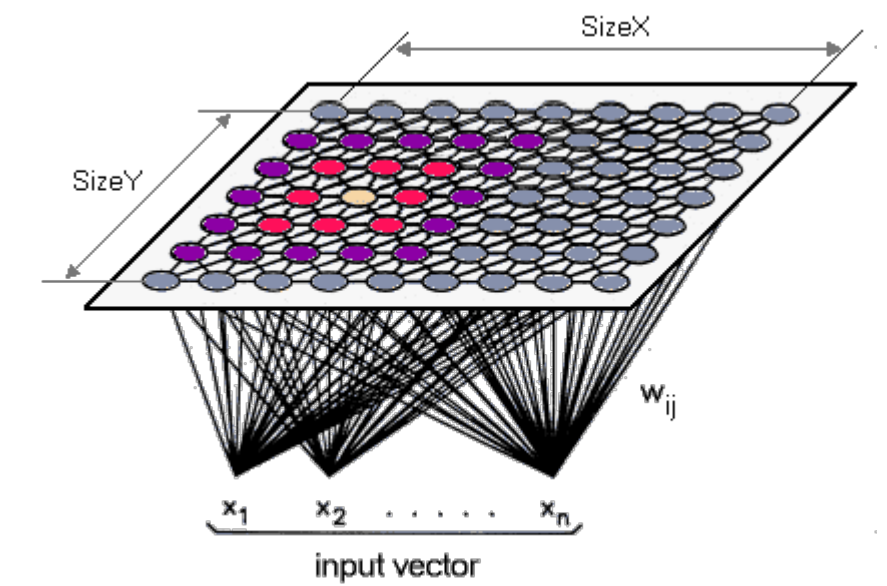




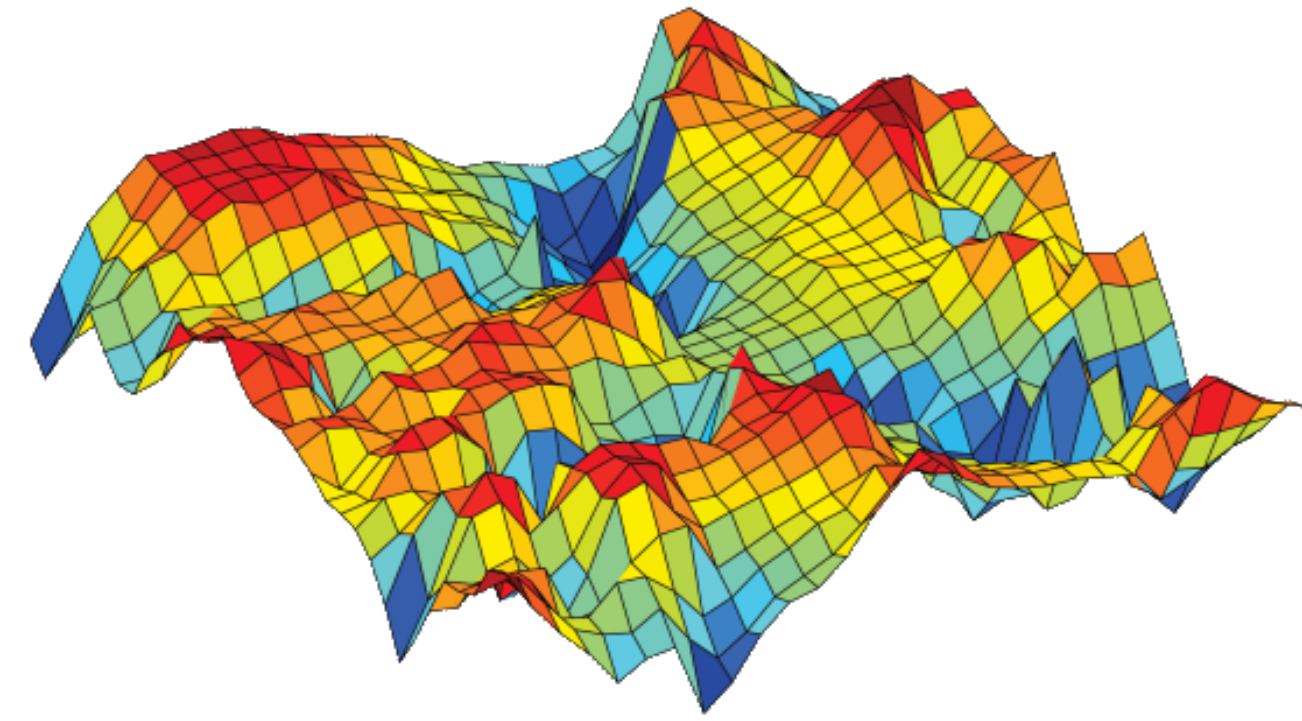
SELF-ORGANIZING MAP (SOM) FIXED TOPOLOGY ONLY FEATURES ARE CHANGING



SELF-ORGANIZING MAP (SOM) ITERATIVE TRAINING



SELF-ORGANIZING MAP (SOM) DISTANCES AS TOPOLOGY, BUT REGULAR GRID IN 2D



SELF-ORGANIZING MAP (SOM) ADVANTAGE: READABILITY

ONTOLOGIC
DISCUSSION
OF SOM.
INTRINSIC VS.
CONTEXTUAL
MEANING

EIGENFACES

KIND OF PCA VIA
EIGENVECTORS



EIGENFACES
A FEW FACES ARE
ENOUGH
TO IDENTIFY ALL
OTHERS,
WITHOUT PARAMETERS

Schedule

Mondays 10:00 - 12:00
051-0726-16L | 2 ECTS*

Creative Data Mining Intuitively Analysing Design Ideas

The goal of this course is to introduce various data mining techniques for design and urban planning applications. Students will learn how to select relevant data sources and collect their own data using a “sensor backpack”. Various methods will be applied to a common project to evaluate the predominant influencing factors of the urban environment on our perceptual experiences. A select neighborhood in the city will be used as a case study. Final results will be presented in the last class.

The course will start with an initial overview to data mining and the relevant mathematics as well as an introduction to the programming tool (RStudio). Then students will learn how to use and interpret results from a machine-learning tool to cluster self-made design sketches, which automatically generate qualitative collages. Finally, students will collect data using a “sensor backpack” with environmental sensors such as noise, temperature, illuminance, and air particulates. Students will also generate the data for perceptual quality in this neighborhood through time-stamped and geo-referenced surveys and biofeedback wristbands. Students will be given a work-flow to collect, process, analyze and interpret this data which may be used in their final projects.

Where
HIT H 12

Supervision
Danielle Griego
Matthias Standfest

griego@arch.ethz.ch
standfest@arch.ethz.ch

22.02.2016 Course Introduction
Introduce data-mining techniques and case study

29.02.2016 Introduction to the Environment
Introduction to R Studio and clustering

07.03.2016 From analog to digital analysis
Use hand-drawn sketched to auto-generated collages

14.03.2016 Seminar week (No lecture)

21.03.2016 Analysis and interpretation I
Evaluate auto-generated collages

28.03.2016 Holiday (No lecture)

04.04.2016 Time-series data analysis and Urban Planning
Introduction to time-series analysis

11.04.2016 Data collection with sensor backpack
Collect data and introduce workflows

18.04.2016 Holiday (No lecture)

25.04.2016 Analysis and interpretation II
Evaluate sensor backpack data

02.05.2016 Q&A Feedback Workshop
Finalise semester projects

09.05.2016 Final iA critique
Combined critique with the other iA courses
(14:00 - 16:00)

Requirement Former knowledge of any digital tool or coding language is most welcome but NOT required. You only need to provide a reasonable amount of motivation and of course a notebook.

*** Total 60 h = 2 ECTS**

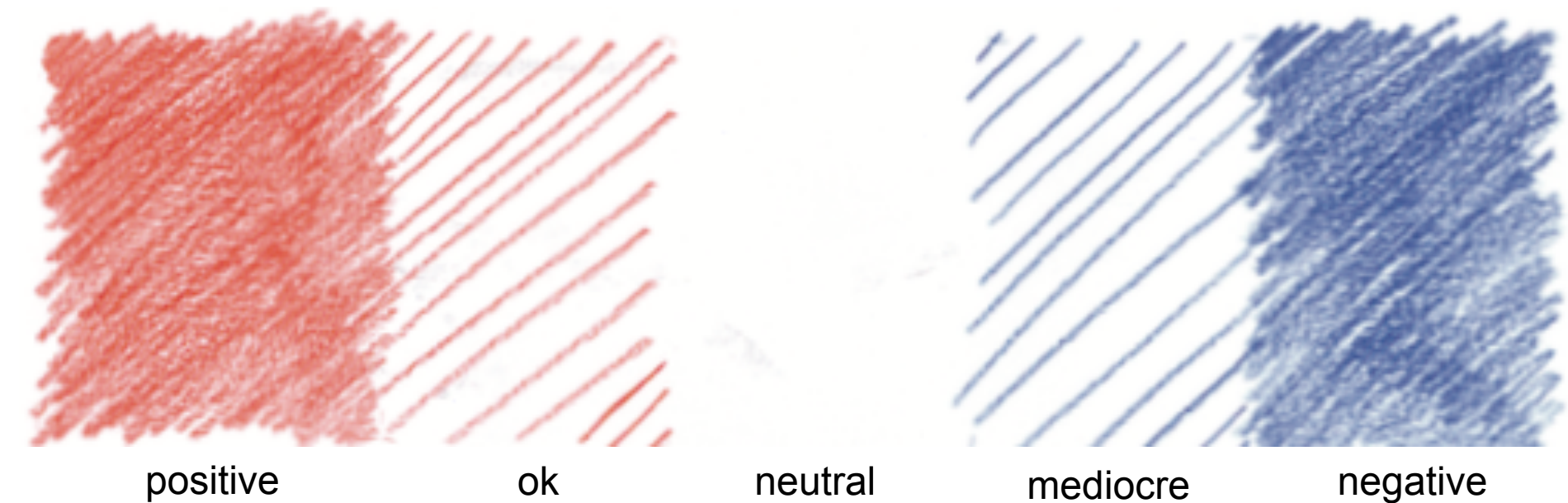
Exercises 40% (documentations)
Final Presentation 40% (Final project)
Attendance 20%

The most recent outline will be found on www.ia.arch.ethz.ch

Homework:

Part 1:

1. Color the segment maps of Alt-Wiedikon we provided today according to your perception of the urban space from a 2D plan. Hand in the hardcopies by Friday 18 March before 5pm.
2. Keep the following in mind:
 1. Use the shaded diagram below as a guide,
 2. Keep this task simple and work intuitively
 3. Be consistent for all 9 plans
 4. Reference Google satellite image to better understand the actual urban layout.



Homework:

Part 2:

1. Review the R-tutorials lecture_2_05 through lecture_2_06
2. Use a different built in dataset such as (mtcars, rock, trees, LifeCycleSavings) and visualize the clustering analysis using the improved visualization from lecture_2_06 tutorial
3. Submit a pdf of the final image by Monday 21 March