



UNIVERSIDAD
DE BURGOS

Sistemas Inteligentes
Curso 2023-2024

Reconocimiento de patrones. Fundamentos.

Curso 2024-2025

Pedro Latorre Carmona
Dept. Ingeniería Informática – UBU

Outline

- Introduction
- Classifiers based on Bayes decision theory
- Support Vector Machines (SVMs)
- Clustering
- *Houston, we have a few problems!*

Introduction

Introduction

How can we define what **pattern recognition** is?

From Wikipedia:

“**Pattern recognition** is a branch of **machine learning** that focuses on the recognition of patterns and regularities in data, **although it is in some cases considered to be nearly synonymous with machine learning.**

Pattern recognition systems are in many cases trained from labeled **training** data (supervised learning), but when no labeled data are available other algorithms can be used to discover previously unknown patterns (unsupervised learning).

The terms **pattern recognition**, **machine learning**, **data mining** and **knowledge discovery in databases** (KDD) are hard to separate...”.

Introduction

What is a **pattern**?

A series of measurements whose combination will define the “object under study”.

What are the steps in a “pattern recognition” system?

We can consider the following ones:



Introduction

There are different *taxonomies* for Pattern recognition.

One of them is based on whether we have information about the class the data samples belong to or not.

A pattern recognition method is called:

- *Supervised*: if we have access to the labels (class information) associated to the data samples.
- *Unsupervised*: if this information is not available.

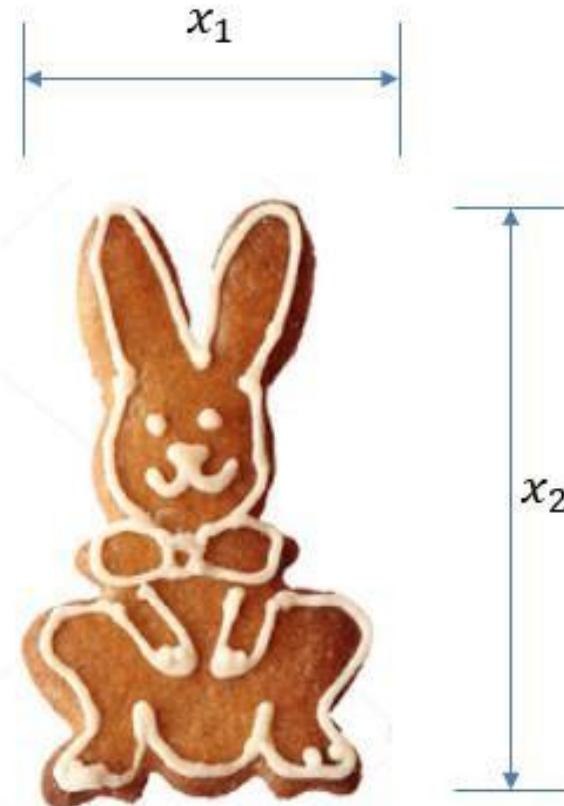
Introduction

Imagine we have the following situation. We need to distinguish two types of animal crackers for food: *rabbits* from *sheeps*.



Introduction

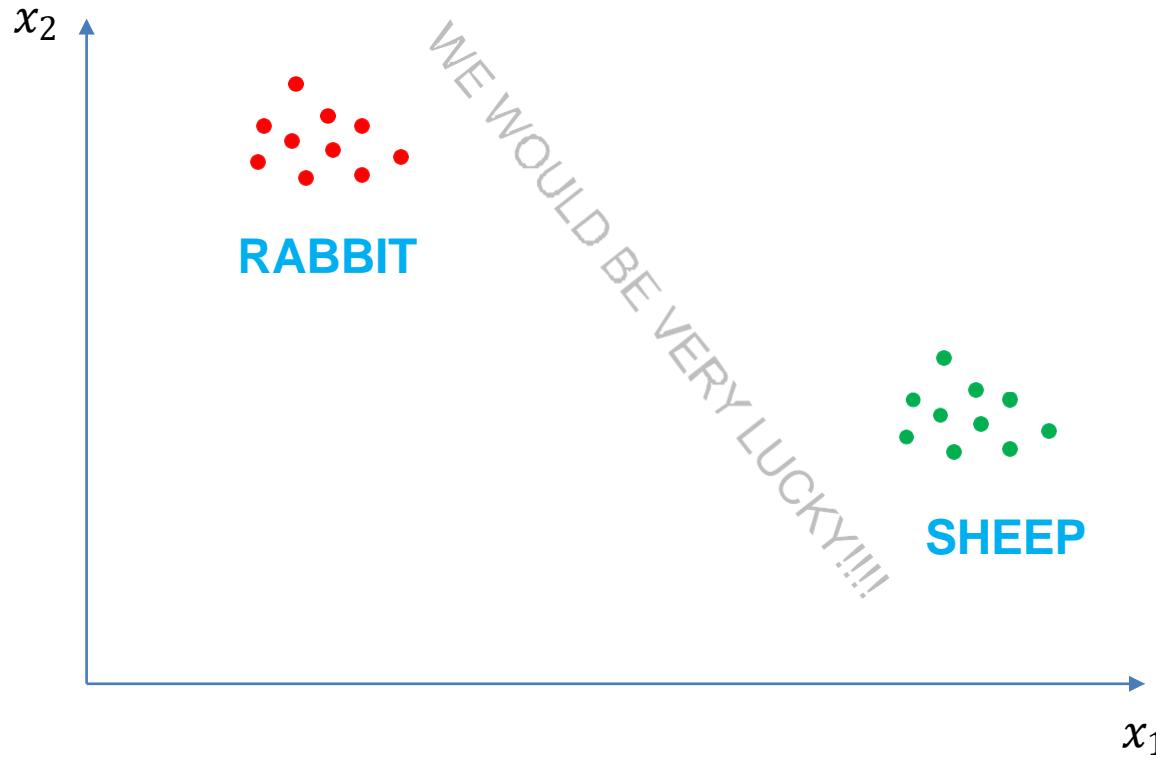
We can define the following two (simple) features:



and define $X = (x_1, x_2)$ as a **feature vector**.

Introduction

What if we do the following plot and find...?



Introduction

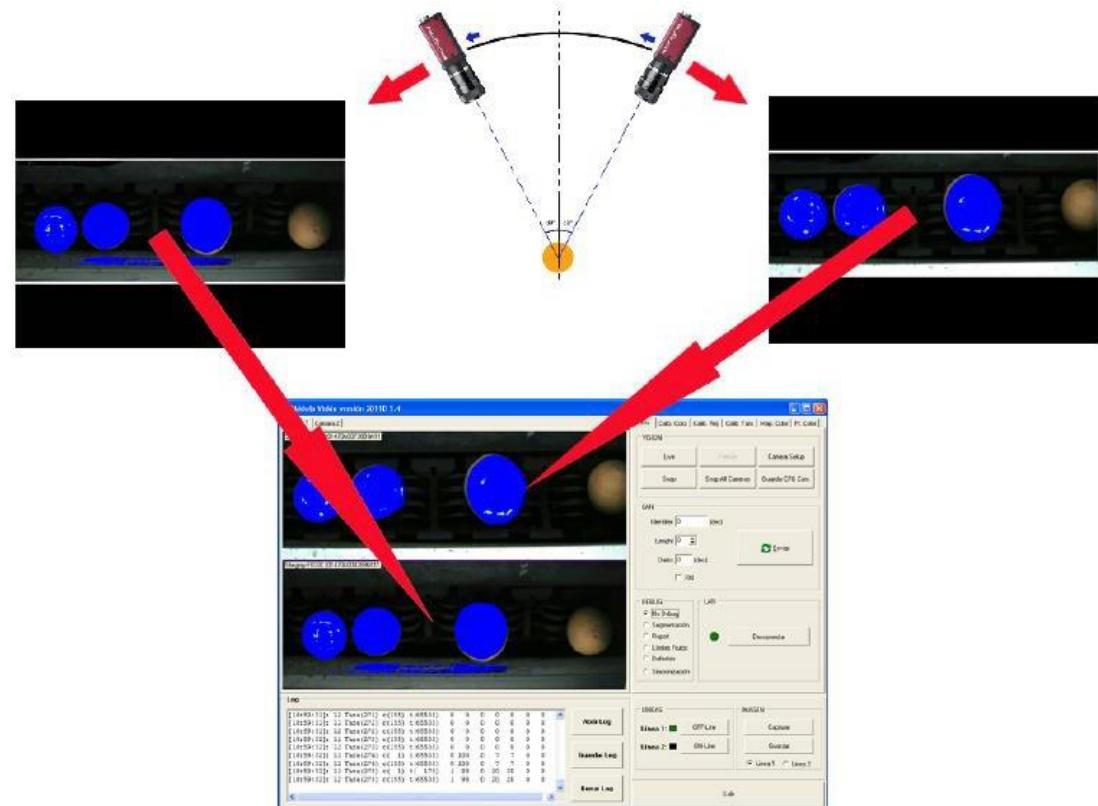
Imagine an on-line fruit inspection/sorting system:



Introduction

Imagine an on-line fruit inspection/sorting system:

We could “infer”
a *center* and
a *radius* of/in the fruit



The feature vector could be:

$$x_i \equiv (c_i, r_i);$$

$$c_i \equiv (x_i, y_i)$$

Introduction

A *supervised* pattern recognition dataset is formed by (X_i, y_i) pairs, where:

$X_i \equiv \{x_{1i}, x_{2i}, x_{3i}, \dots, x_{di}\}$;

$y_i \equiv$ Class the data sample i belongs to.

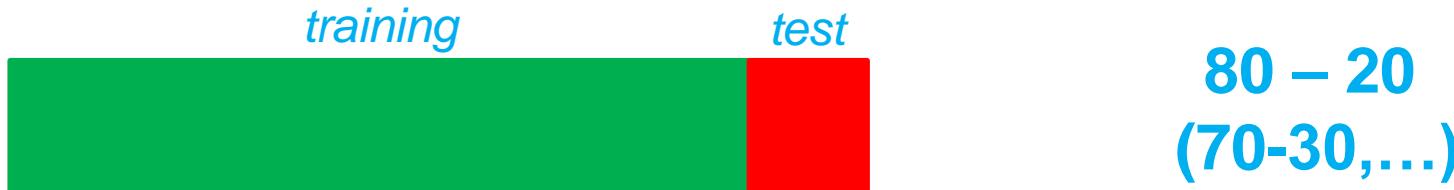
where d represents the so-called problem **DIMENSIONALITY**.

Any pattern recognition dataset is divided into:

- A *Training* dataset
- A *Test* dataset

Introduction

A complete dataset must **ALWAYS** be divided into a *training* and a *test* dataset, and **size(*training_set*) > size(*test_set*)**



Introduction

How can we measure classification accuracy?

$$X \equiv \{x_1, x_2, x_3, \dots, x_N\}; \quad Y = \{y_1, y_2, \dots, y_N\} \in \{+1, -1\}$$

	“Illness”	No “illness”
Classified as “Illness”	True Positives (TP)	False Positives (FP)
Classified as “No illness”	False Negatives (FN)	True Negatives (TN)

$$Acc\ (\%) = \left(\frac{TP+TN}{TP+TN+FP+FN} \right) \cdot 100$$

Introduction

What if there are more than **2** classes?

$$X \equiv \{x_1, x_2, x_3, \dots, x_N\}; \quad Y = \{y_1, y_2, \dots, y_N\} \in \{1, 2, 3, \dots\}$$

$$i = 1 \quad i = 2 \quad i = 3 \quad \dots \quad i = l$$

1 **2 3 4 5 6 ...** **2** **1 3 4 5 6 ...** **3** **1 2 4 5 6**

$$\overline{Acc} = \frac{\sum_{i=1}^l \left(\frac{TP_i + TN_i}{TP_i + TN_i + FP_i + FN_i} \right)}{l}$$

Supervised classification (SC)

SC- Bayes decision rules

Classifiers based on **Bayes** *decision theory*

SC- Bayes decision rules

Let's imagine the following situation:

A school is formed by 60% boys and 40% girls. We randomly pick up a student. What is the probability we chose a boy?

Answer: **$p(Boy) \equiv p(\omega = \omega_1) = 0.6$**

SC- Bayes decision rules

But....What if....?



We are told all boys (100%) wear pants and half girls (50%) wear pants and half (50%) wear skirts. A camera sees a student with pants, what is the probability that the student is a boy?

$$p(\text{Boy}|\text{Pants}) \equiv p(\omega_1|x); \quad p(\text{Pants}| \text{Boy}) \equiv p(x|\omega_1)$$

$$p(\text{Boy}|\text{Pants}) = \frac{p(\text{Boy}) \cdot p(\text{Pants}|\text{Boy})}{p(\text{Boy}) \cdot p(\text{Pants}|\text{Boy}) + p(\text{Girl}) \cdot p(\text{Pants}|\text{Girl})} = \frac{0.6 \cdot 1}{(0.6 \cdot 1) + (0.4 \cdot 0.5)} = 0.75$$

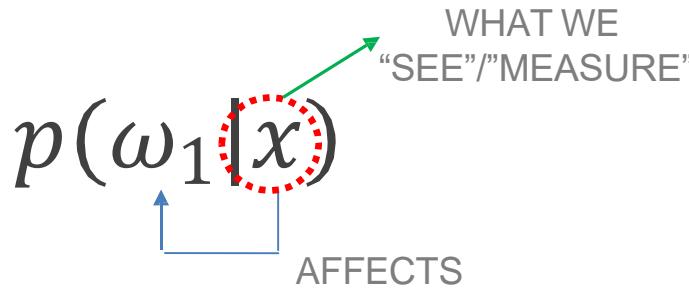
(Before): $p(\text{Boy}) \equiv p(\omega_1) = 0.6$

FURTHER information changes our ability to better estimate probabilities.

SC- Bayes decision rules

BEFORE (ANY MEASUREMENT): $p(\omega_1)$

AFTER (ANY MEASUREMENT): $p(\omega_1 | x)$



Therefore, Bayes rule is **COOL!**

SC- Bayes decision rules

Let us define a variable considering a label associated to each class (rabbit, sheep).

Rabbit: $\omega = \omega_1$

Sheep: $\omega = \omega_2$

The ***a priori probability*** reflects our knowledge about how likely are classes to appear, **BEFORE** any measurement is taken.

These probabilities are:

$$p(\omega_1); p(\omega_2)$$

SC- Bayes decision rules

In general, for **C** classes, we would have:

$$\sum_{i=1}^C p(\omega_i) = 1$$

What do you think about the following plan (2 classes for now)?:

If $p(\omega_1) > p(\omega_2)$ → **Decide** that the **sample** belongs to **class ω_1**

If $p(\omega_1) \leq p(\omega_2)$ → **Decide** that the **sample** belongs to **class ω_2**

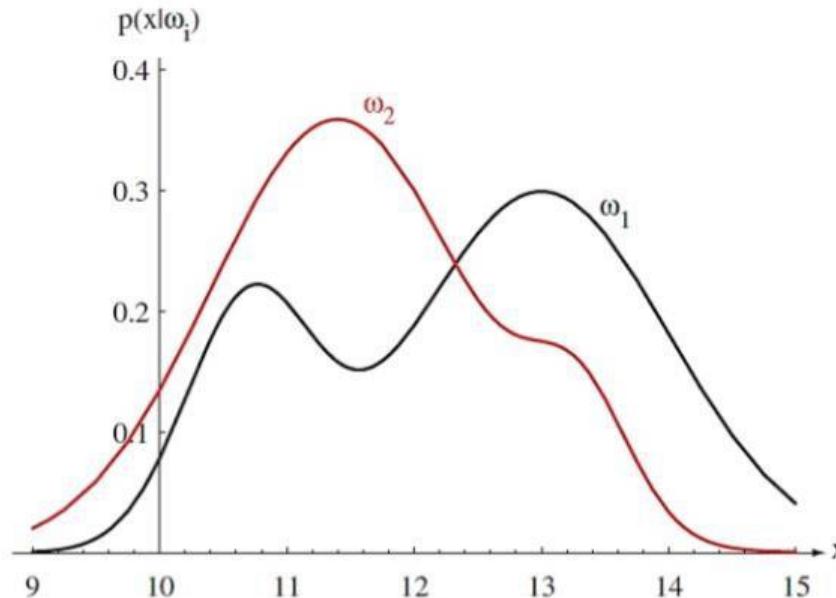
Answer: We would always be selecting the same class!!

SC- Bayes decision rules

Therefore, somehow, **we need to introduce the data into the model.**

A **feature** is a variable you can **measure**. A feature can be: (a) a scalar, $x \in \mathbb{R}$; (b) a vector, $x \in \mathbb{R}^d$

The **class-conditional probability density**, $p(x|\omega_i)$ reflects our **assumptions** about the **underlying distribution** for the data in the class.



SC- Bayes decision rules

What do we really need, then?

$$p(\omega_i | x)$$

How can we obtain it? Using **Bayes rule**:

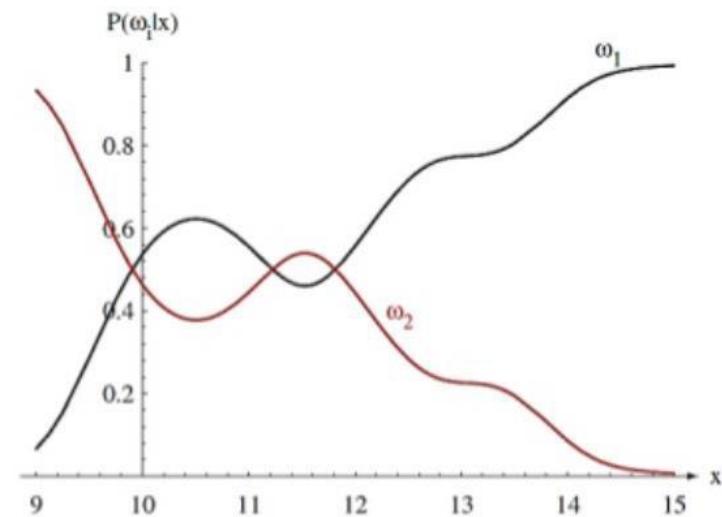
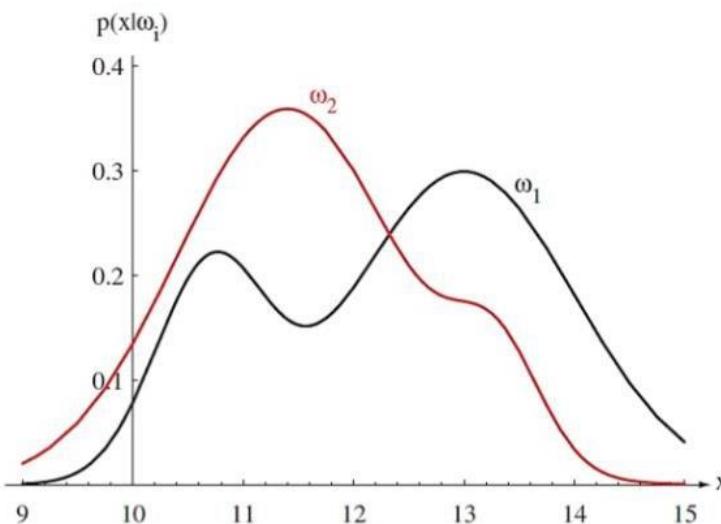
$$p(\omega_i, x) = p(x|\omega_i) \cdot p(\omega_i) = p(\omega_i|x) \cdot p(x) \Rightarrow$$

$$\Rightarrow p(\omega_i|x) = \frac{p(x|\omega_i) \cdot p(\omega_i)}{p(x)} = \boxed{\frac{p(x|\omega_i) \cdot p(\omega_i)}{\sum_{j=1}^C p(x|\omega_j) \cdot p(\omega_j)} = p(\omega_i|x)}$$

SC- Bayes decision rules

Then:

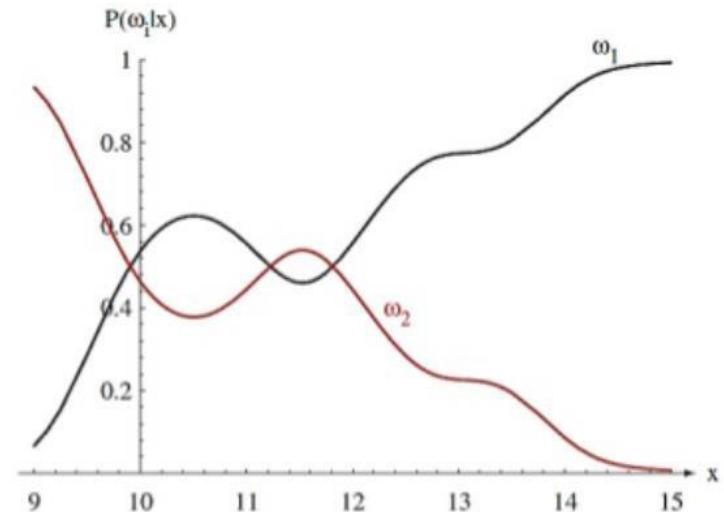
$$p(\omega_i|x) = \frac{p(x|\omega_i) \cdot p(\omega_i)}{\sum_{j=1}^c p(x|\omega_j) \cdot p(\omega_j)}$$



SC- Bayes decision rules

From the previous plot:

- If $p(\omega_1|x) > p(\omega_2|x) \rightarrow$
Sample x belongs to class ω_1
- If $p(\omega_2|x) > p(\omega_1|x) \rightarrow$
Sample x belongs to class ω_2



SC- Bayes decision rules

Therefore, the class would be:

$$\omega^*(\mathbf{x}) = \arg \max_i [p(\omega_i | \mathbf{x})]$$



Maximum A Posteriori (**MAP**)
criterion

The probability of error (for the two class situation) would be:

$$p(\varepsilon | \mathbf{x}) = \begin{cases} p(\omega_1 | \mathbf{x}) & \text{if we decide } \omega_2 \\ p(\omega_2 | \mathbf{x}) & \text{if we decide } \omega_1 \end{cases}$$

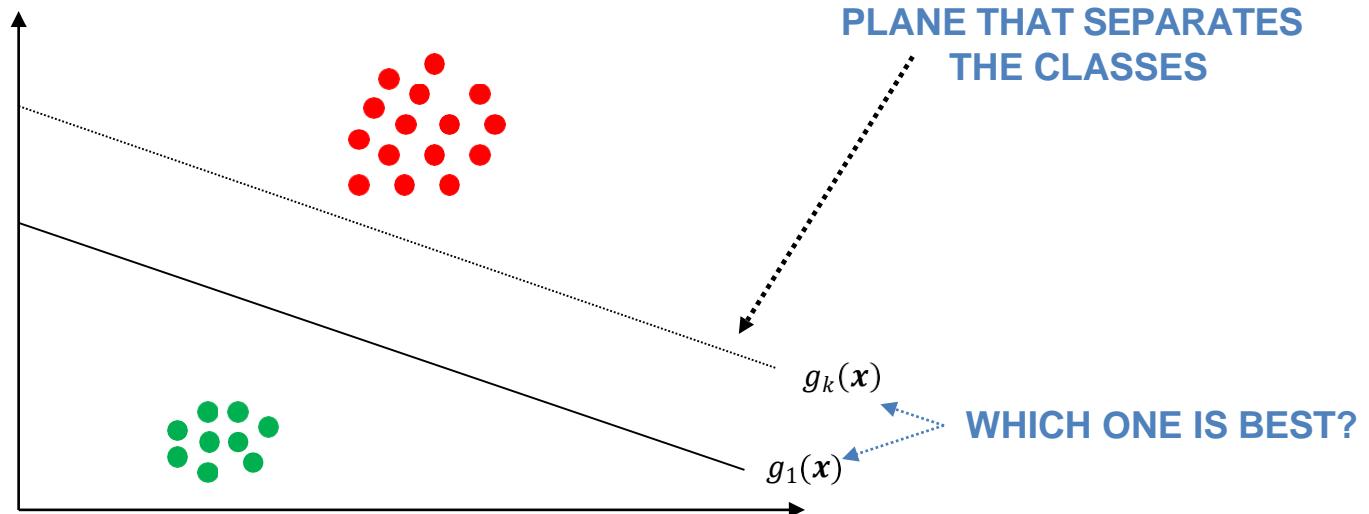
Support Vector Machines (SVMs)

SC (Linear)-SVMs

Let us consider a group of points $\{x_i, i = 1, \dots, N\}$ that belong to one of two classes $\{C_1, C_2\}$ assumed to be linearly separable by the hyperplane:

$$g(\mathbf{x}) \equiv \boldsymbol{\omega}^T \mathbf{x} + \omega_0 = 0$$

We may have many solutions $\{g_1(\mathbf{x}), \dots, g_k(\mathbf{x})\}$



SC (Linear)-SVMs

The distance between a point and a line should be given by:

In the case of a line in the plane given by the equation $ax + by + c = 0$, where a, b and c are **real** constants with a and b not both zero, the distance from the line to a point (x_0, y_0) is^{[1][2]: p.14}

$$\text{distance}(ax + by + c = 0, (x_0, y_0)) = \frac{|ax_0 + by_0 + c|}{\sqrt{a^2 + b^2}}.$$

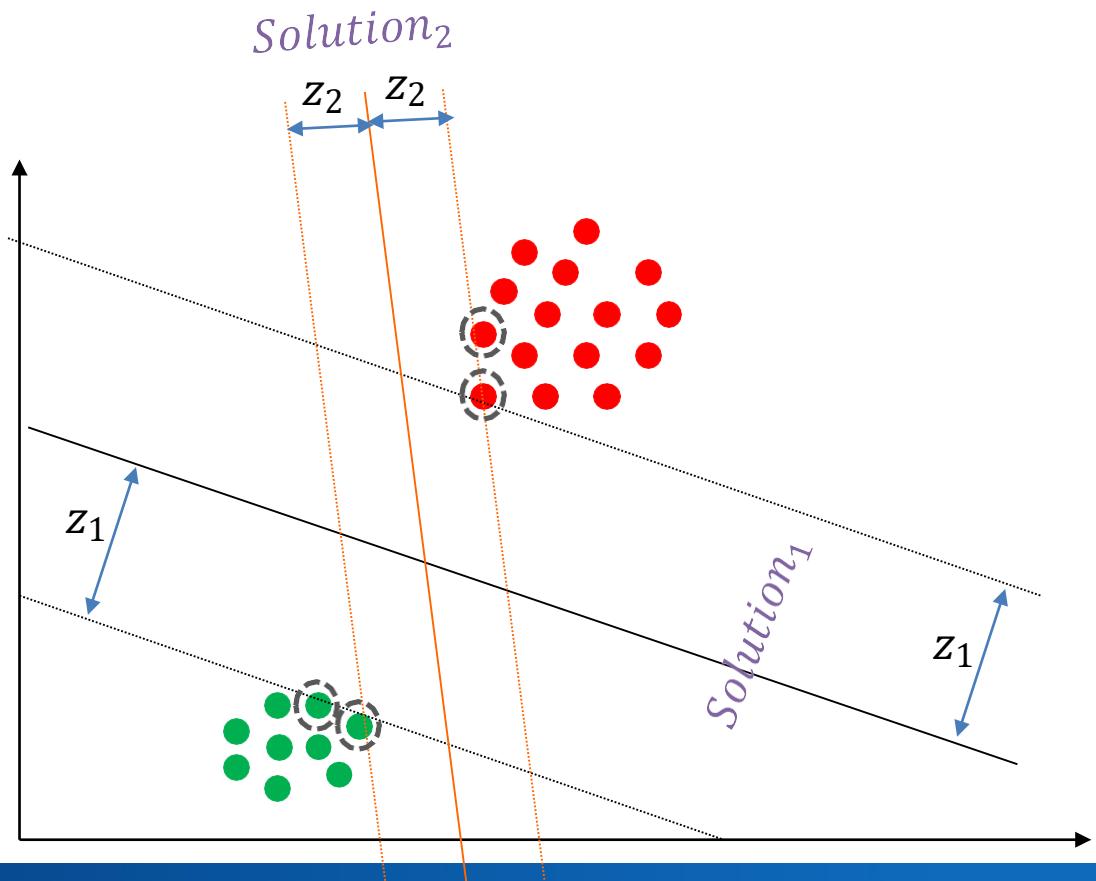
https://en.wikipedia.org/wiki/Distance_from_a_point_to_a_line#:~:text=In%20Euclidean%20geometry%2C%20the%20distance,nearest%20point%20on%20the%20line.

SC (Linear)-SVMs

Our aim would be to obtain $g(x) \equiv \omega^T x + \omega_0$ that maximizes z :

$$\omega = (\omega_1, \omega_2)$$

$$z = \frac{|g(x)|}{\|\omega\|} = \frac{|g(x)|}{\sqrt{\omega_1^2 + \omega_2^2}}$$



SC (Linear)-SVMs

This would be equivalent to:

$$\text{maximize } \left(\frac{2}{\|\omega\|} \right)$$

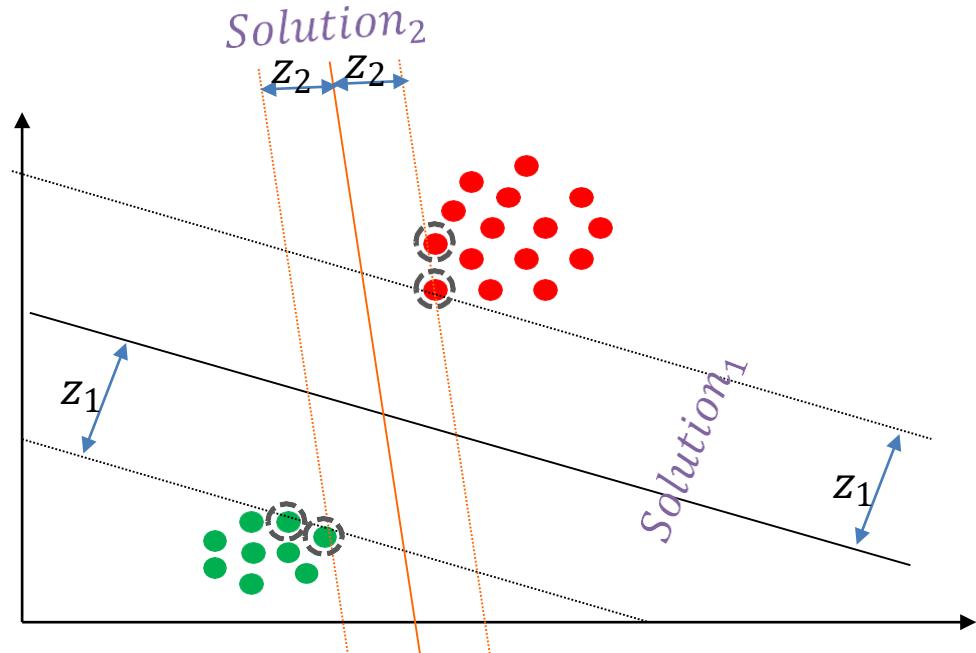
Or

$$\text{minimize} \left(\frac{1}{2} \|\omega\|^2 \right)$$

subject to:

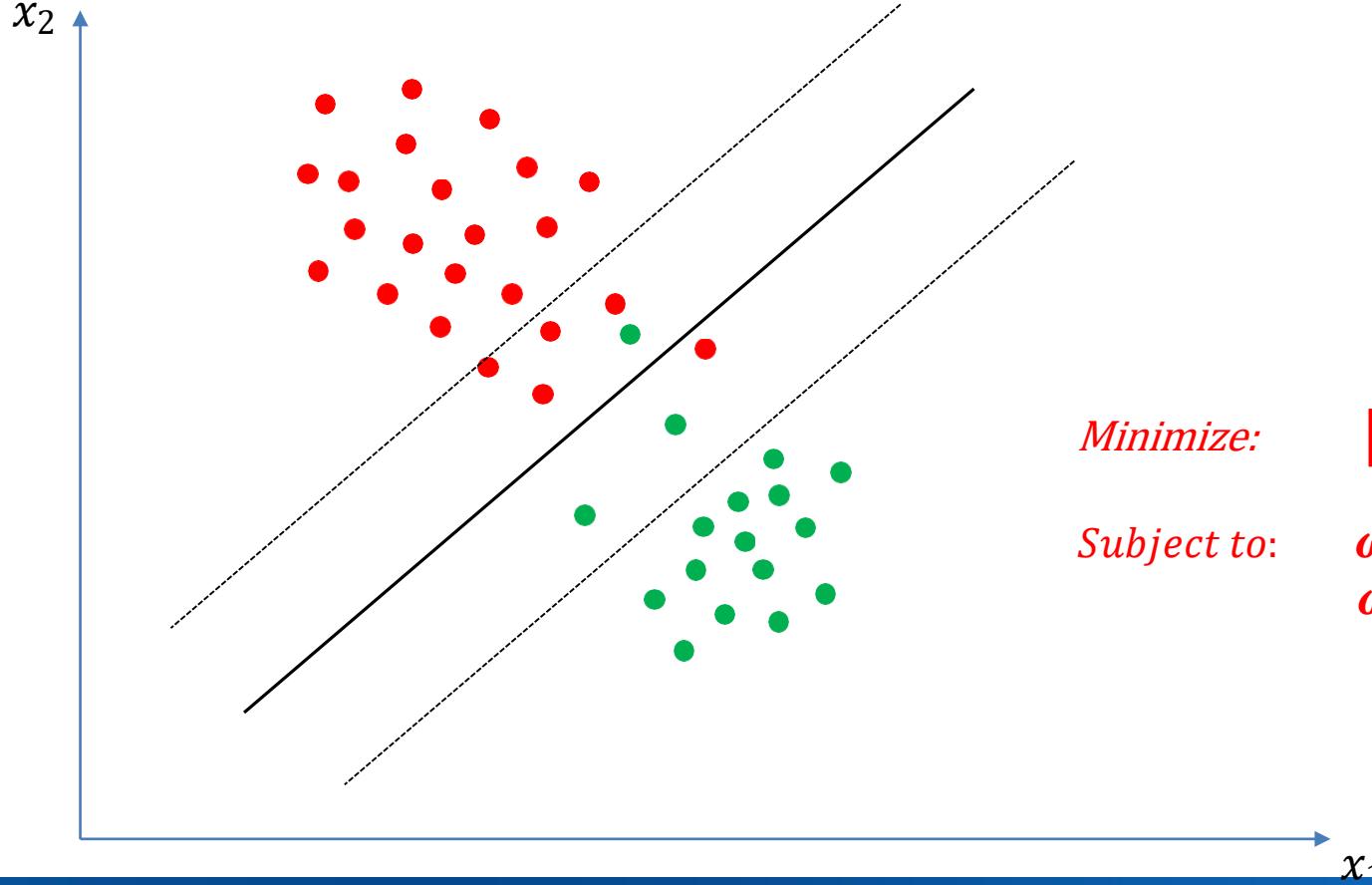
$$\omega^T x + \omega_0 \geq 1 \quad \forall x \in C_1$$

$$\omega^T x + \omega_0 \leq -1 \quad \forall x \in C_2$$



SC (Linear)-SVMs

What if we have unseparable classes? We allow for **some (N) errors!**



Minimize:

$$\left[\frac{1}{2} \|\omega\|^2 + C \cdot \sum_{i=1}^N \xi_i \right]$$

Subject to:

$$\begin{aligned}\omega^T x_i + \omega_0 &\geq 1 - \xi_i \\ \omega^T x_i + \omega_0 &\leq -1 + \xi_i\end{aligned}$$

SC (Linear)-SVMs

It can be shown that:

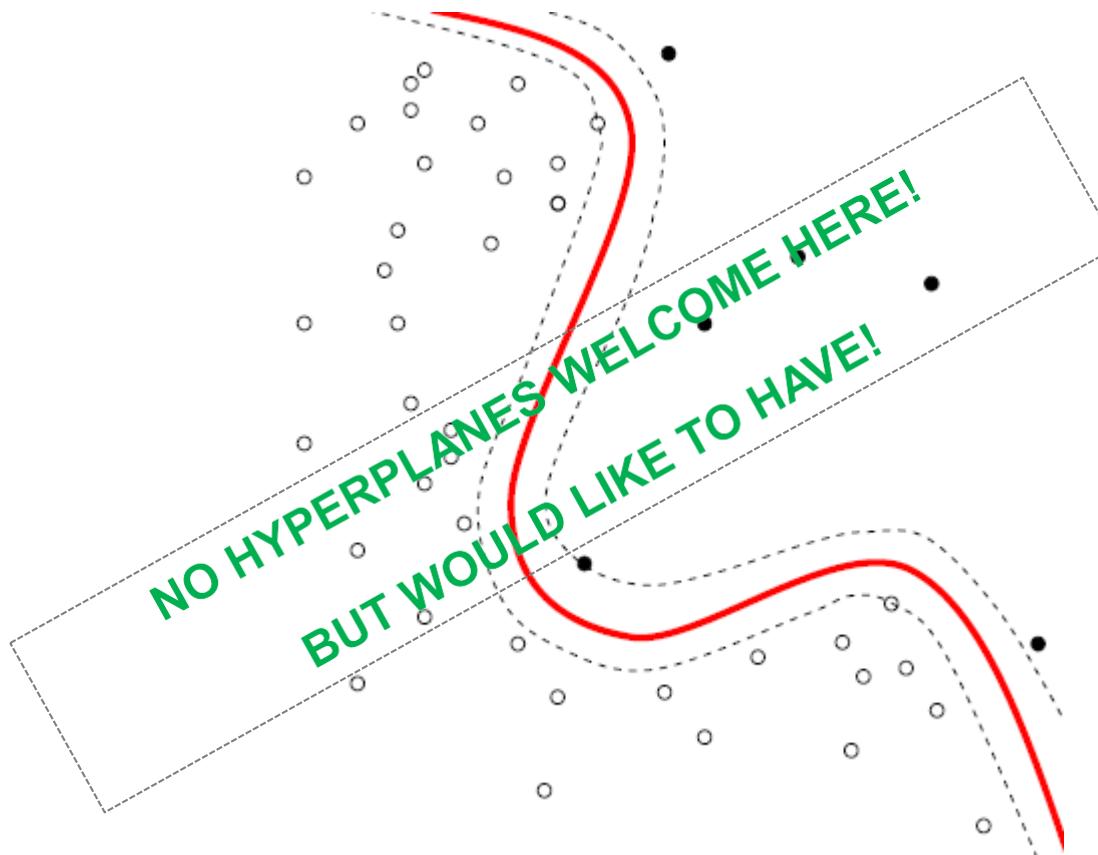
$$\omega = \sum_{i=1}^{N_{SV}} \lambda_i \cdot y_i \cdot x_i$$

where: $y_i \in \{-1, +1\}$ and N_{SV} is the number of training samples that define the hyperplane slope.

This is a CAPITAL result because it says the solution is a linear combination of some training samples, called SUPPORT VECTORS (SV) (and therefore the NAME!).

SC (Non linear)-SVMs

How could we achieve the following (needed by the data features)?



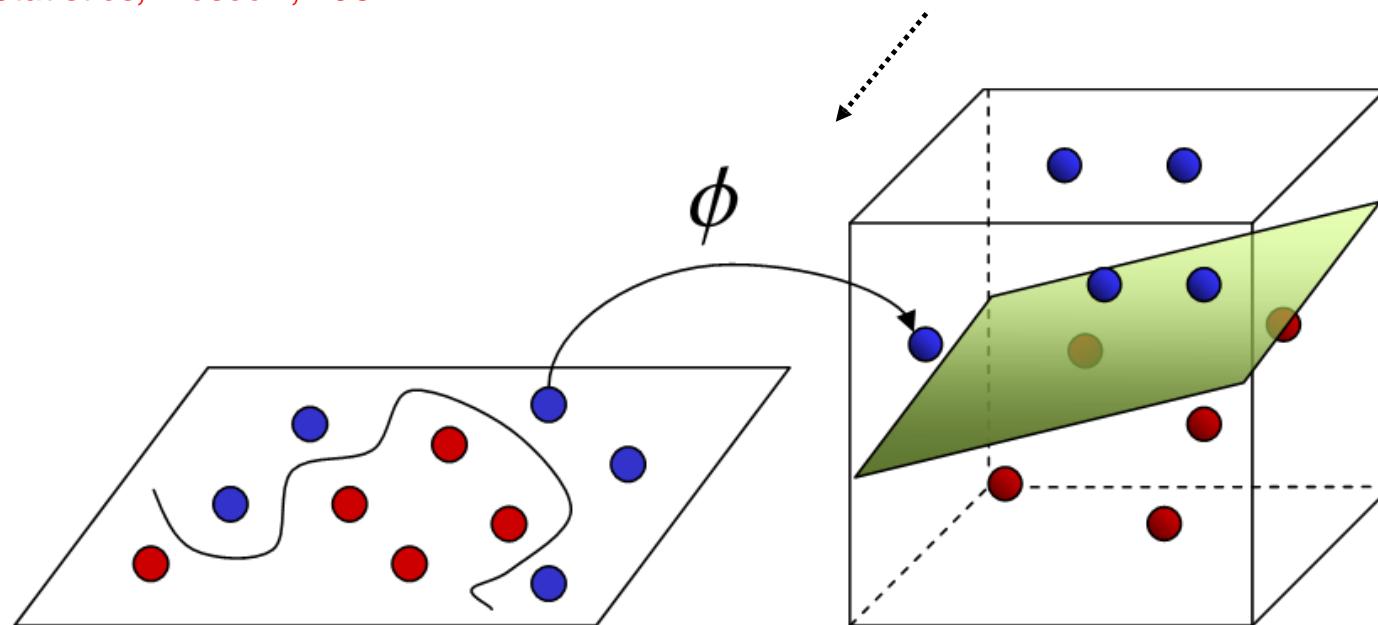
SC (Non linear)-SVMs

Vladimir Vapnik

(1936 -)

PhD Statistics, Moscow, 1964.

INCREASE THE NUMBER OF DIMENSIONS
OF THE SPACE TO MAKE THE DATA POINTS
LINEARLY SEPARABLE



Input Space

Feature Space

SC (Non linear)-SVMs

In a similar way it can be shown that:

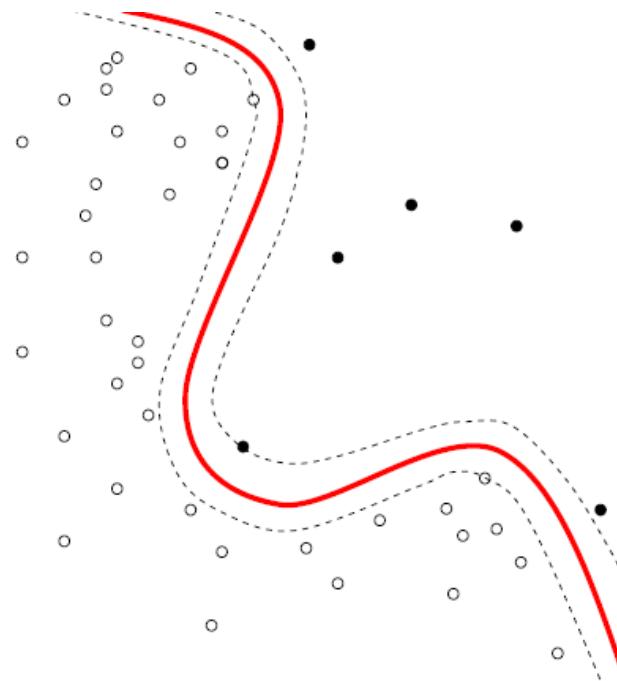
$$g(x) = \sum_{i=1}^{N_{SV}} \lambda_i \cdot y_i \cdot K(x_i, x) + \omega_0$$

Kernel

where :

$$y_i \in \{-1, +1\}$$

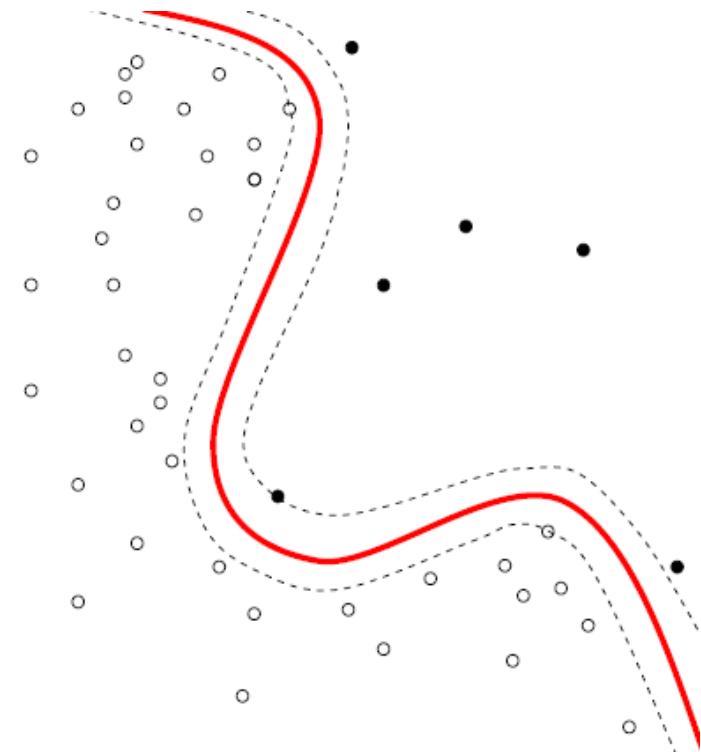
N_{SV} is the number of **training** samples that the solution considers.



SC (Non linear)-SVMs

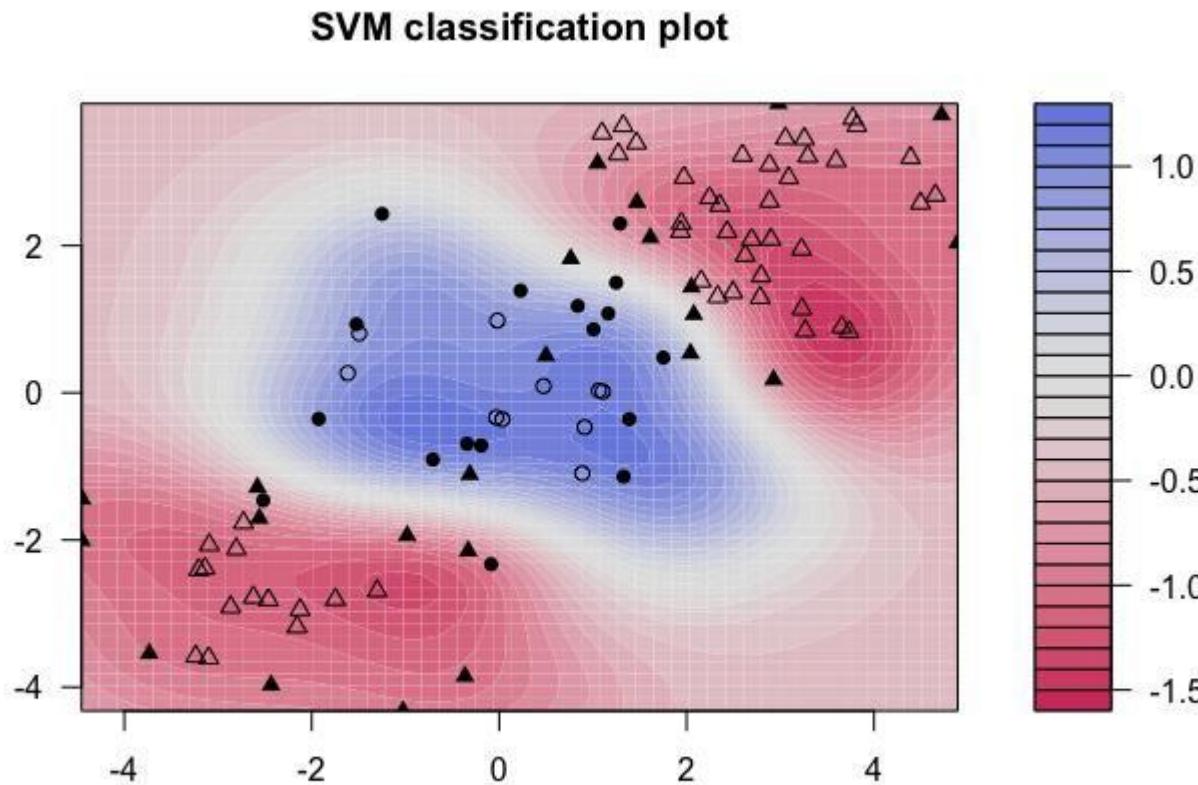
The most used Kernel is called “**Radial Basis Function**” (**RBF**)

$$K(x_i, x) \propto e^{-\frac{\|x_i - x\|^2}{\gamma^2}}$$

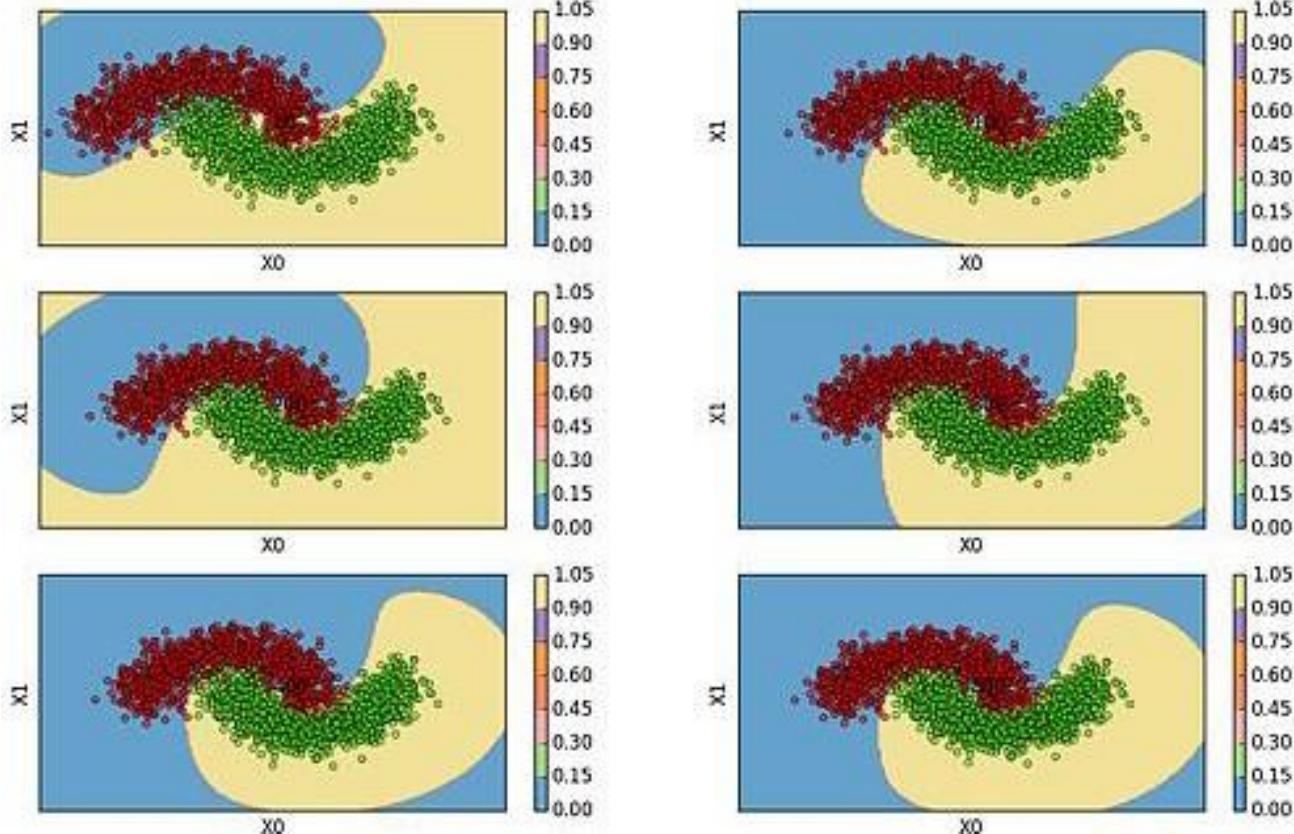


SC (Non linear)-SVMs

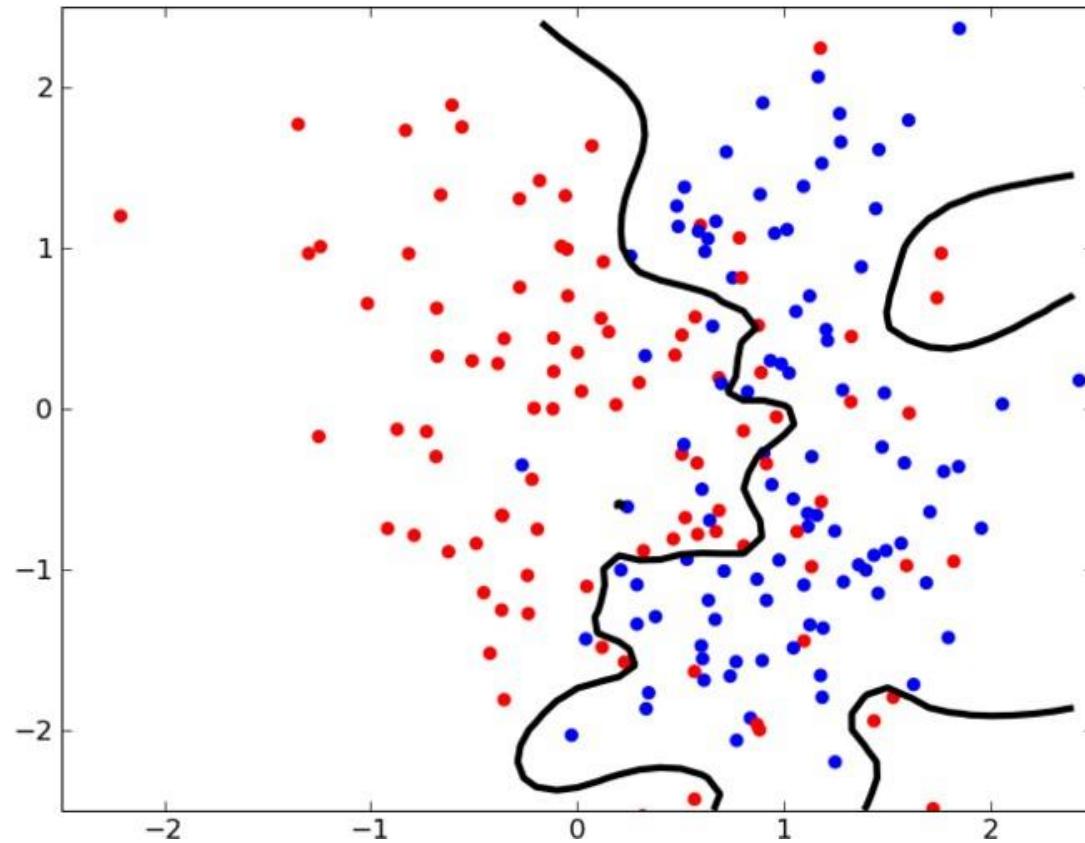
In general, the *border* can be very complicated, even *closed*.



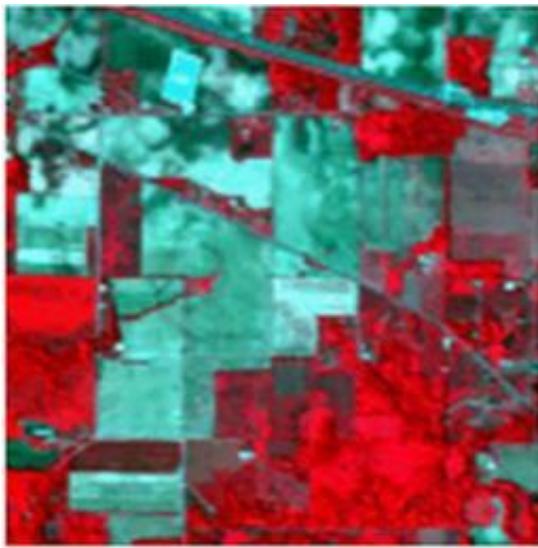
SC (Non linear)-SVMs



SC (Non linear)-SVMs

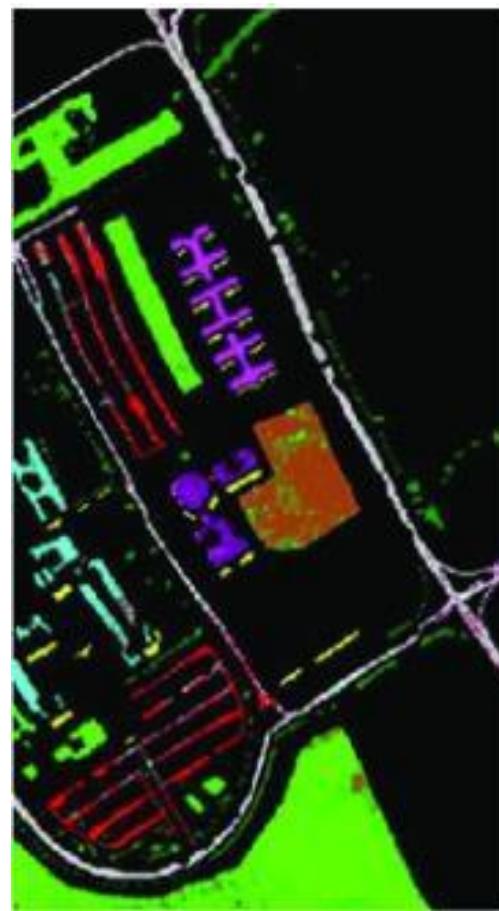
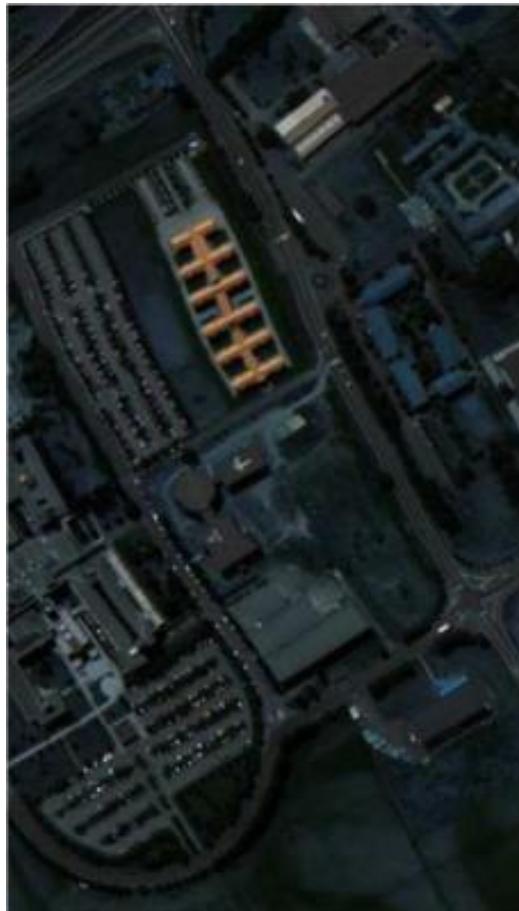


SC (Non linear)-SVMs

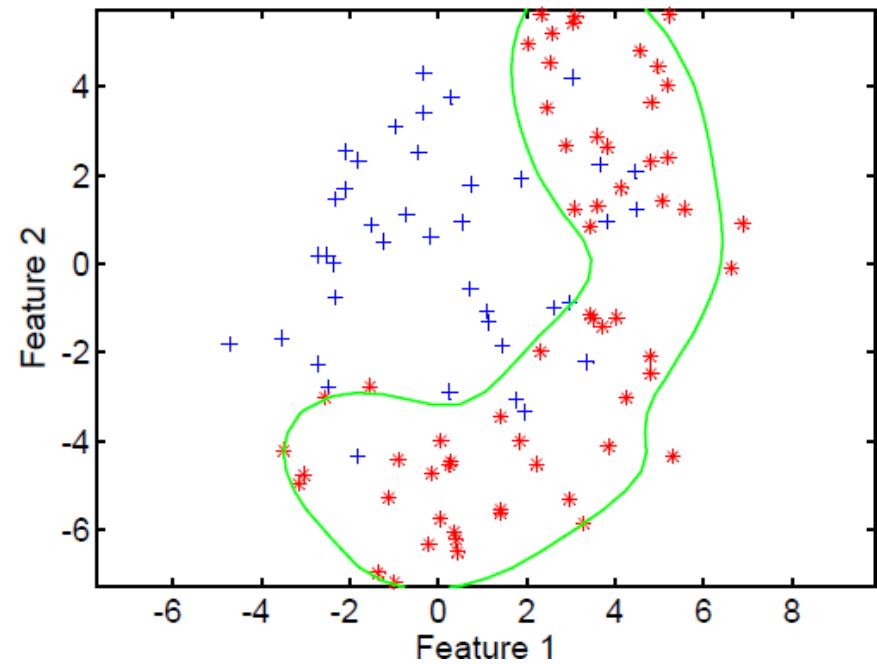
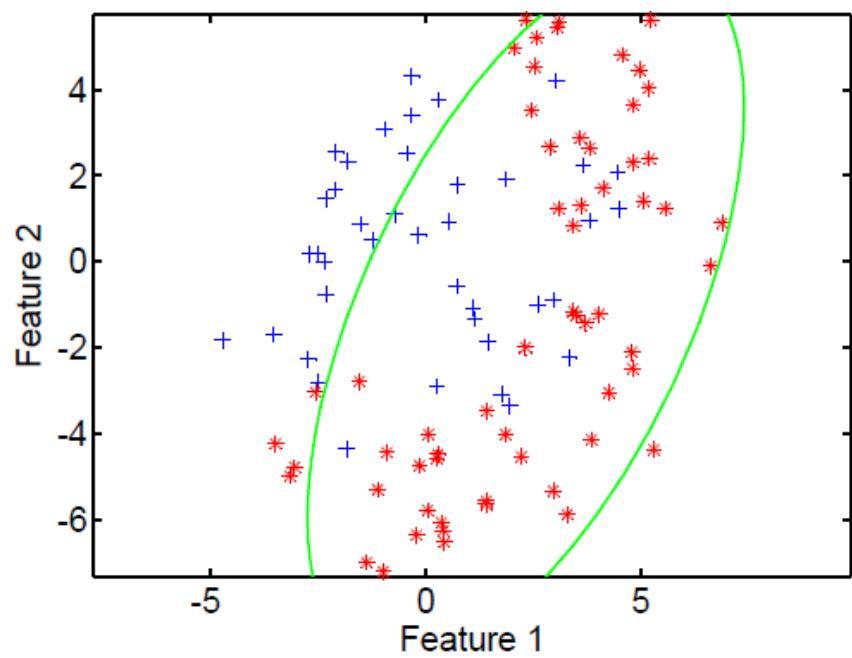


Alfalfa	Soybeans-notill
Corn-notill	Soybeans-min till
Corn-min till	Soybeans-clean till
Corn	Wheat
Hay-windowed	Woods
Grass/trees	Bldg
Grass/pasture-mowed	Stone-steel towers
Grass/pasture	Oats

SC (Non linear)-SVMs



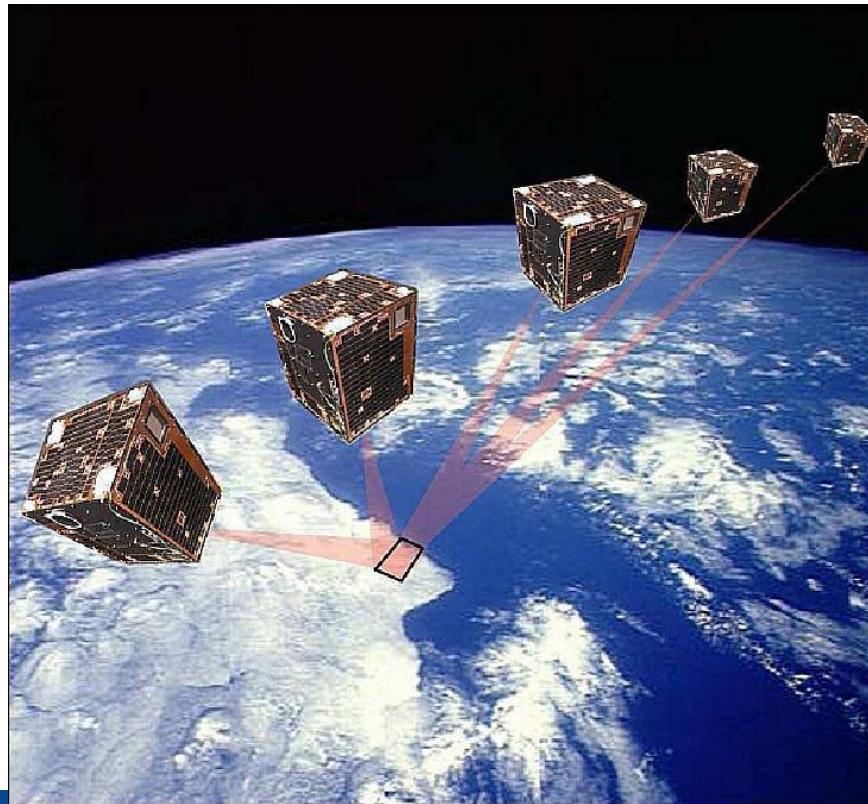
SC One class-SVM (OC-SVM)



SC One class-SVM (OC-SVM)

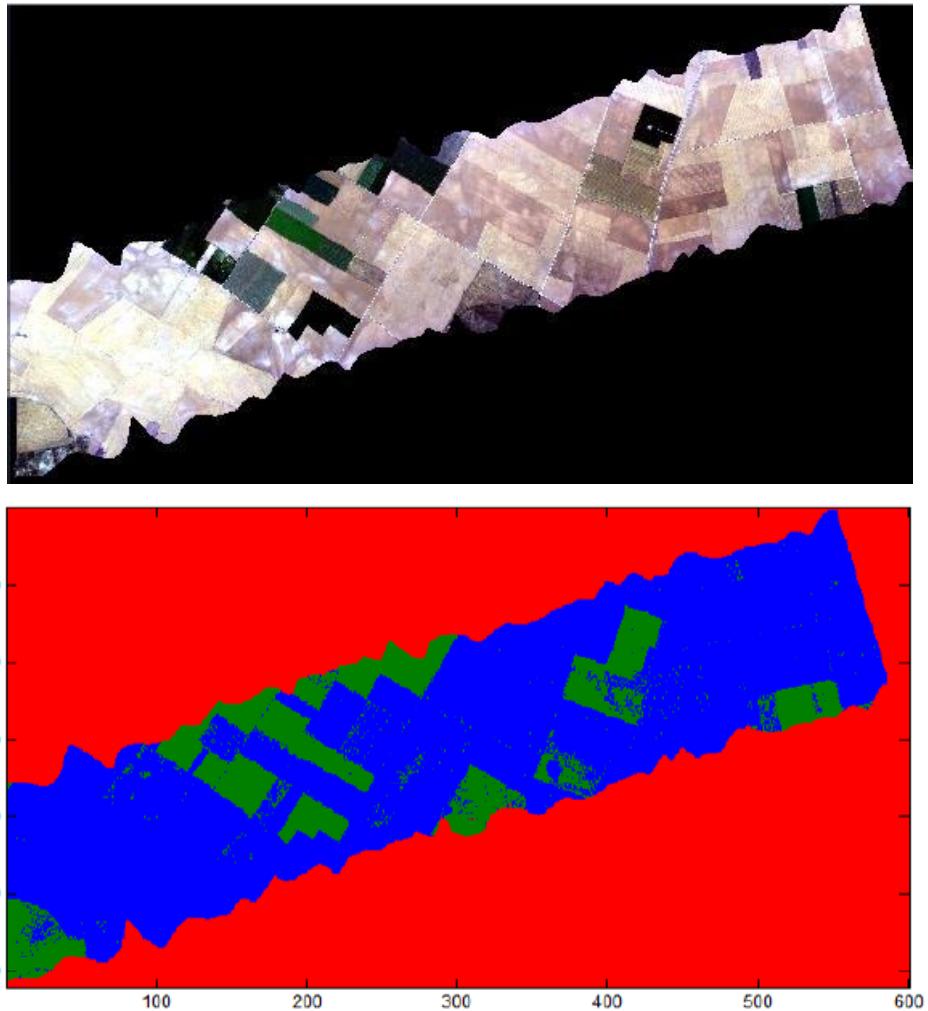
Compact High Resolution Imaging Spectrometer (CHRIS) on PROBA platform

Multispectral + Multiangular capabilities



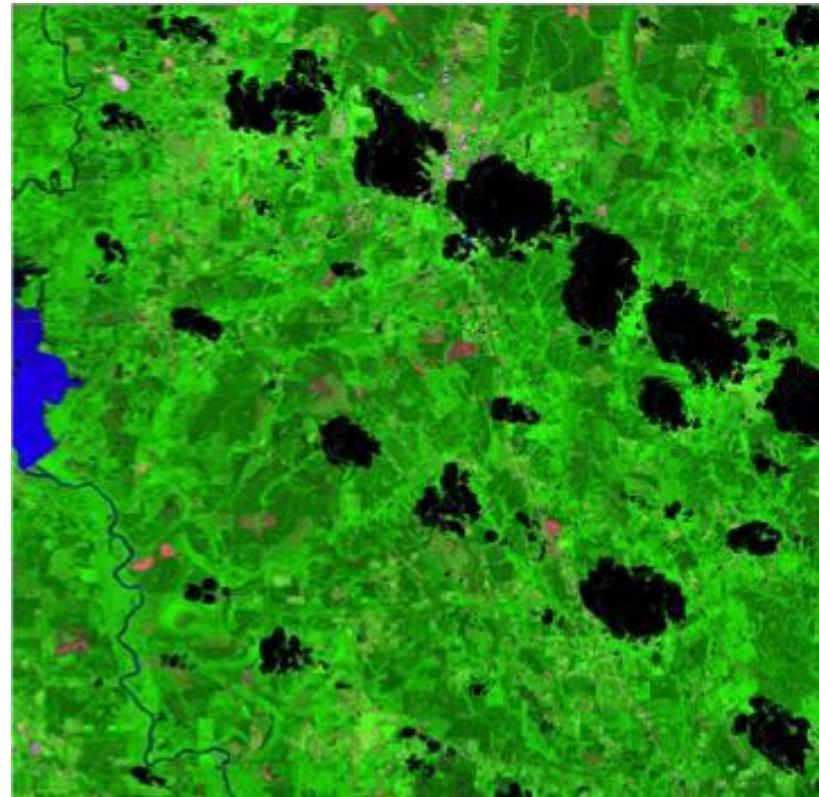
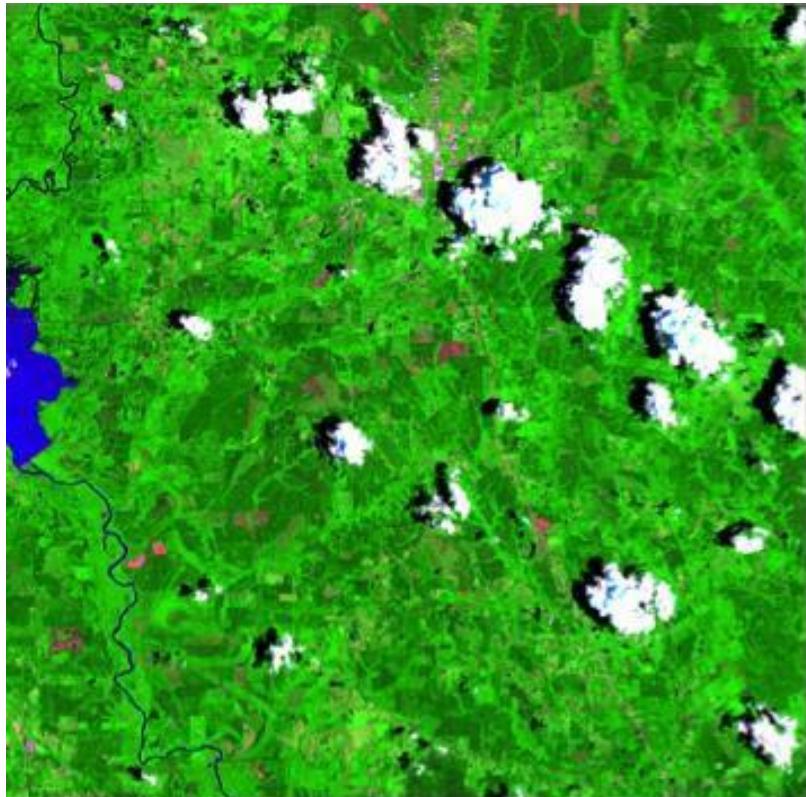
SC One class-SVM (OC-SVM)

Class of interest:
Vegetation



SC One class-SVM (OC-SVM)

A cloud mask



SC One class-SVM (OC-SVM)

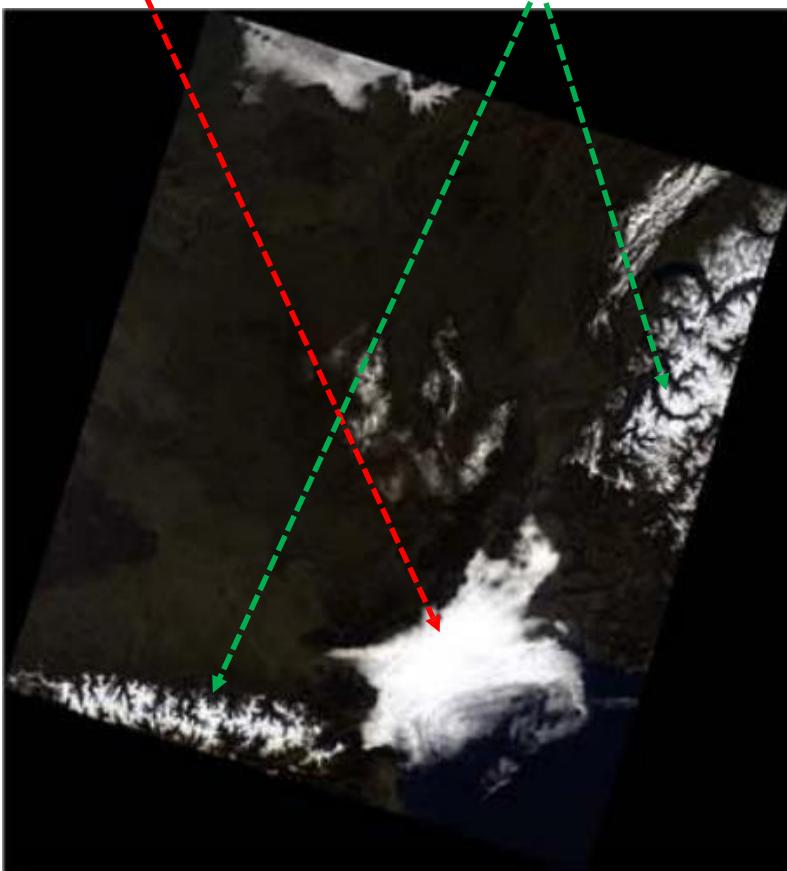
It is more complicated than that, but...



Shadows produced by clouds

SC One class-SVM (OC-SVM)

Clouds & No-clouds



Unsupervised classification (UC)

Clustering

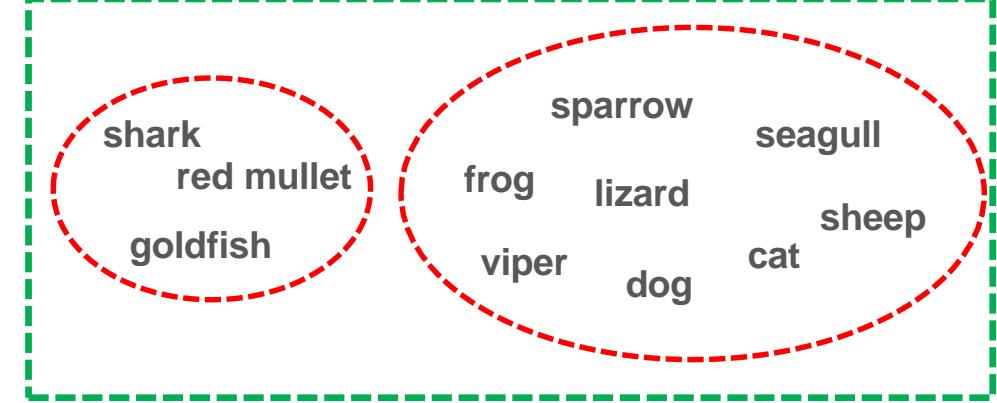
UC - Clustering

Imagine:

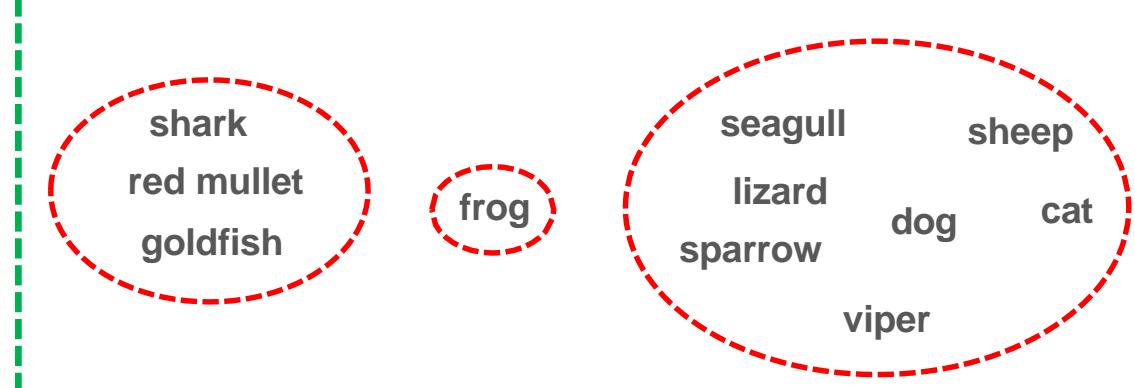
CRITERIA 1: PROGENY



CRITERIA 2: EXISTENCE OF LUNGS



CRITERIA 3: ENVIRONMENT WHERE THE ANIMALS LIVE



UC - Clustering

Clustering means “do clusters”.

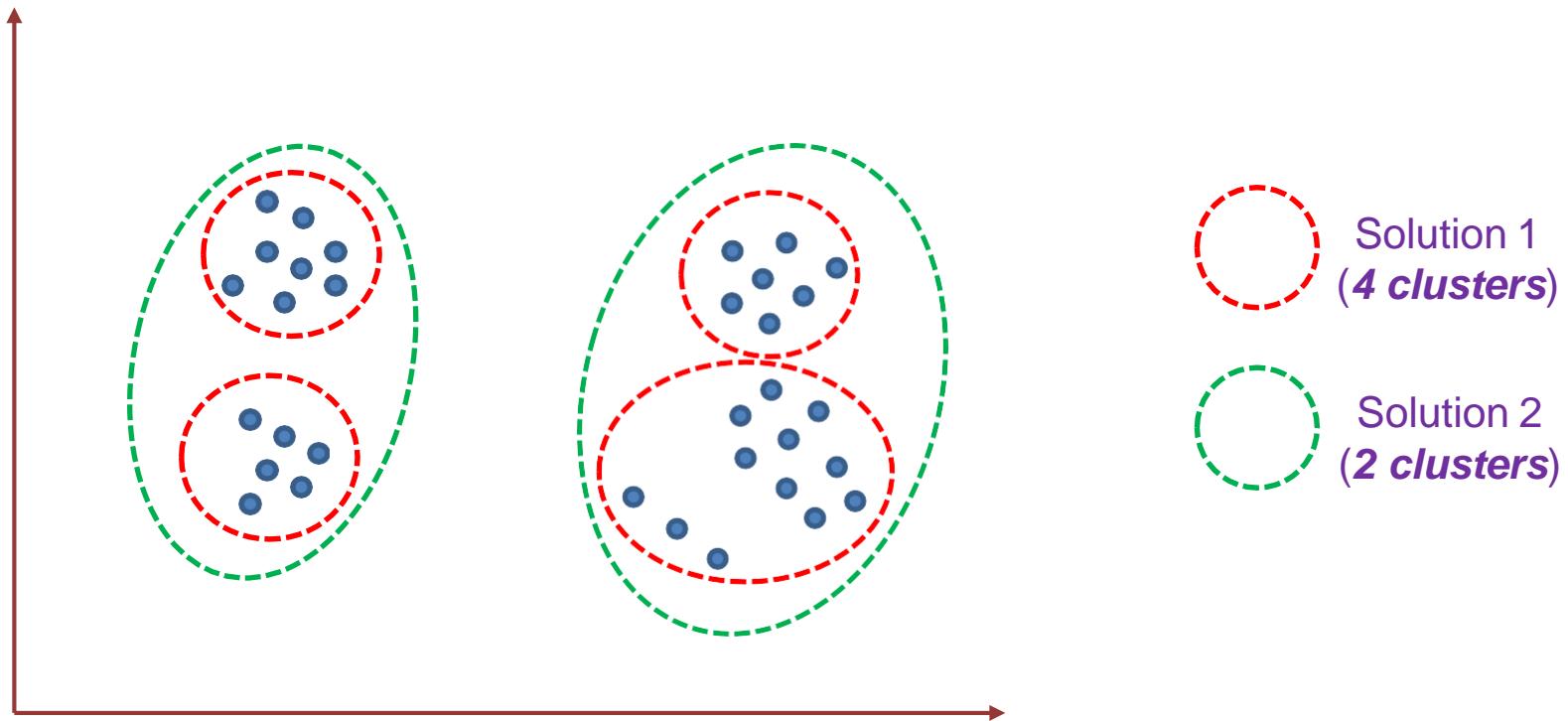
To do that:

- We need a **criterion**.
- Once the clusters are clustered we need a **validation** rule.
- Afterwards, we need to **interpret** the results.

UC - Clustering

Why do we need to validate the results?

Answer: Because sometimes several solutions are possible.



UC Clustering-Examples

K-Means

The idea is to minimise the distance from each data simple to a **representative** or **centroid**, i. e.:

Assume S sets: $S = \{S_1, S_2, \dots, S_k\}$. The aim is to find the **partition** that accomplishes:

$$\arg \min_S \left\{ \sum_{i=1}^k \|x - \mu_i\|^2 \right\}$$

UC Clustering-Examples

K-Means

This is accomplished **iteratively**, in two steps

$$S_i^{(t)} = \left\{ x_p : \|x_p - \mu_i^{(t)}\|^2 \leq \|x_p - \mu_j^{(t)}\|^2 ; \quad \forall j, 1 \leq j \leq k \right\}$$

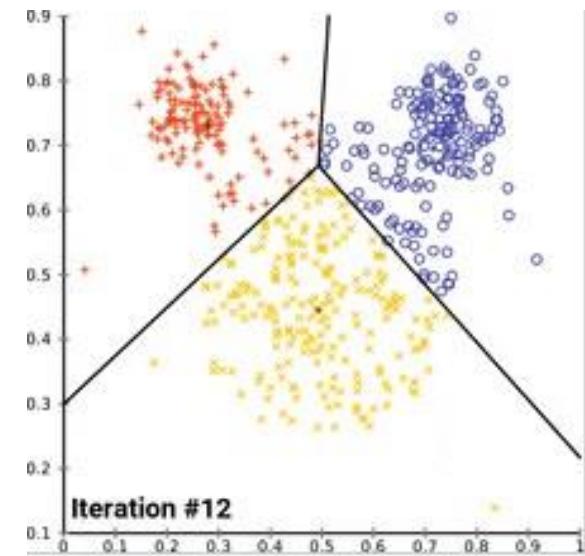
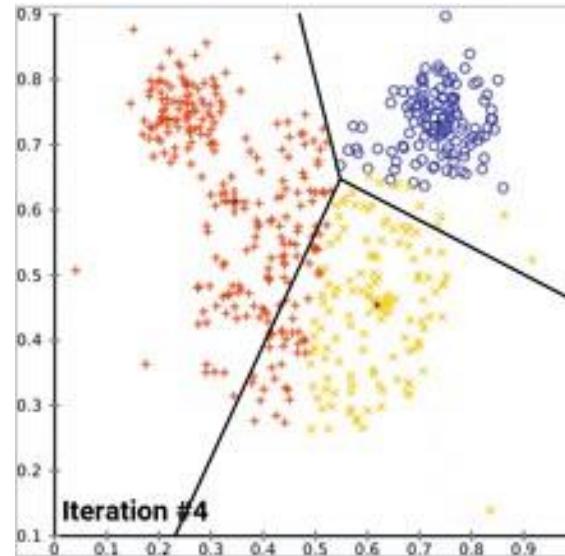
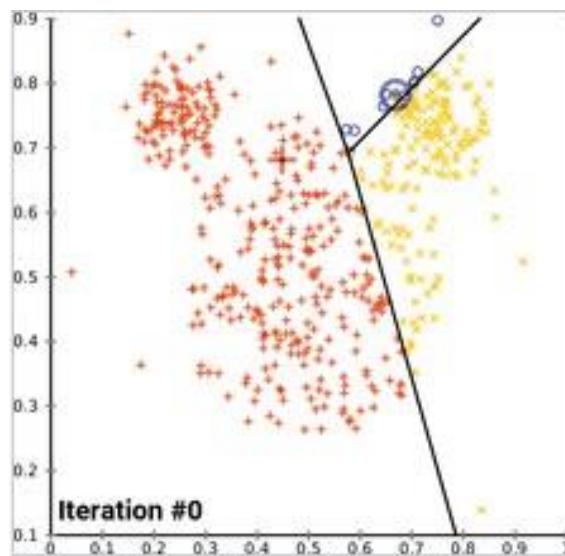
Assignment step:

$$\mu_j^{(t+1)} = \frac{1}{|S_i^{(t)}|} \sum_{x_j \in S_i^{(t)}} x_j$$

UC Clustering-Examples

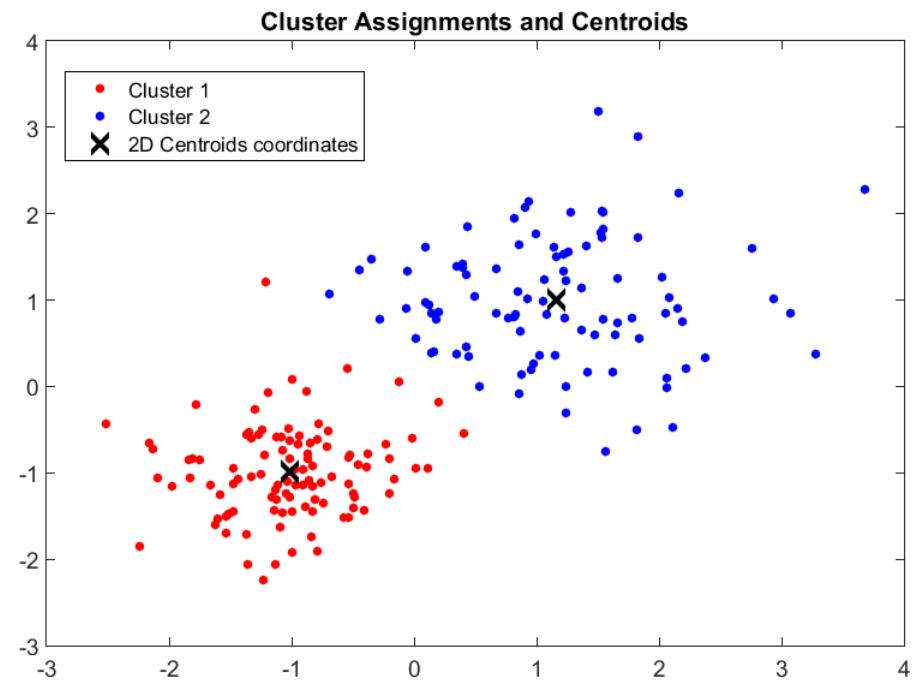
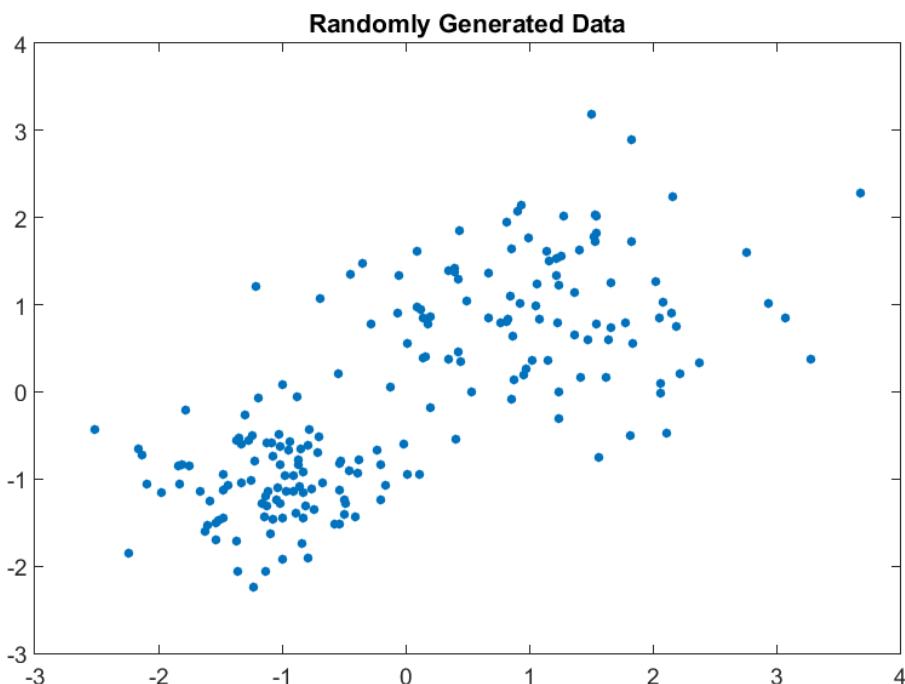
K-Means

Memberships and *centroids positions* are iteratively updated



UC Clustering-Examples

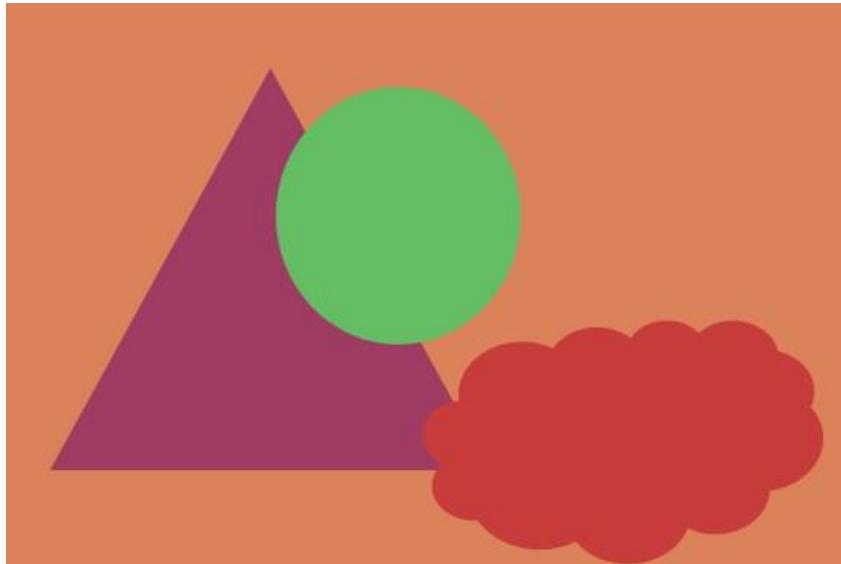
K-Means



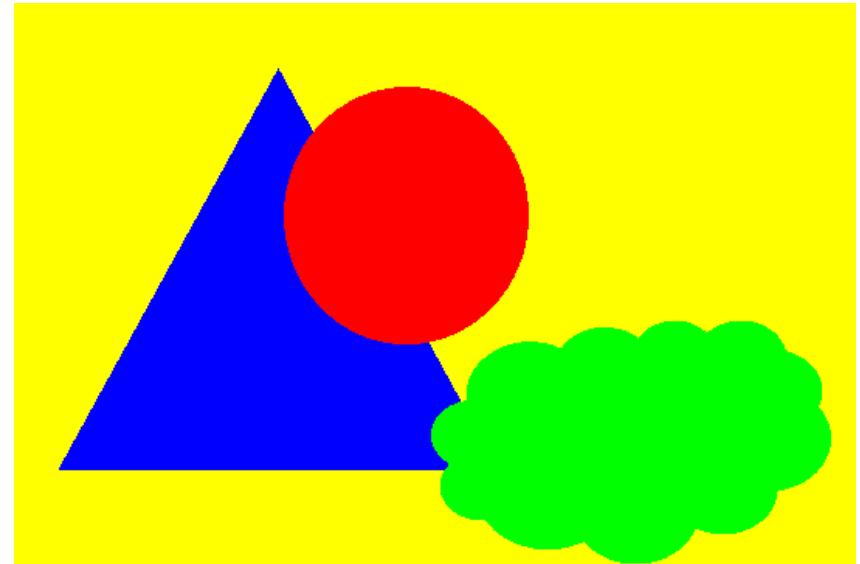
UC Clustering-Examples

K-Means – Noise free

Original – test image
4 clusters



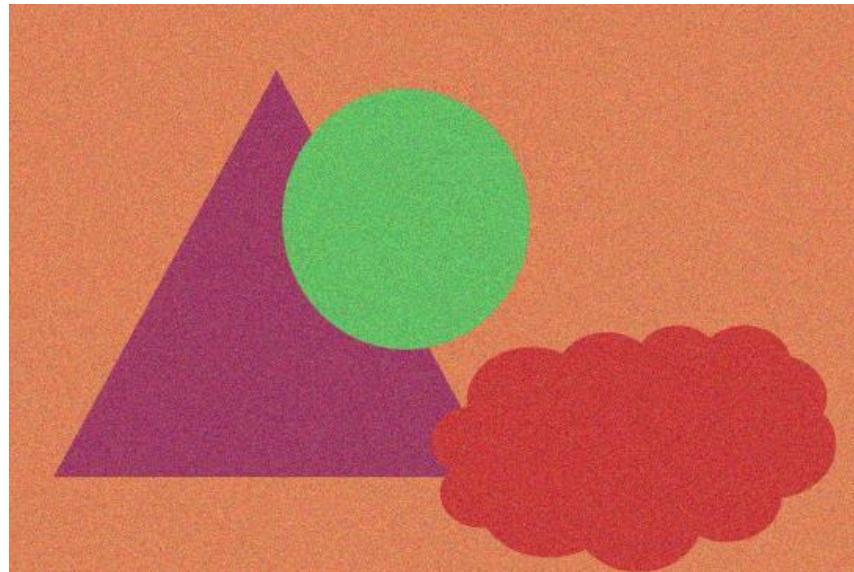
Segmentation / labeling result



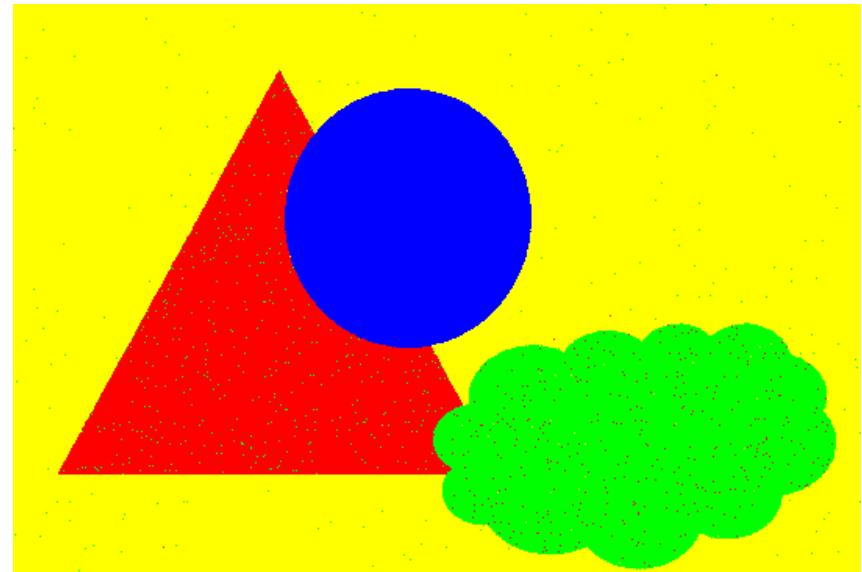
UC Clustering-Examples

K-Means – Gaussian noise

$\sigma = 0.05$



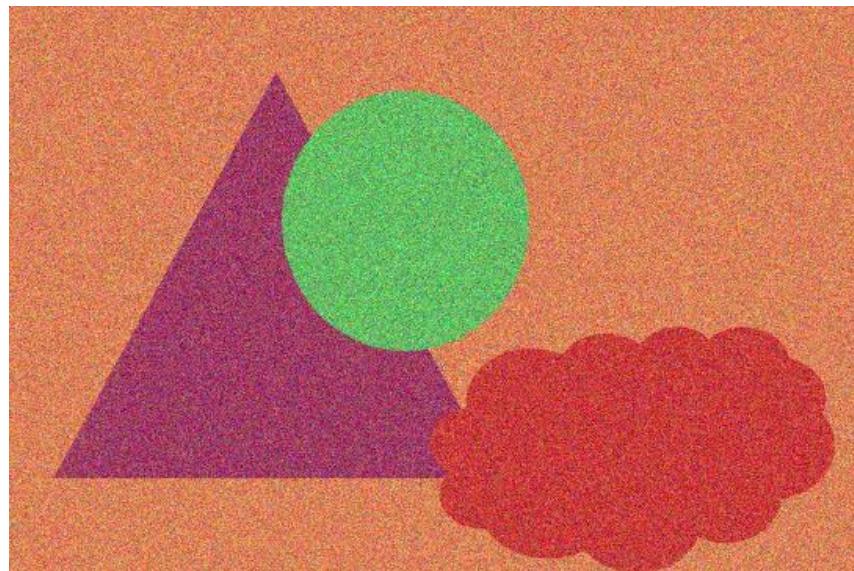
Segmentation / labeling result



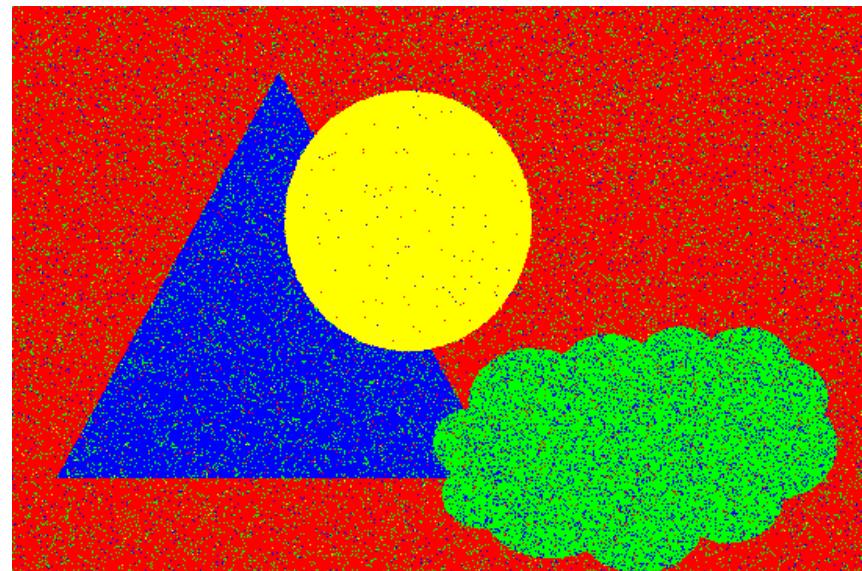
UC Clustering-Examples

K-Means – Gaussian noise

$\sigma = 0.1$



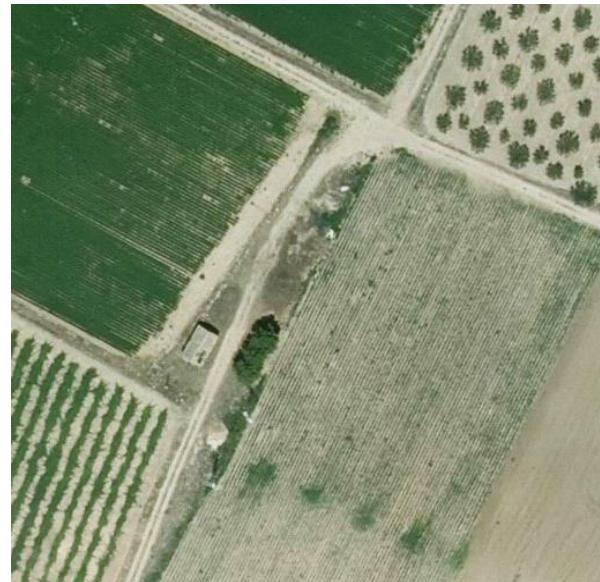
Segmentation / labeling result



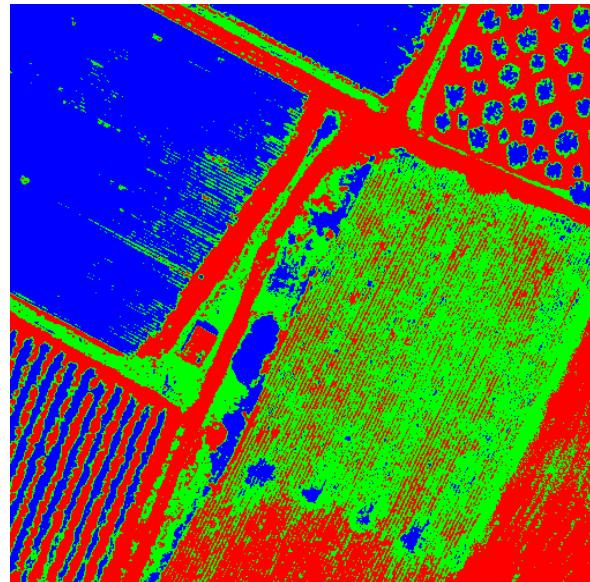
UC Clustering-Examples

K-Means – Examples from orthophotos

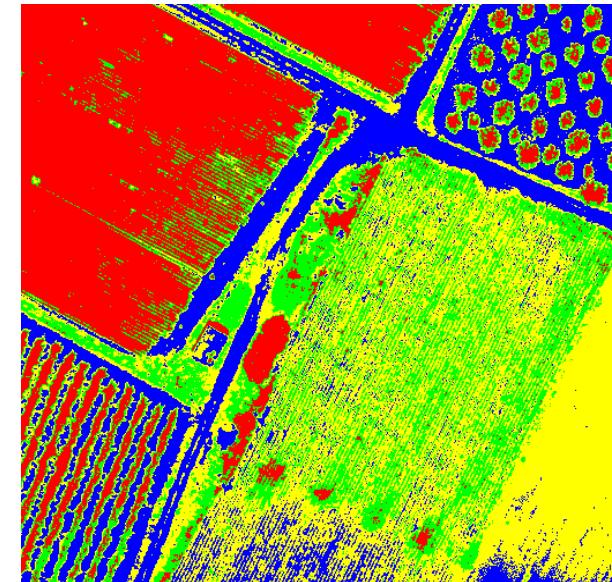
Original RGB image



3 clusters



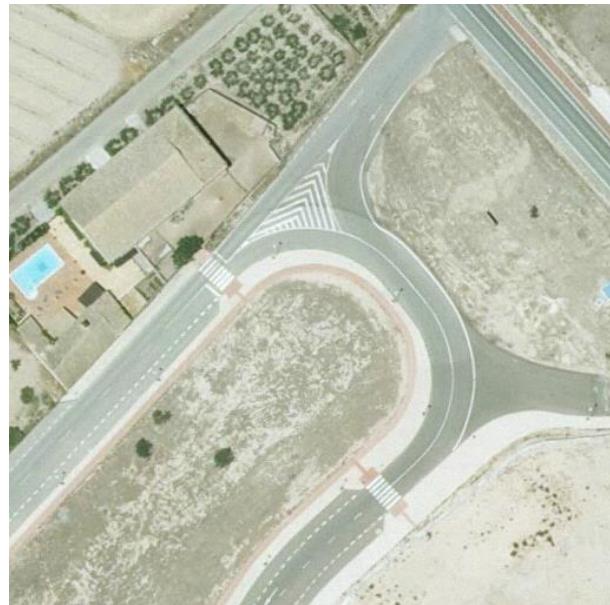
4 clusters



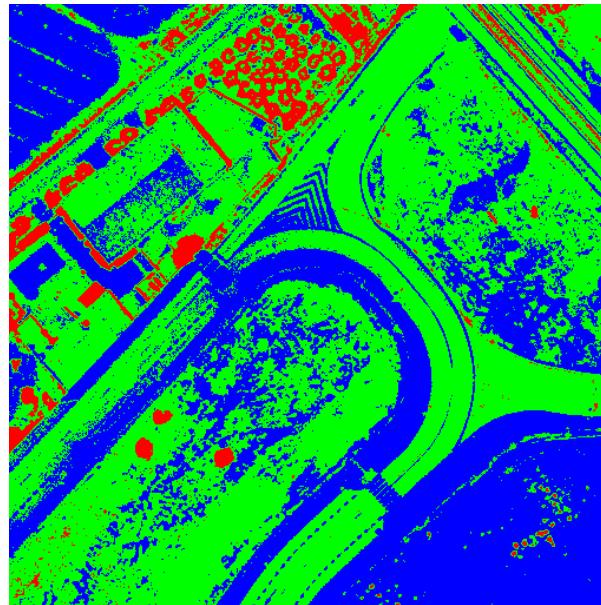
UC Clustering-Examples

K-Means – Examples from orthophotos

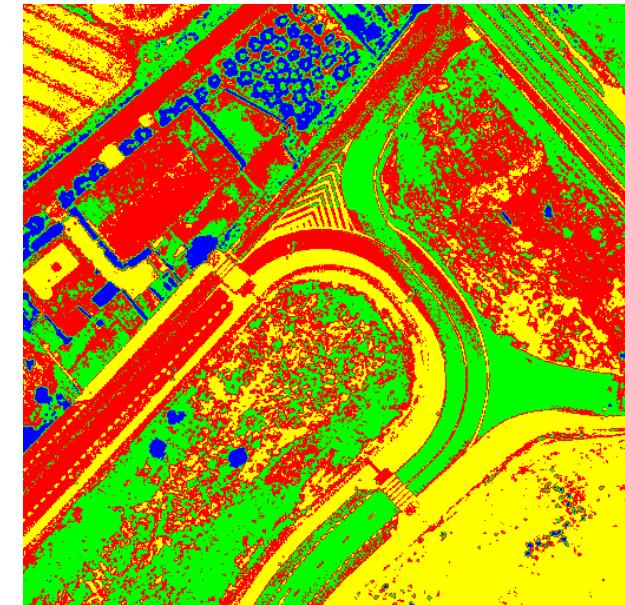
Original RGB image



3 clusters



4 clusters



UC Clustering-Examples

Fuzzy C-Means (FCM)

In K-Means each data point belongs to a cluster with *probability* equal to 1.

In FCM the method allows each data point i to have a certain *membership* value to each cluster, j (represented by $u_{ij} \equiv u_i(x_j)$).

The function to minimize would be:

$$J_m(U, v) = \sum_{j=1}^n \sum_{i=1}^c (u_{ij})^m \|x_j - \mu_i\|^2$$

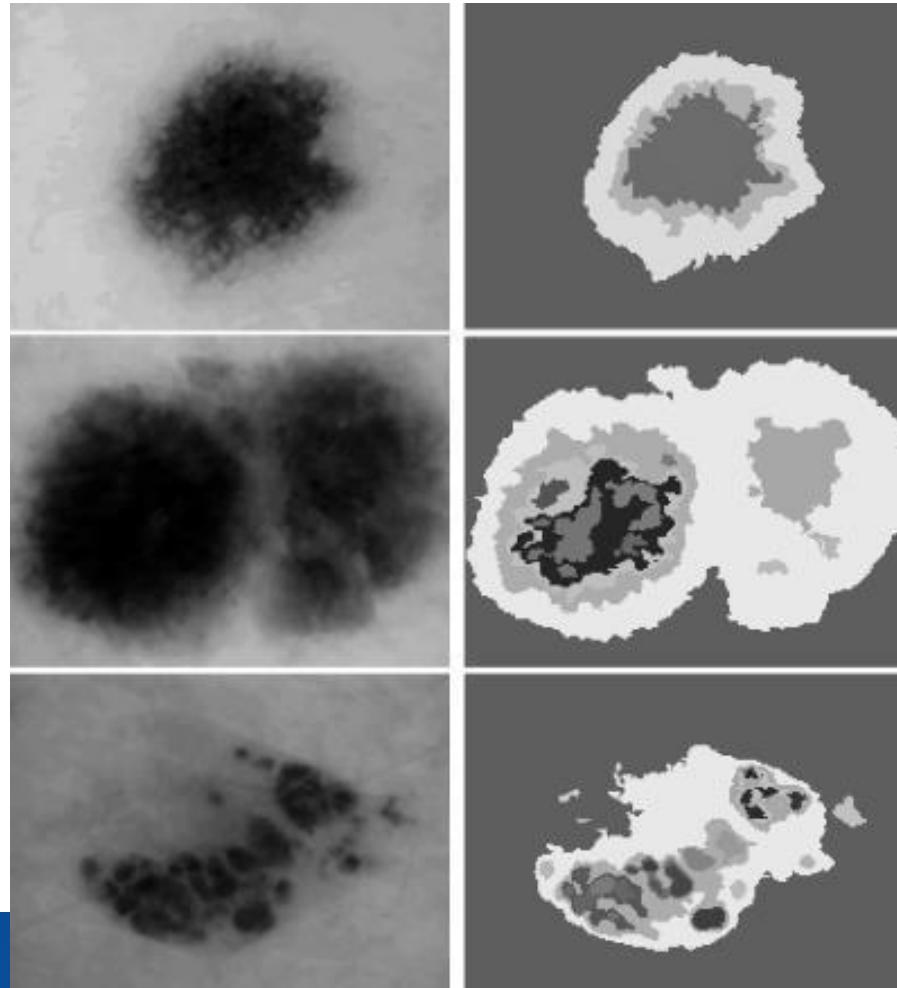
$$0 \leq u_{ij} \leq 1; \quad \sum_{i=1}^c u_{ij} = 1$$

UC Clustering-Examples

Fuzzy C-Means (FCM)

Where could FCM work well?

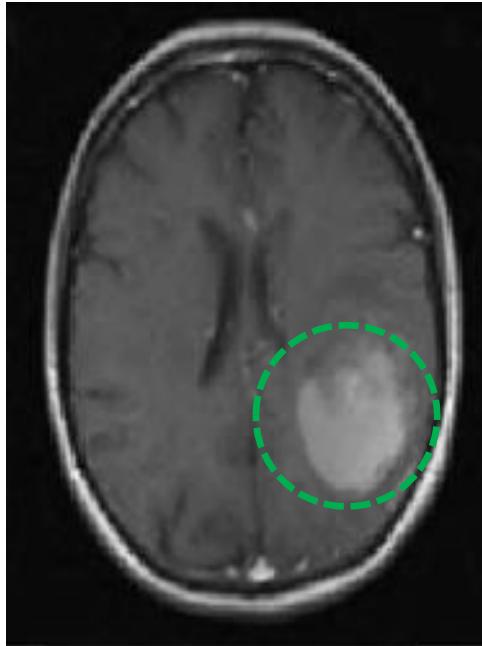
Answer: For instance, in objects
with ***diffused borders***



UC Clustering-Examples

Fuzzy C-Means (FCM)

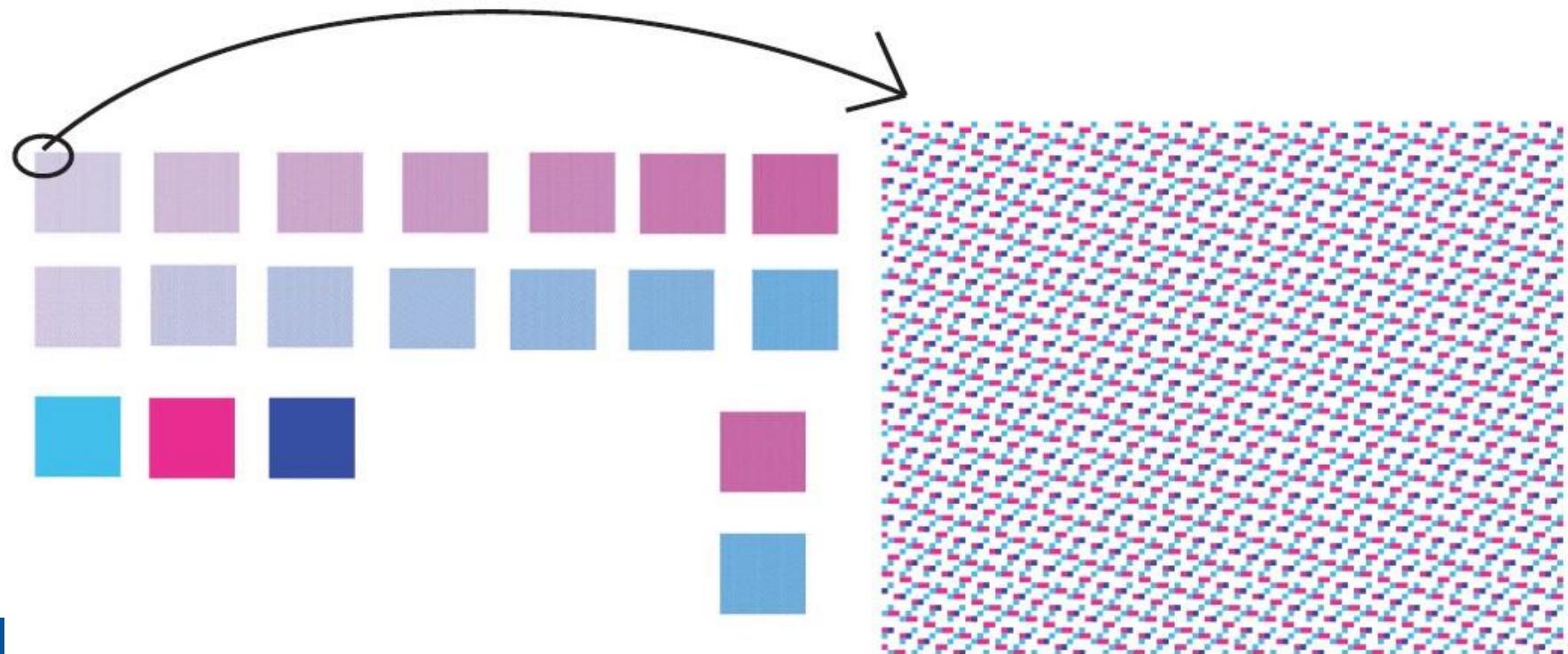
BRAIN



UC Clustering-Examples

Fuzzy C-Means (*FCM*)

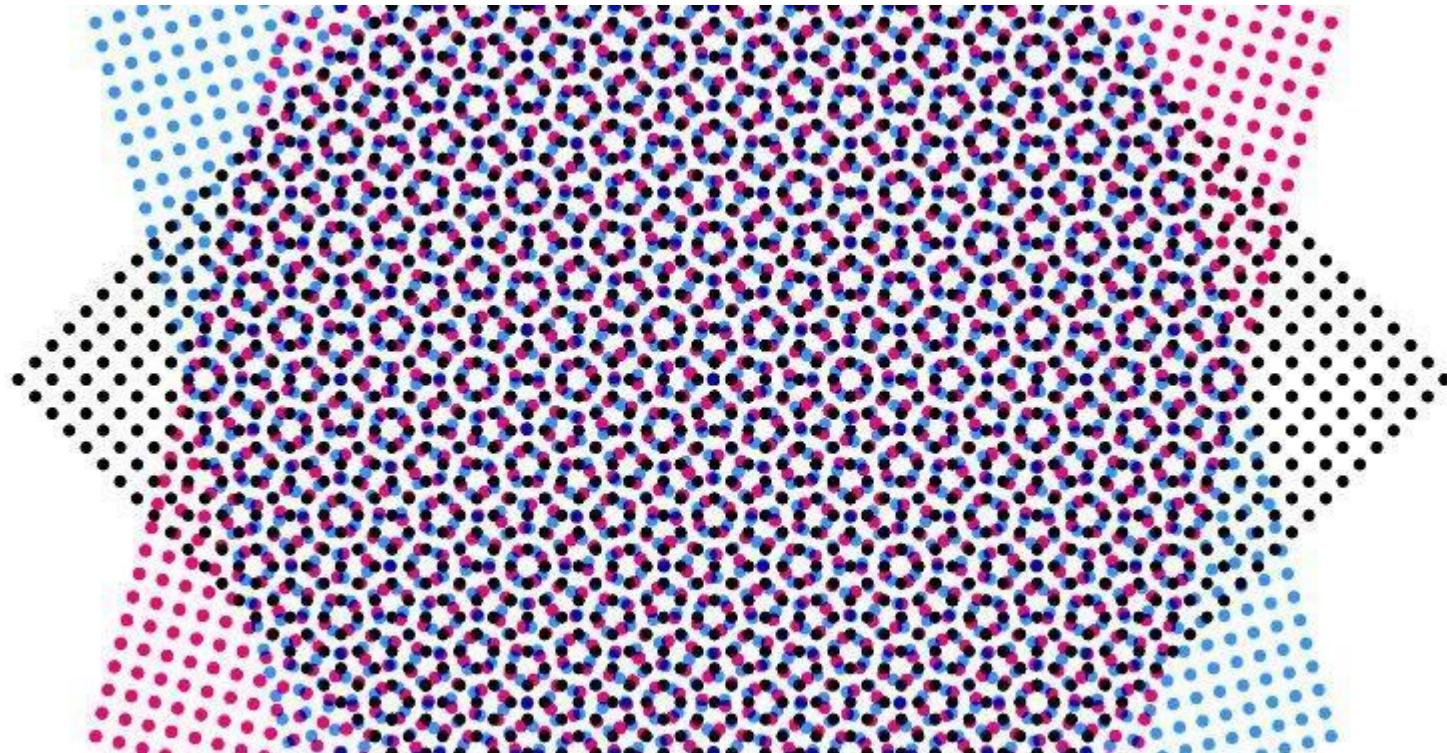
Digital file of a **halftone** (i. e., print/no print) CMYK design



UC Clustering-Examples

Fuzzy C-Means (FCM)

Be careful with the/a **MOIRÉ** pattern!



UC Clustering-Examples

Fuzzy C-Means (FCM)

A binary (CMYK) ink jet printer was used.



UC Clustering-Examples

Fuzzy C-Means (FCM)



UC Clustering-Examples

Fuzzy C-Means (FCM)

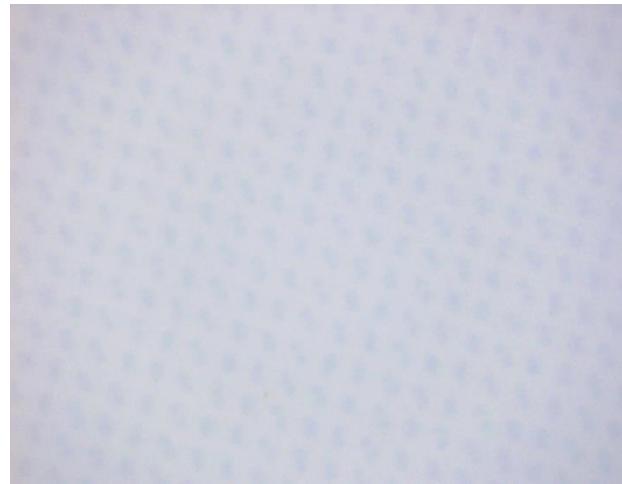
Images of the printed ceramic tile at microscopic scale were acquired:



UC Clustering-Examples

Fuzzy C-Means (FCM)

$C = 20\%$



$C = 50\%$



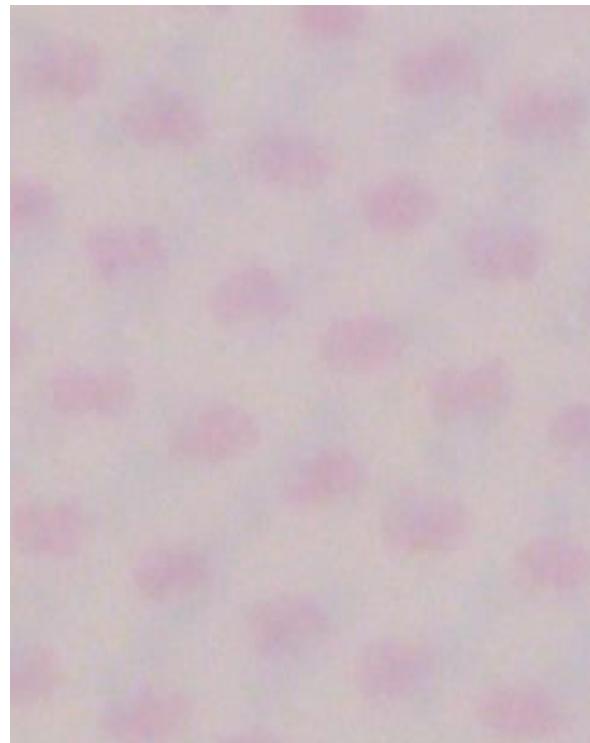
$C = 80\%$



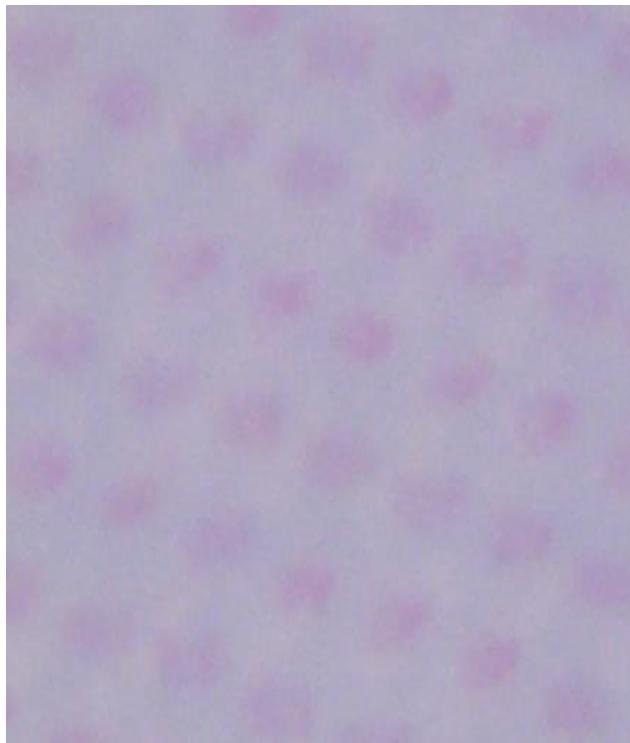
UC Clustering-Examples

Fuzzy C-Means (FCM)

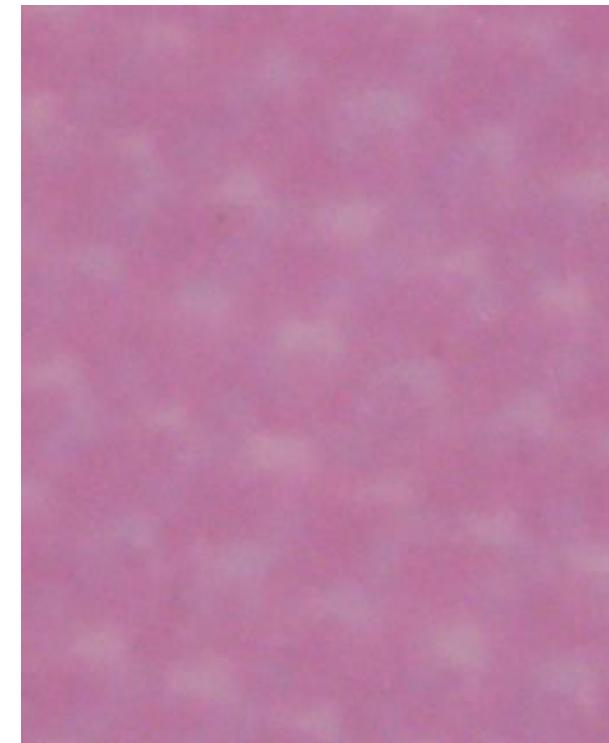
$C = 20\%$; $M = 20\%$



$C = 80\%$; $M = 20\%$



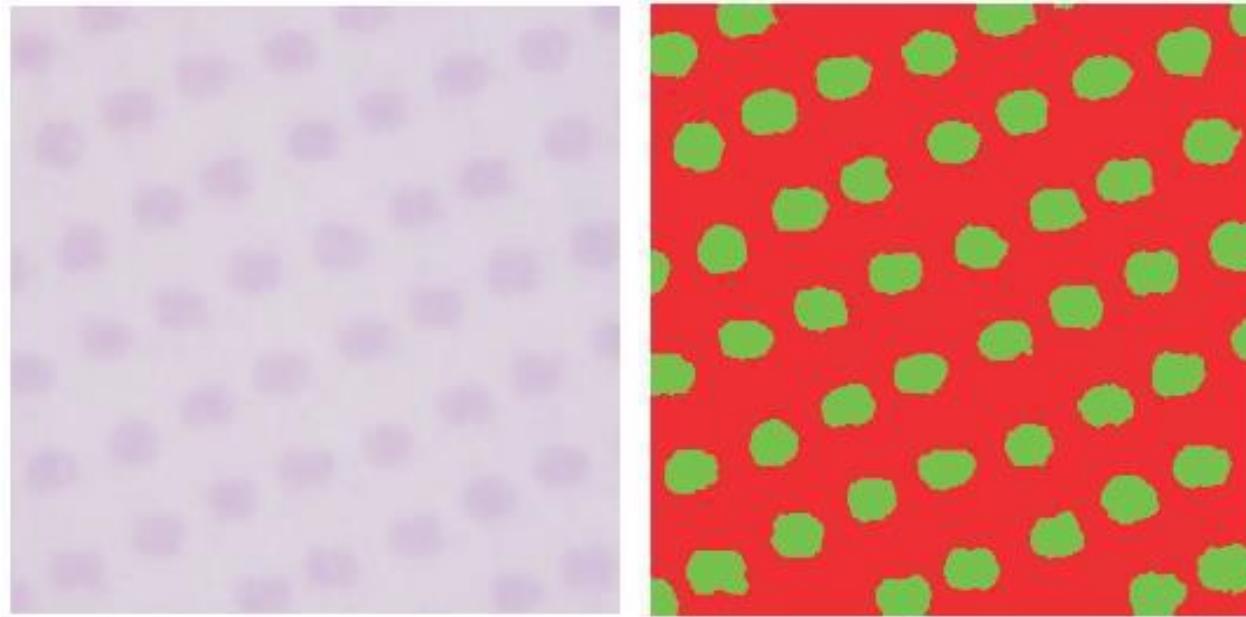
$C = 20\%$; $M = 80\%$



UC Clustering-Examples

Fuzzy C-Means (FCM)

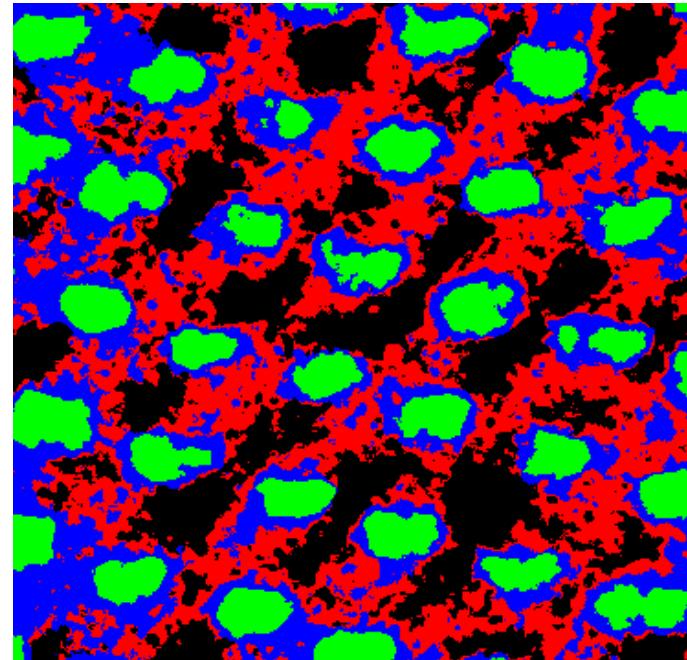
Microscopic image of *Magenta* (color) printed dots on a ceramic tile



UC Clustering-Examples

Fuzzy C-Means (FCM)

$C = 70\%$; $M = 20\%$



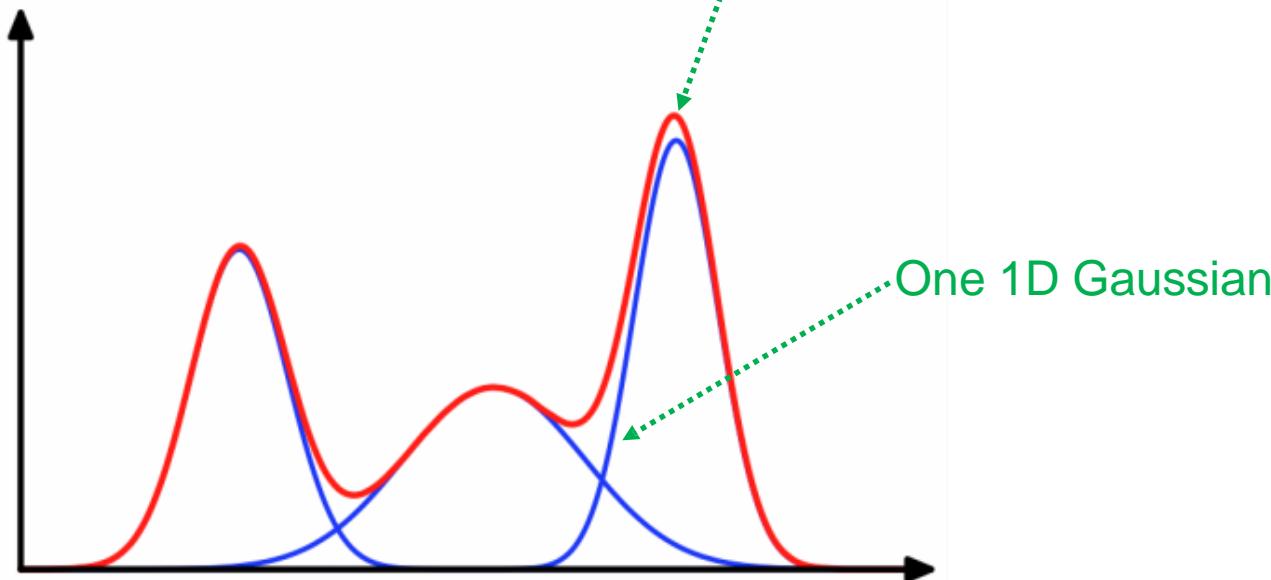
UC Clustering-Examples

Mixtures of Gaussians – Expectation Maximization (EM)

What is a **Mixture of Gaussians (MoGs)**?

$$p(x) = \sum_{j=1}^M p(\Omega_j) \cdot p(x|\Omega_j)$$

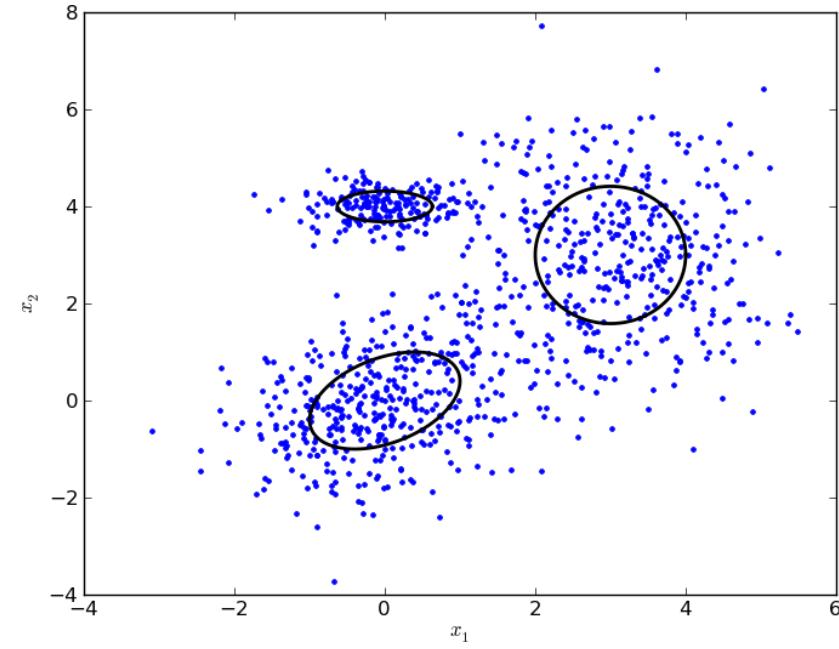
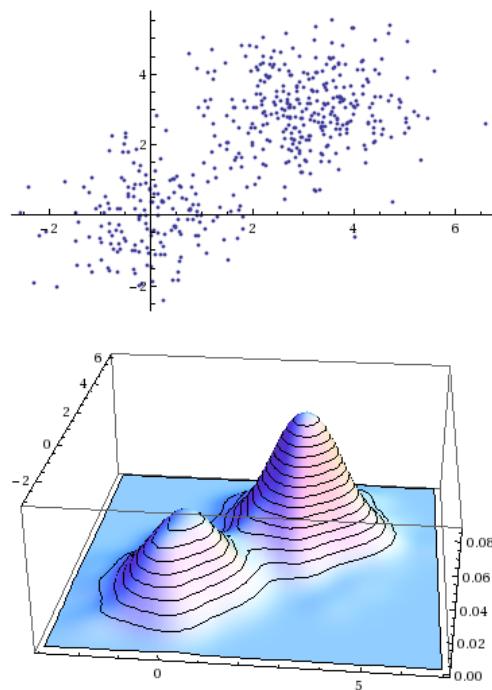
Evolving curve –
Linear combination of Gaussians



UC Clustering-Examples

Mixtures of Gaussians – Expectation Maximization (EM)

What would we have in 2D?



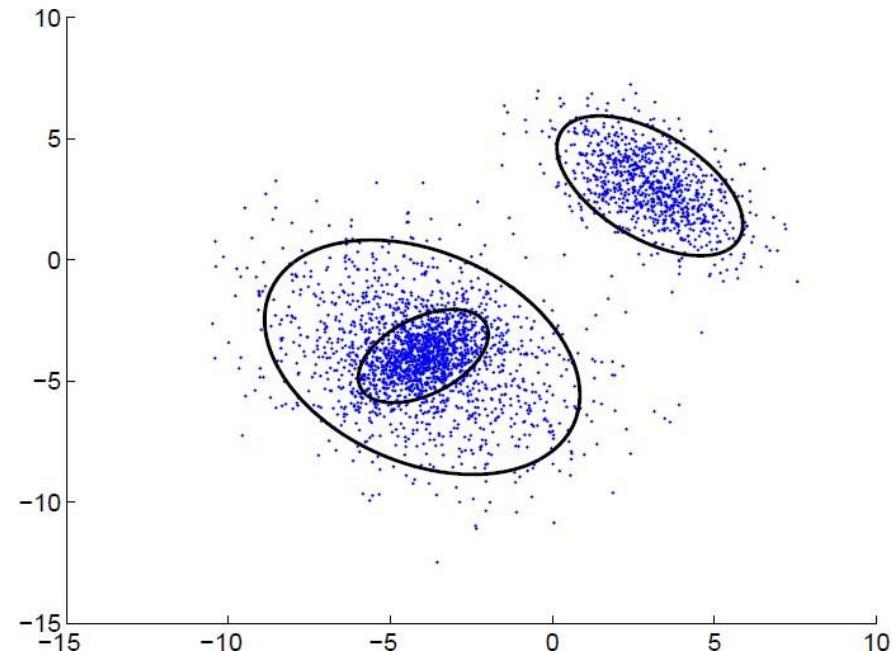
UC Clustering-Examples

Mixtures of Gaussians – Expectation Maximization (EM)

In a mixture of gaussians, we have:

$$p(x) = \sum_{j=1}^M p(\Omega_j) \cdot p(x|\Omega_j)$$

$$p(x|\Omega_j) = \frac{1}{(2\pi)^{d/2} |\Sigma_j|^{1/2}} e^{-\frac{1}{2} \left((x-\mu_j)^T \Sigma_j^{-1} (x-\mu_j) \right)}$$



UC Clustering-Examples

Mixtures of Gaussians – Expectation Maximization (EM)

From Bayes rule and

$$p(x) = \sum_{j=1}^M p(\Omega_j) \cdot p(x|\Omega_j); \quad p(x|\Omega_j) = \frac{1}{(2\pi)^{d/2} |\Sigma_j|^{1/2}} e^{\left[-\frac{1}{2} \cdot ((x - \mu_j)^T \Sigma_j^{-1} (x - \mu_j)) \right]}$$

we have:

$$p(\Omega_j|x) = \frac{p(\Omega_j) \cdot p(x|\Omega_j)}{\sum_{k=1}^M p(\Omega_k) \cdot p(x|\Omega_k)}$$

UC Clustering-Examples

Mixtures of Gaussians – Expectation Maximization (EM)

Expectation (E) step:

$$p(\Omega_j | x_i)^{(t)} = \frac{p(\Omega_j)^{(t)} \cdot p(x_i | \Omega_j)^{(t)}}{\sum_k p(\Omega_k)^{(t)} \cdot p(x_i | \Omega_k)^{(t)}}$$

Maximization (M) step:

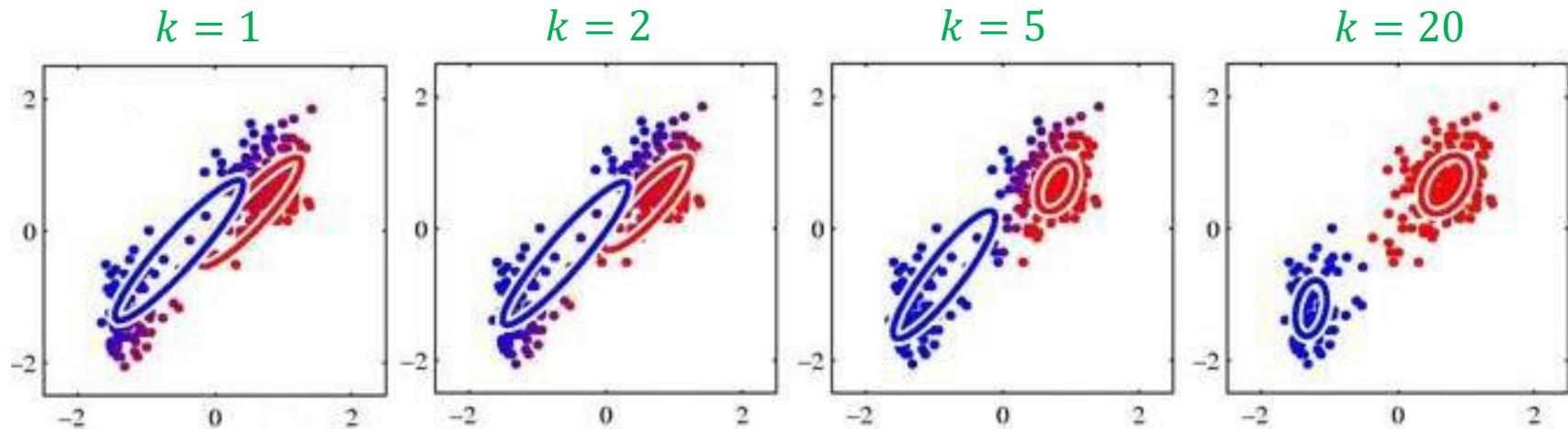
$$p(\Omega_j)^{(t+1)} = \frac{1}{n} \sum_{i=1}^n \{ p(\Omega_j | x_i)^{(t)} \}$$

$$\mu_j^{(t+1)} = \frac{\sum_{i=1}^n \{ x_i p(\Omega_j | x_i)^{(t)} \}}{\sum_{i=1}^n \{ p(\Omega_j | x_i)^{(t)} \}}$$

$$\Sigma_j^{(t+1)} = \frac{\sum_{i=1}^n \{ p(\Omega_j | x_i)^{(t)} \} (x_i - \mu_j^{(t+1)}) (x_i - \mu_j^{(t+1)})^T}{\sum_{i=1}^n \{ p(\Omega_j | x_i)^{(t)} \}}$$

UC Clustering-Examples

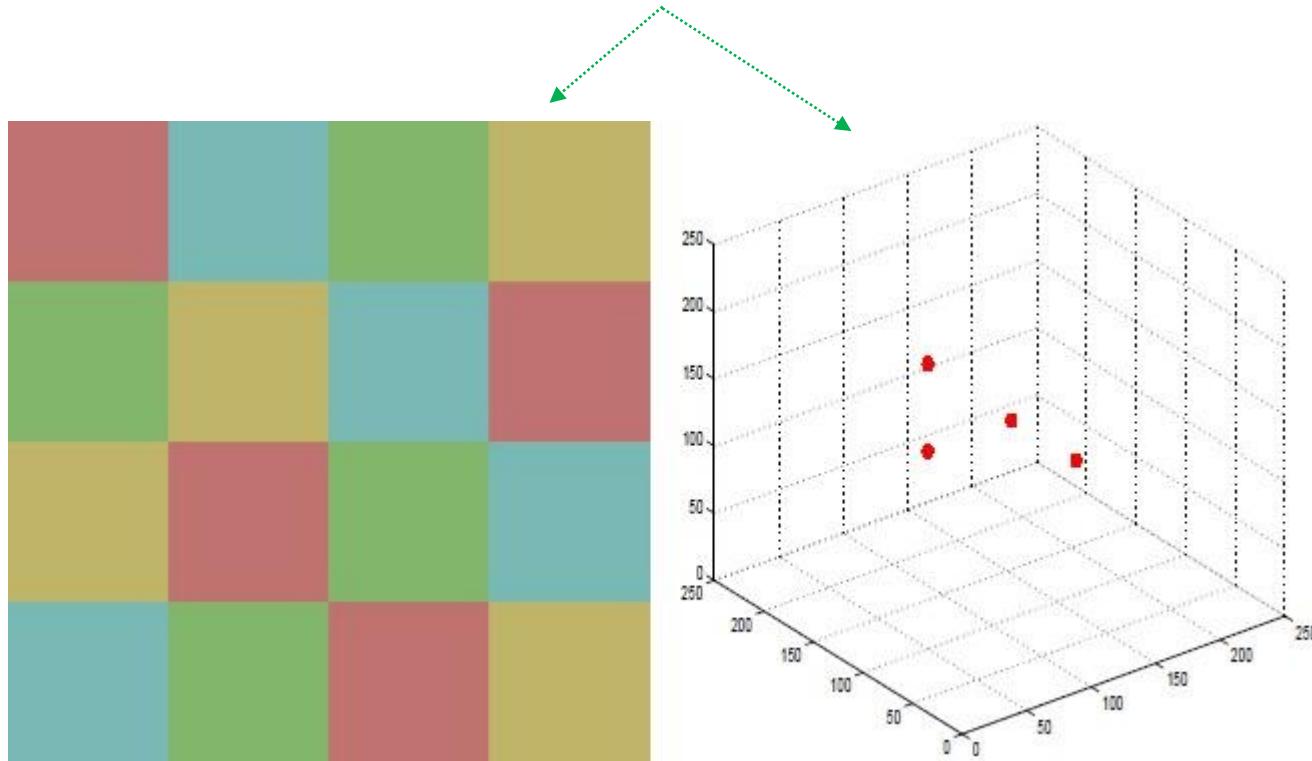
Mixtures of Gaussians – Expectation Maximization (EM)



UC Clustering-Examples

Mixtures of Gaussians – Expectation Maximization (EM)

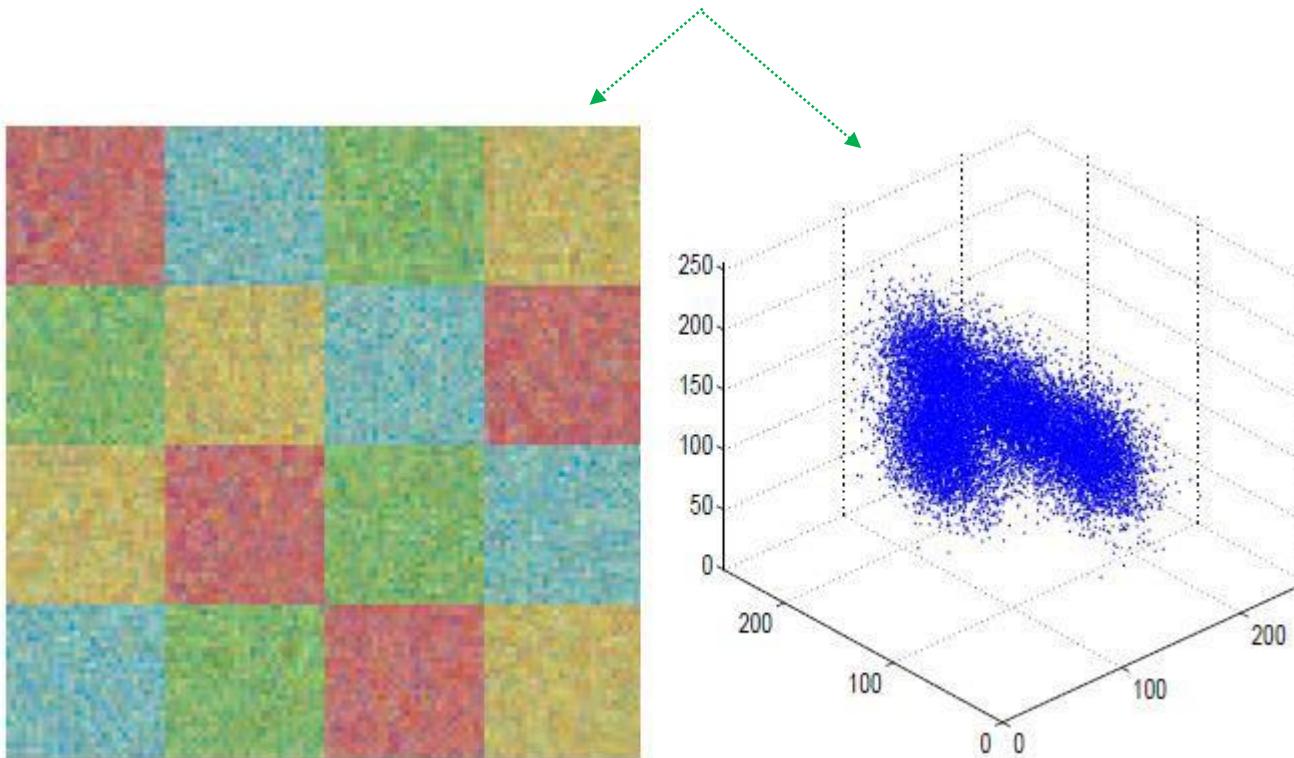
We have a “clean image”



UC Clustering-Examples

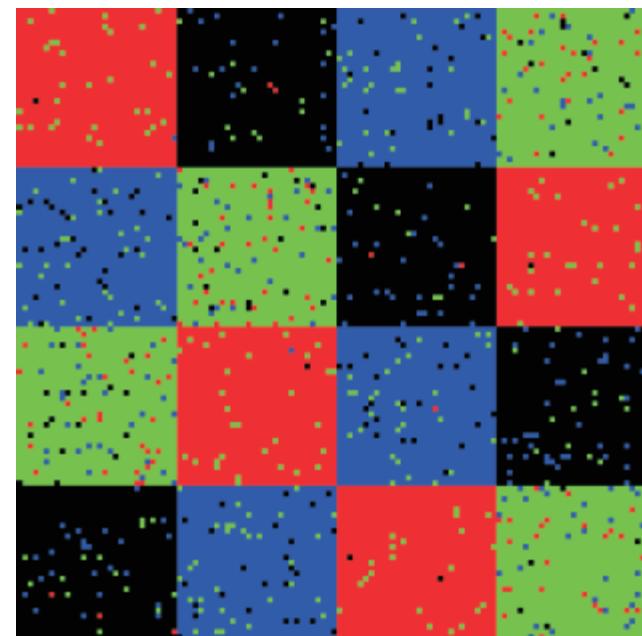
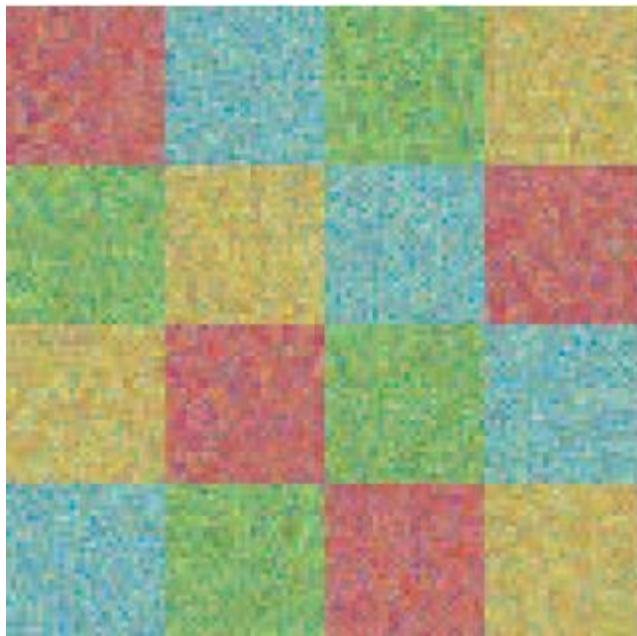
Mixtures of Gaussians – Expectation Maximization (EM)

We introduce “gaussian noise”



UC Clustering-Examples

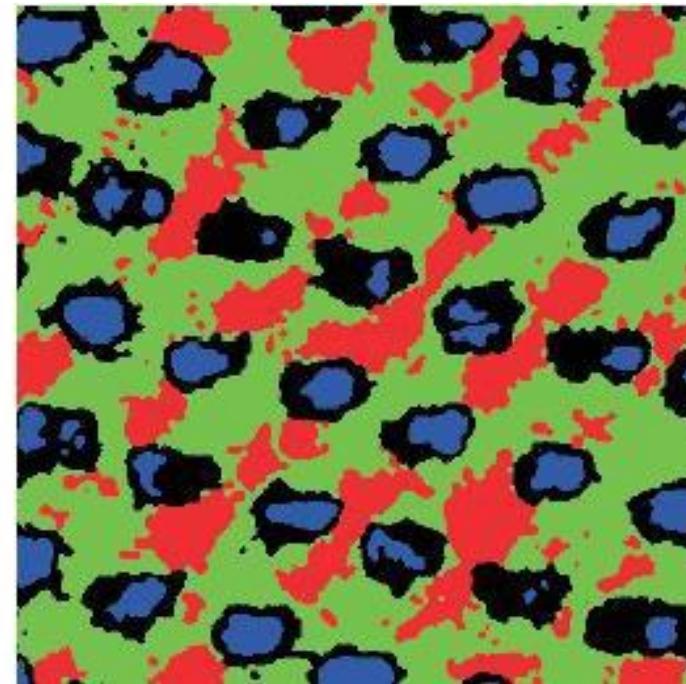
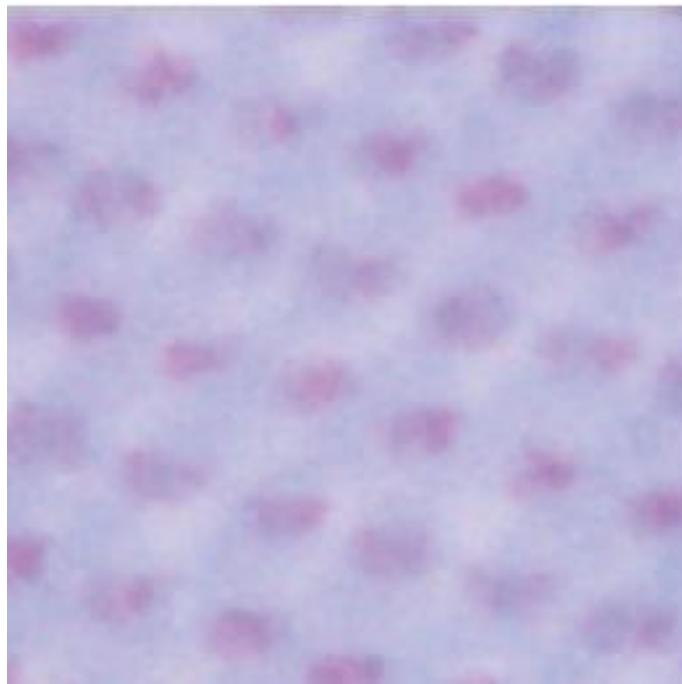
Mixtures of Gaussians – Expectation Maximization (EM)



UC Clustering-Examples

Mixtures of Gaussians – Expectation Maximization (EM)

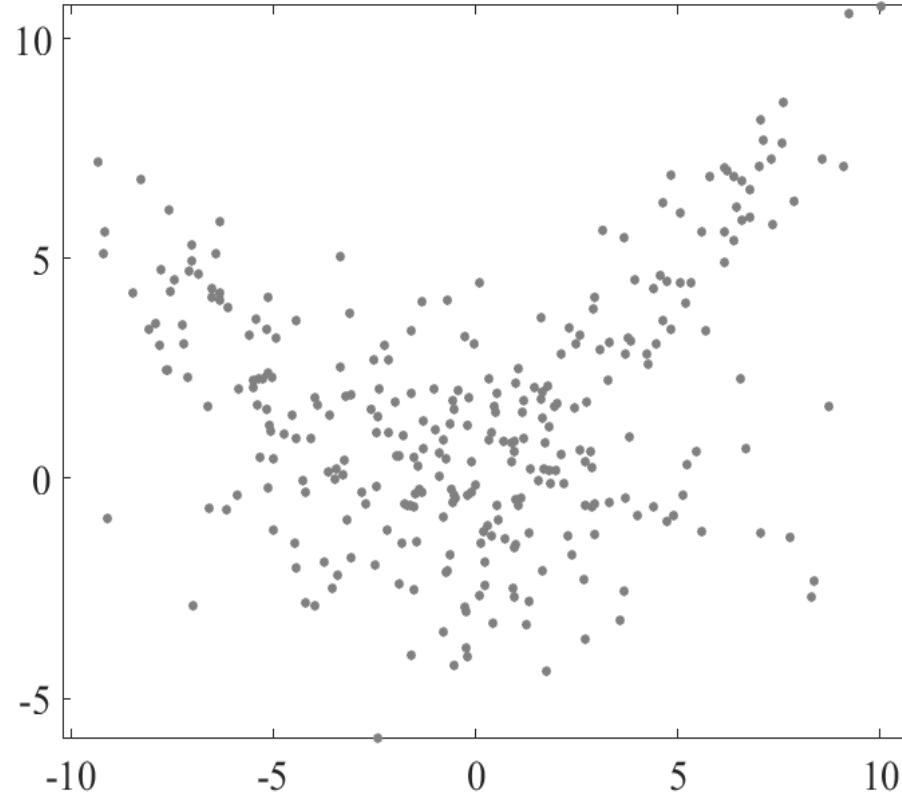
$C = 80\%$; $M = 20\%$



UC Clustering-Examples

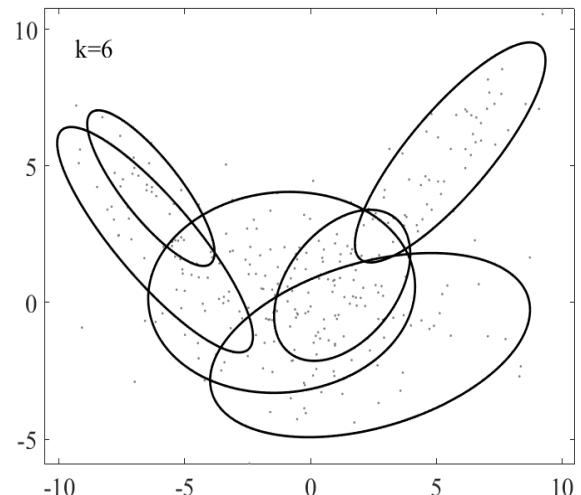
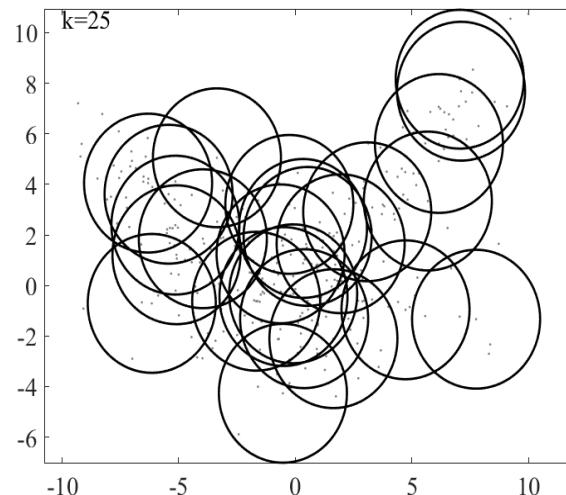
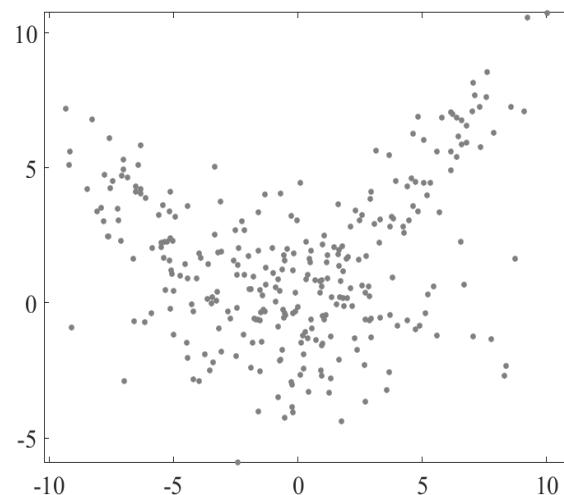
Mixtures of Gaussians – Expectation Maximization (EM)

What if we can *tune* or *select* or *infer* the number of gaussians?



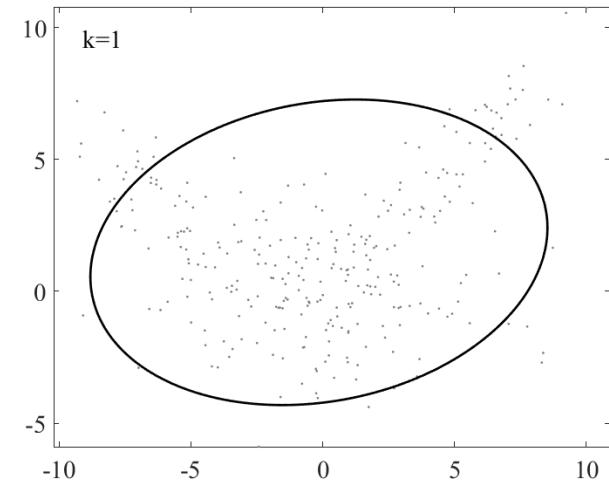
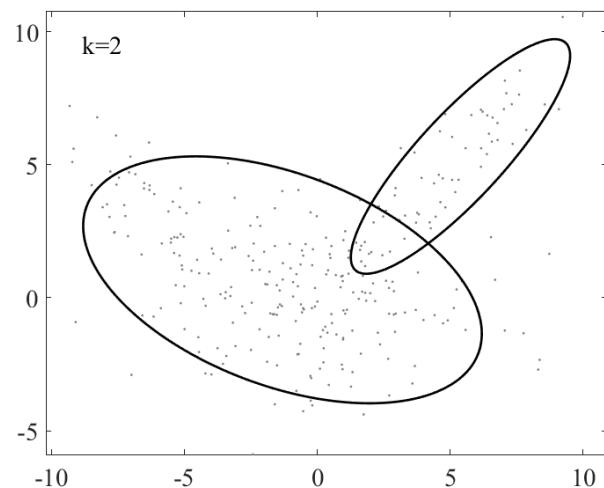
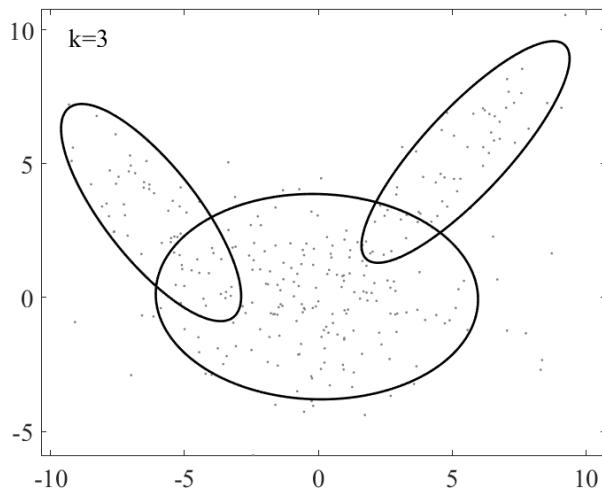
UC Clustering-Examples

Mixtures of Gaussians – Expectation Maximization (EM)



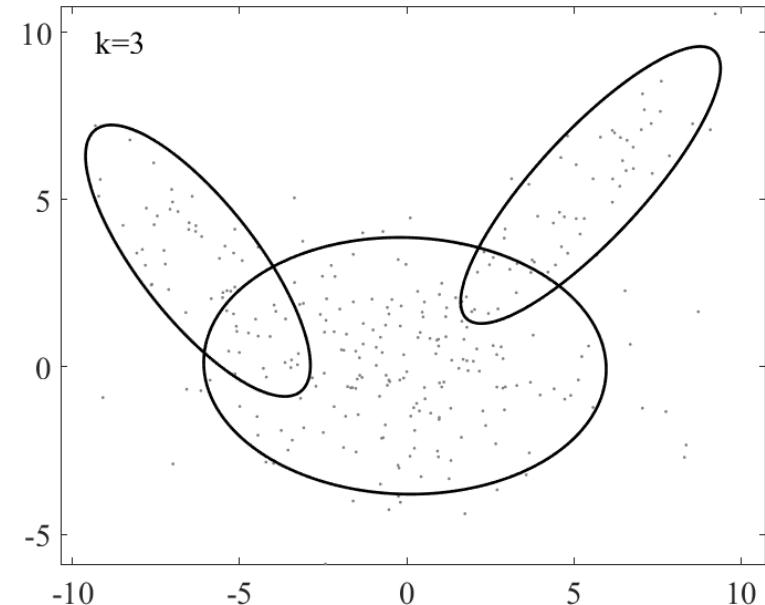
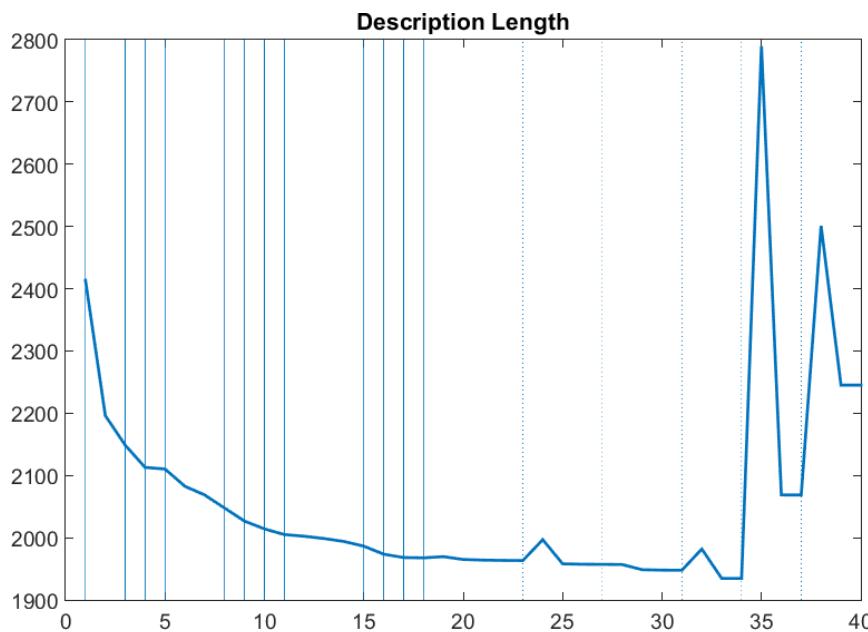
UC Clustering-Examples

Mixtures of Gaussians – Expectation Maximization (EM)



UC Clustering-Examples

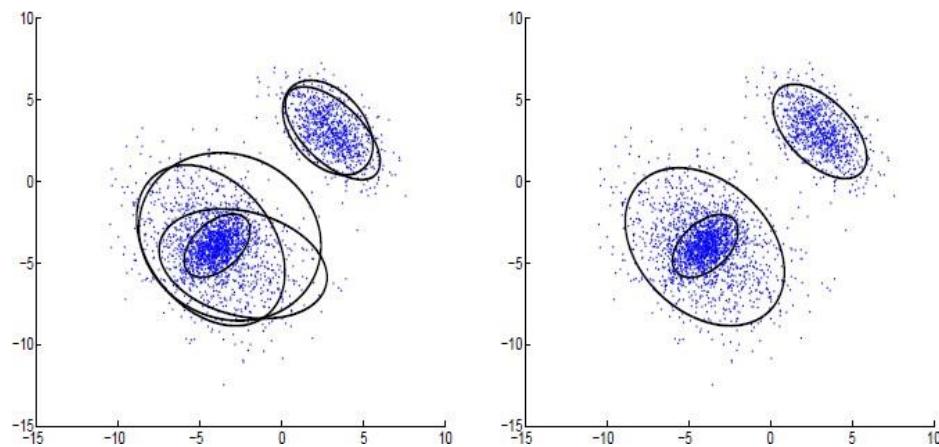
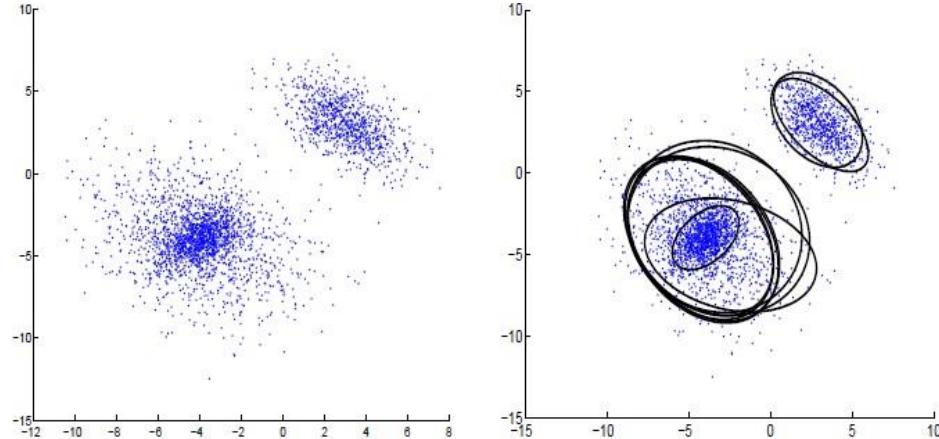
Mixtures of Gaussians – Expectation Maximization (EM)



UC Clustering-Examples

Mixtures of Gaussians – Expectation Maximization (EM)

No *a priori* knowledge about
the data!

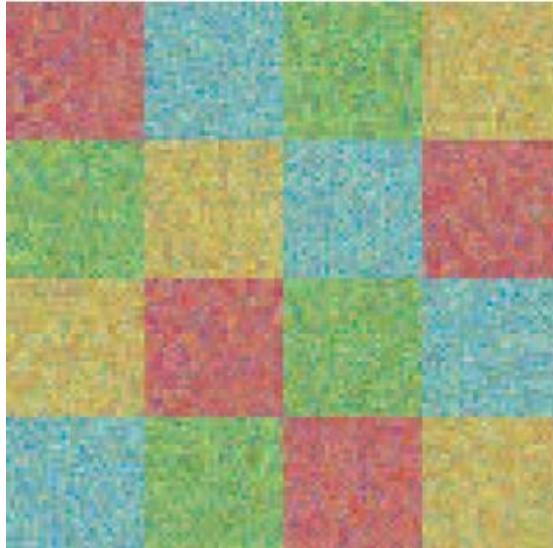


UC Clustering-Examples

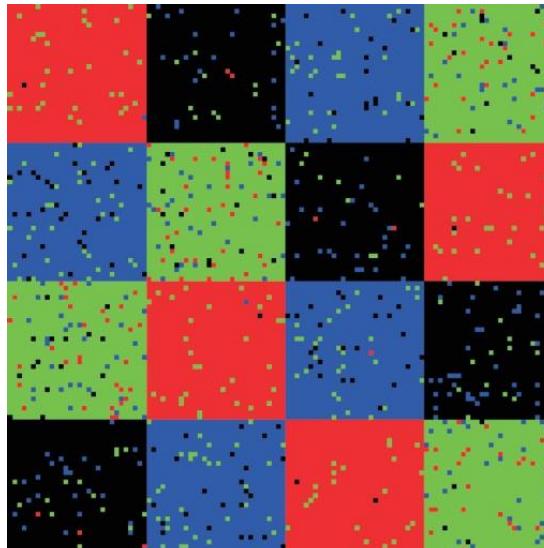
Can we improve the segmentation results?

YES! If we include information from the neighbourhood.

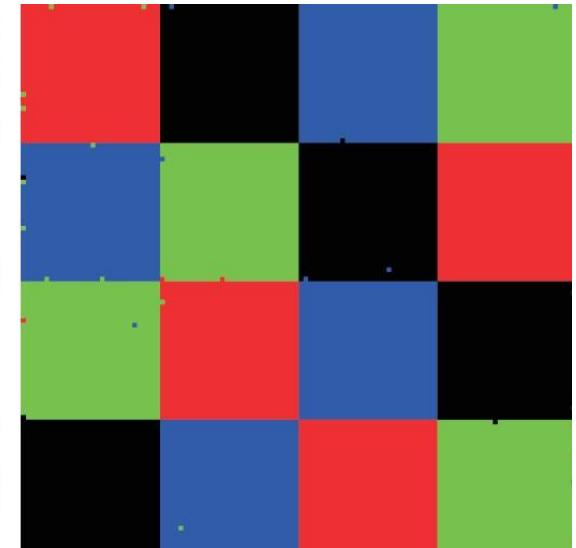
Original



EM



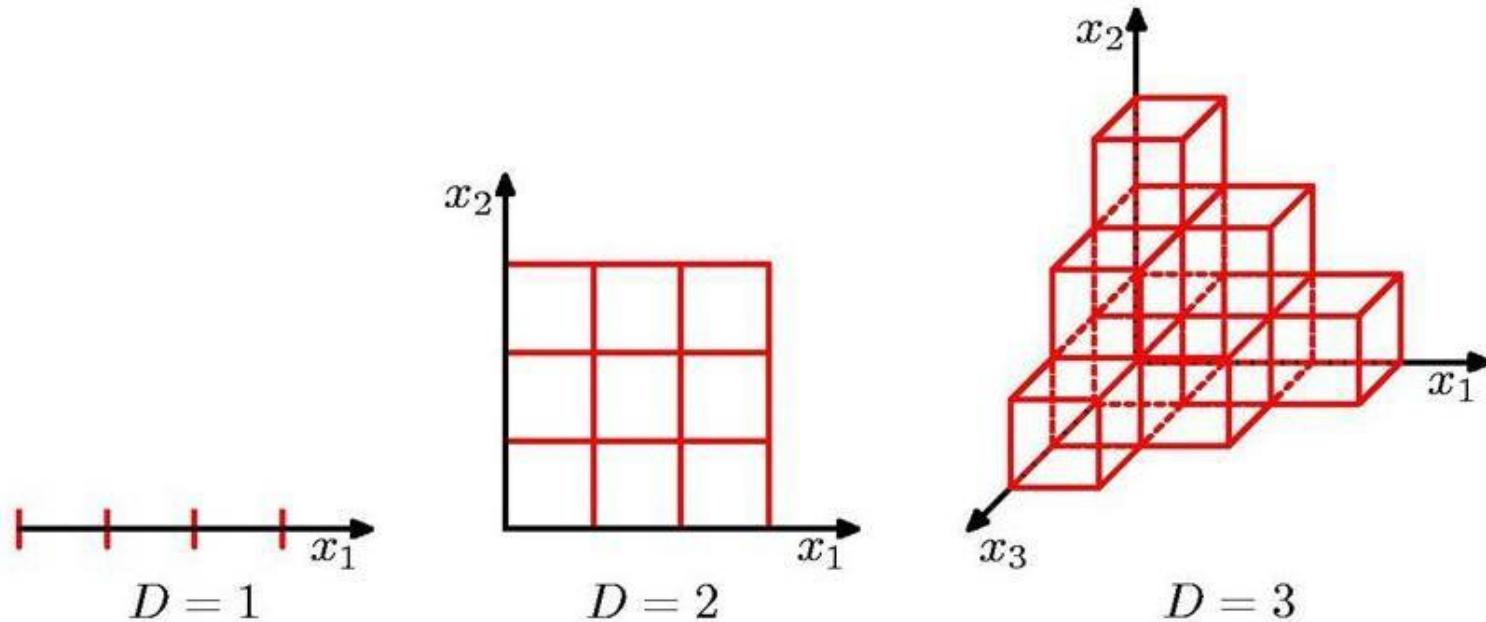
Iterative Conditional Modes
(ICM)



*Houston, we have a few
problems!*

Problems

1) Curse of dimensionality (I)



Problems

Curse of dimensionality

The **Golden rule** would be:

$$n \approx 10^d$$

So that the **number** of points (n) needed for a reliable training (and testing):

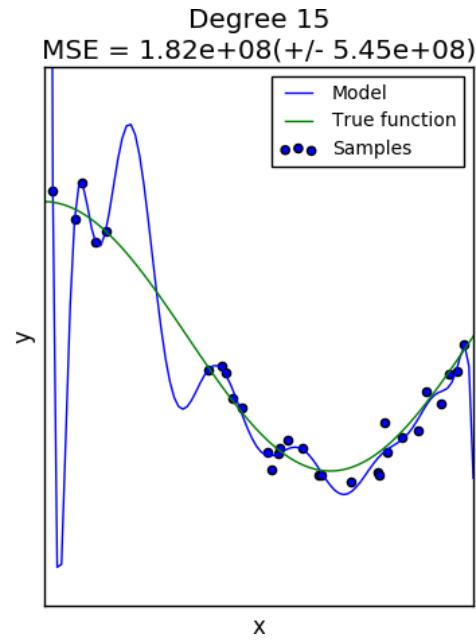
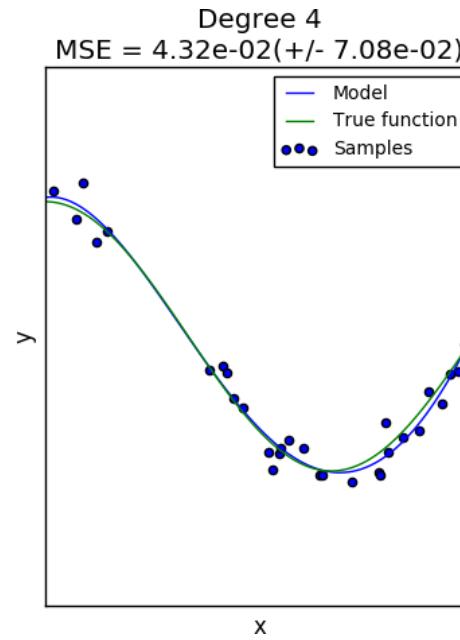
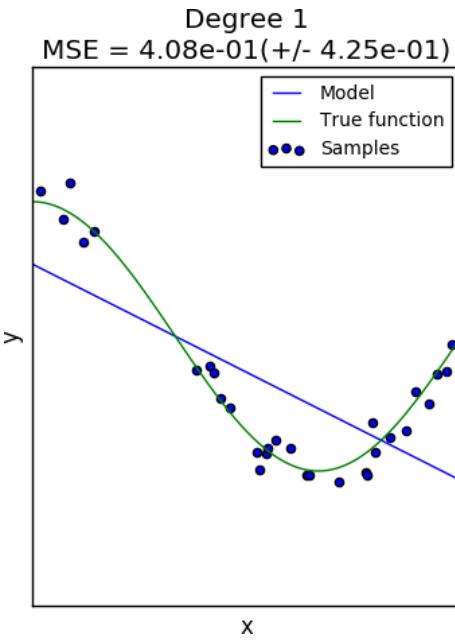
$$d = 1 \rightarrow n \geq 10 \text{ data points}$$

$$d = 2 \rightarrow n \geq 10^2 = 100 \text{ data points}$$

$$d = 3 \rightarrow n \geq 10^3 = 1000 \text{ data points}$$

Problems

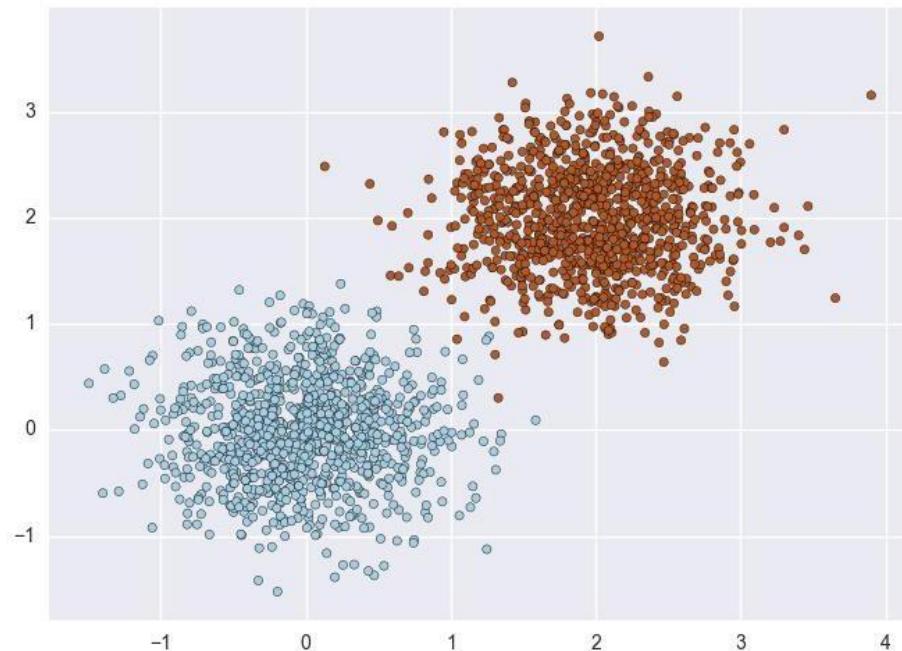
2) Underfitting – Fair fitting - Overfitting



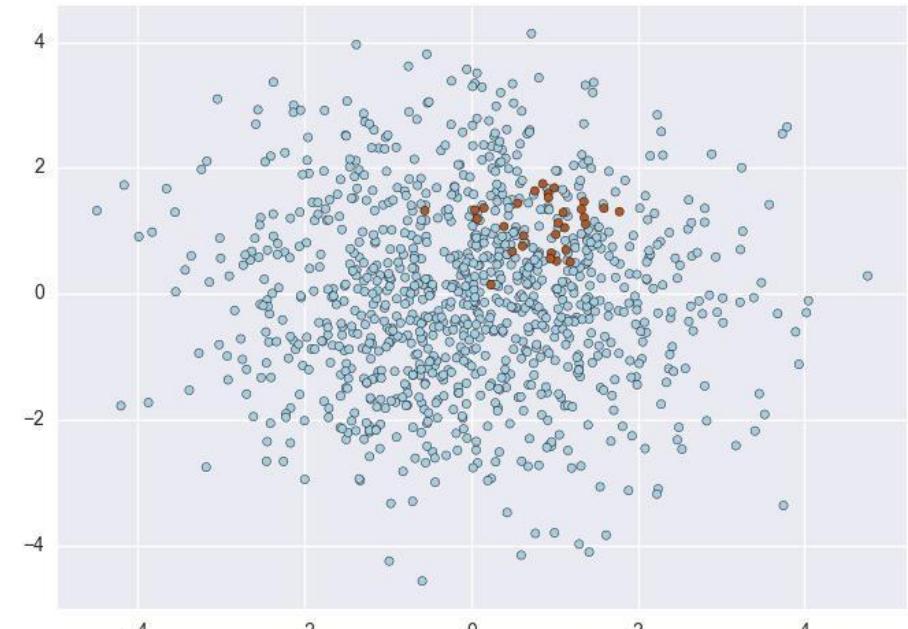
Problems

3) Unbalance and overlapping

BEST CASE SCENARIO



WORST CASE SCENARIO



Problems

Unbalance (I)

What can we do?

We can define new measures of “Accuracy” that are better behaved for unbalance.

$$G_m(\%) = [\sqrt{TP_r \cdot TN_r}] \cdot 100$$

$$A_{cw}(\%) = \{[\omega \cdot TP_r] + [(1 - \omega) \cdot TN_r]\} \cdot 100$$

You can give more importance to
“P” or to “N”

Here: $TP_r = \frac{TP}{TP+FN}$; $TN_r = \frac{TN}{TN+FP}$

- a. $0 \leq \omega \leq 1$
b. $\omega + (1 - \omega) = 1$

Problems

Unbalance (II)

Imagine a ***melanoma*** detection problem. We acquire a multi-spectral image of a skin lesion.



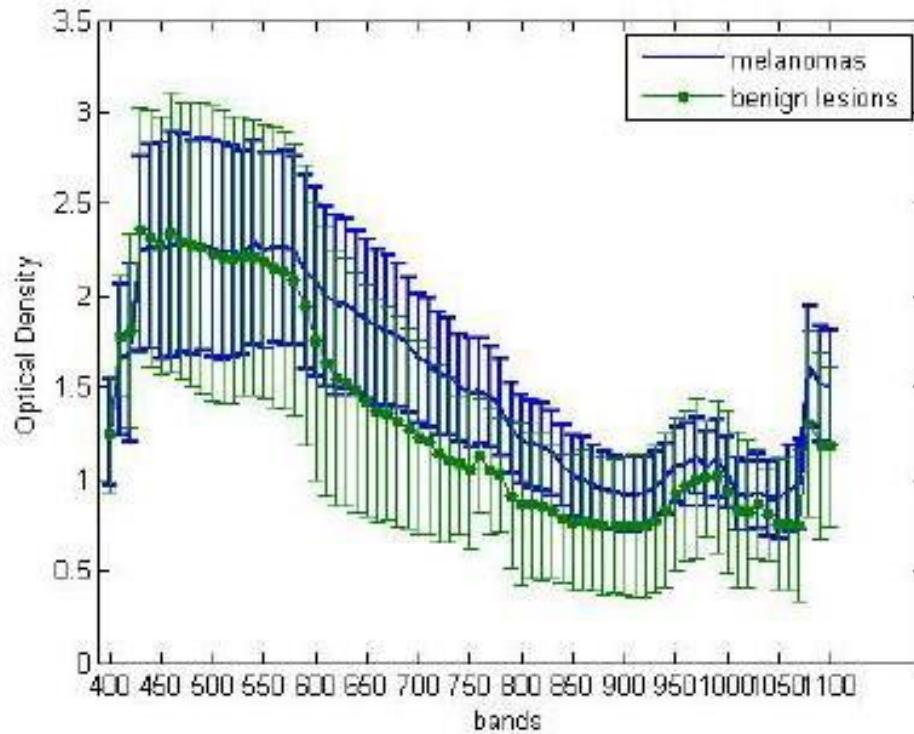
Liquid Crystal Tunable Filter (LCTF)

32 λs in [400, 720]nm; **32 λs** in [650, 1100]nm

Problems

Unbalance (III)

The spectral curves associated to melanoma and other skin lesions are different.



Problems

Unbalance (IV)

Imagine we have 520 patients, 20 of which suffer from melanoma.

Let's consider the following *confusion matrices*:

**30 % of mistake
in the class of interest!**
No, No, No !!!!!!!

	“Melanoma”	No “Melanoma”
Classified as “Melanoma”	TP = 18	FP = 15
Classified as “No Melanoma”	FN = 2	TN = 485

$$Acc (\%) = \left(\frac{TP+TN}{TP+TN+FP+FN} \right) \cdot 100 =$$

$$= \left(\frac{18 + 485}{520} \right) \cdot 100 = \mathbf{96.73 \%}$$

	“Melanoma”	No “Melanoma”
Classified as “Melanoma”	TP = 14	FP = 15
Classified as “No Melanoma”	FN = 6	TN = 485

$$Acc (\%) = \left(\frac{TP+TN}{TP+TN+FP+FN} \right) \cdot 100 =$$

NO WAY!!

$$= \left(\frac{14 + 485}{520} \right) \cdot 100 = \mathbf{95.96 \%}$$

Problems

Unbalance (V)

Remember: 520 patients, 20 of which suffer from melanoma.

Let's consider the following *confusion* matrices:

	“Melanoma”	No “Melanoma”
Classified as “Melanoma”	TP = 18	FP = 15
Classified as “No Melanoma”	FN = 2	TN = 485

	“Melanoma”	No “Melanoma”
Classified as “Melanoma”	TP = 14	FP = 15
Classified as “No Melanoma”	FN = 6	TN = 485

$$TP_r = \frac{TP}{TP+FN}; \quad TN_r = \frac{TN}{TN+FP}$$

$$G_m(\%) = [\sqrt{TP_r \cdot TN_r}] \cdot 100 = \\ = \sqrt{\left(\frac{18}{20}\right) \cdot \left(\frac{485}{500}\right)} \cdot 100 = \mathbf{93.43\%}$$

BETTER
BEHAVIOUR!!!

$$G_m(\%) = [\sqrt{TP_r \cdot TN_r}] \cdot 100 = \\ = \sqrt{\left(\frac{14}{20}\right) \cdot \left(\frac{485}{500}\right)} \cdot 100 = \mathbf{82.40\%}$$

Problems

Unbalance (VI)

Bear in mind: Not all mistakes might be **equally IMPORTANT!!!**

	"Melanoma"	No "Melanoma"
Classified as "Melanoma"	TP = 18	FP = 15
Classified as "No Melanoma"	FN = 2	TN = 485

	"Melanoma"	No "Melanoma"
Classified as "Melanoma"	TP = 14	FP = 15
Classified as "No Melanoma"	FN = 6	TN = 485

↑ NO. OF "FNs"
**IS THE WORST MISTAKE OF THEM ALL
IN THIS CASE**

You are telling the patient he/she does not
have melanoma but he/she has it!

References

- 1 Christopher Bishop, Pattern Recognition and Machine Learning, Information Science and Statistics Series, Springer Verlag, 2006.
- 2 Sergios Theodoridis, Konstantinos koutroumbas, Pattern Recognition, Academic press, 2009.