

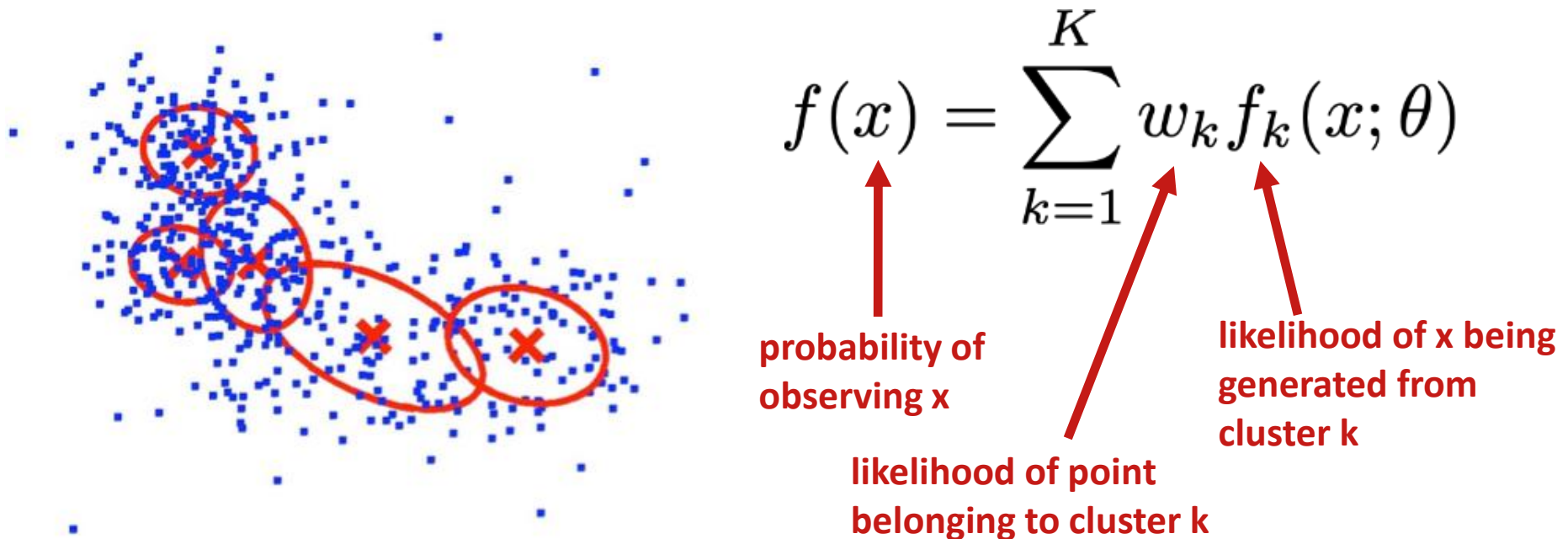
# **Clusters probabilísticos**

Dr. Raimundo Sánchez  
raimundo.sanchez@uai.cl  
@raimun2

# Métodos probabilísticos

Los métodos probabilísticos proporcionan una descripción distribucional completa para cada componente, generando clústeres flexibles.

Es decir, dado un modelo, cada punto tiene un vector de K probabilidades de pertenencia

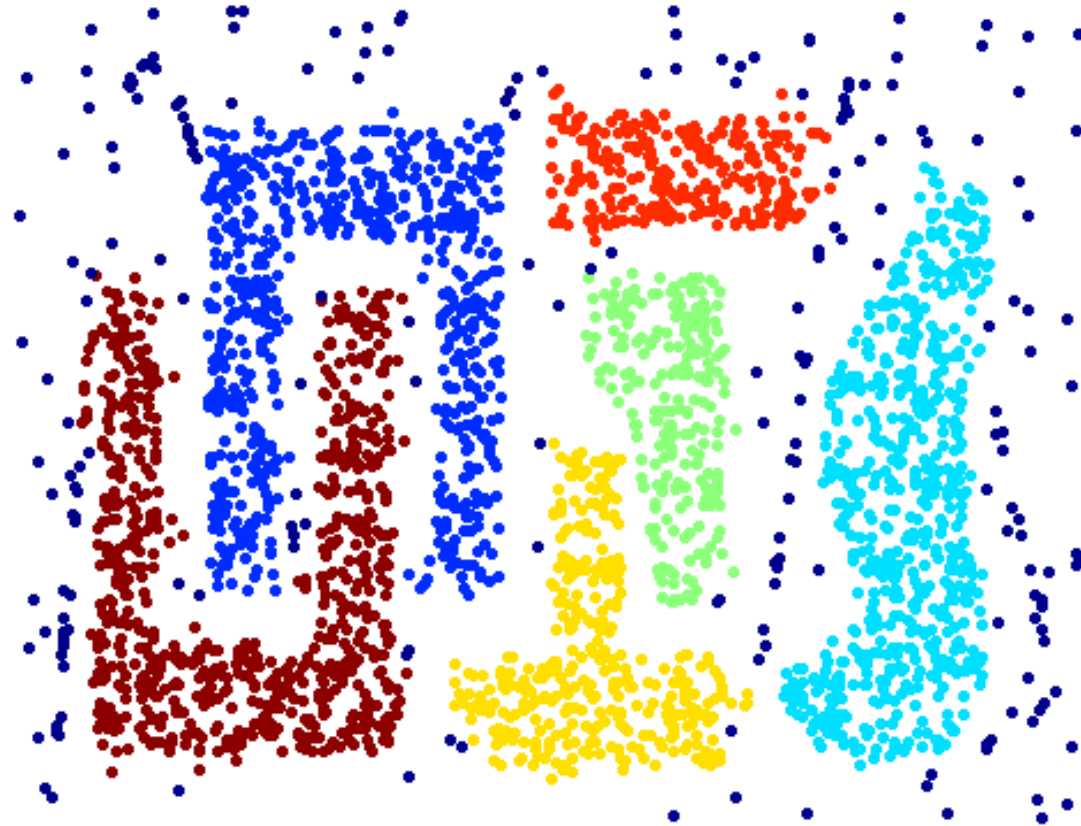


# **Métodos basados en densidad:**

## **DBSCAN**

# DBSCAN

DBSCAN es un algoritmo de clusters de densidad, donde dado un conjunto de puntos en algún espacio, agrupa los puntos que están empaquetados de cerca.



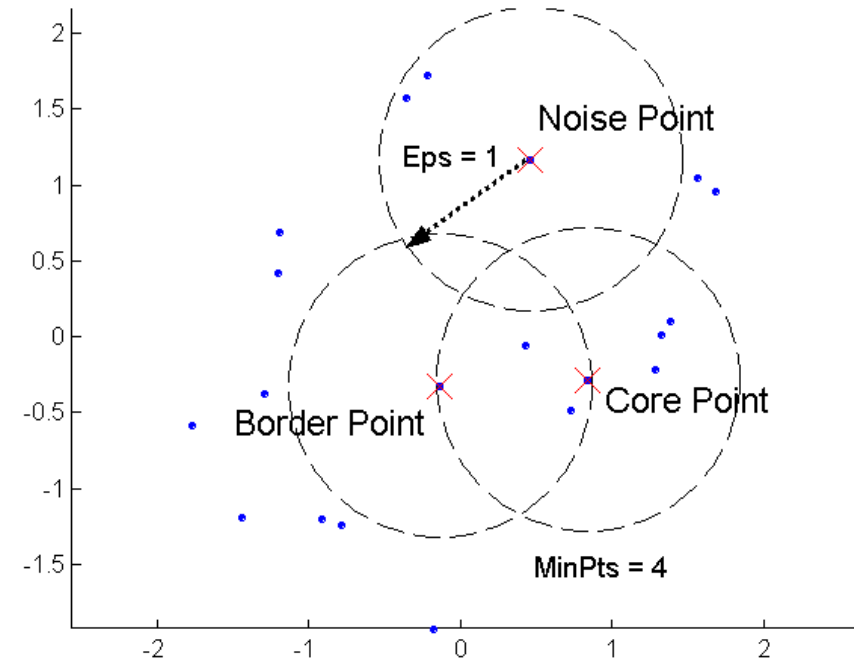
# Definiciones

**Densidad** es el número de puntos dentro de un determinado **radio (Eps)**

Un **punto central** es un punto que tiene el mismo **número de puntos (MinPts)** dentro de una o mas esferas definida por Eps (incluido él mismo).

Un **punto fronterizo** tiene menos puntos que MinPts dentro de Eps, pero está cerca de un punto central

Un **punto de ruido** es cualquier punto que no sea un punto central o un punto fronterizo.



# Algoritmo

- Definir **Eps** y **MinPts**
- Determinar puntos **centro**, **frontera**, y **ruido**
- Eliminar puntos de ruido
- Aplicar el siguiente algoritmo de clusters

$current\_cluster\_label \leftarrow 1$

**for** all core points **do**

**if** the core point has no cluster label **then**

$current\_cluster\_label \leftarrow current\_cluster\_label + 1$

        Label the current core point with cluster label  $current\_cluster\_label$

**end if**

**for** all points in the  $Eps$ -neighborhood, except  $i^{th}$  the point itself **do**

**if** the point does not have a cluster label **then**

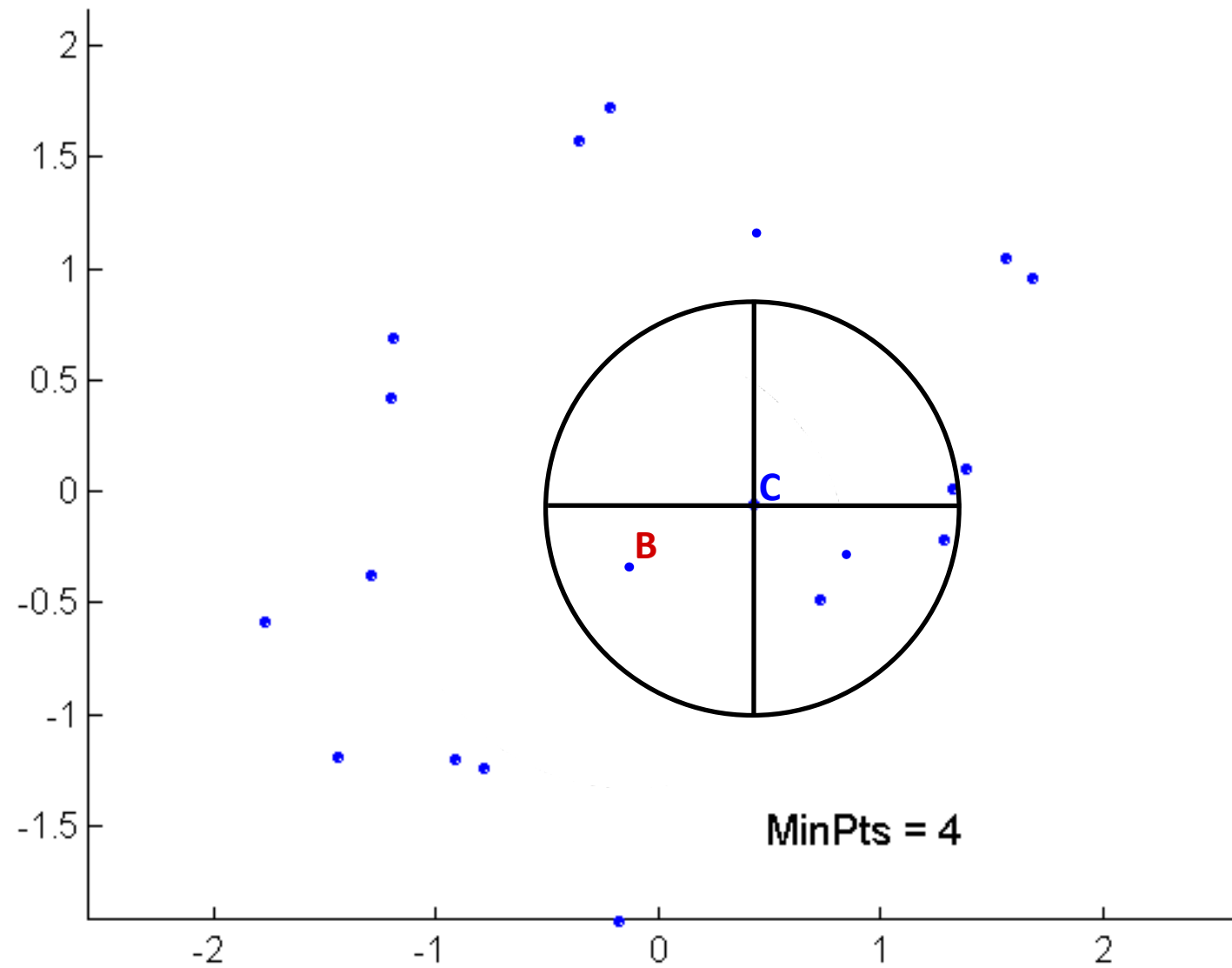
            Label the point with cluster label  $current\_cluster\_label$

**end if**

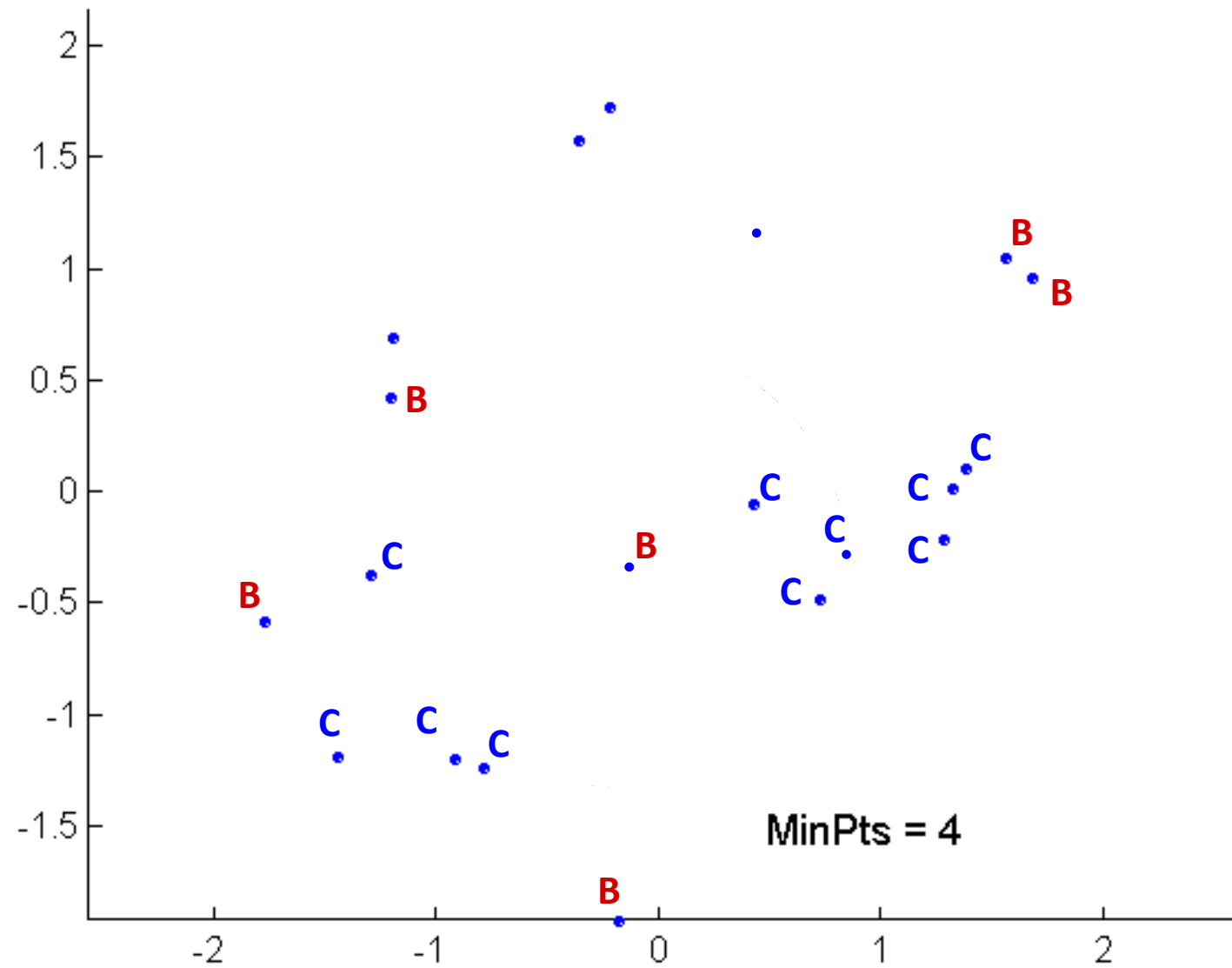
**end for**

**end for**

# Ejemplo

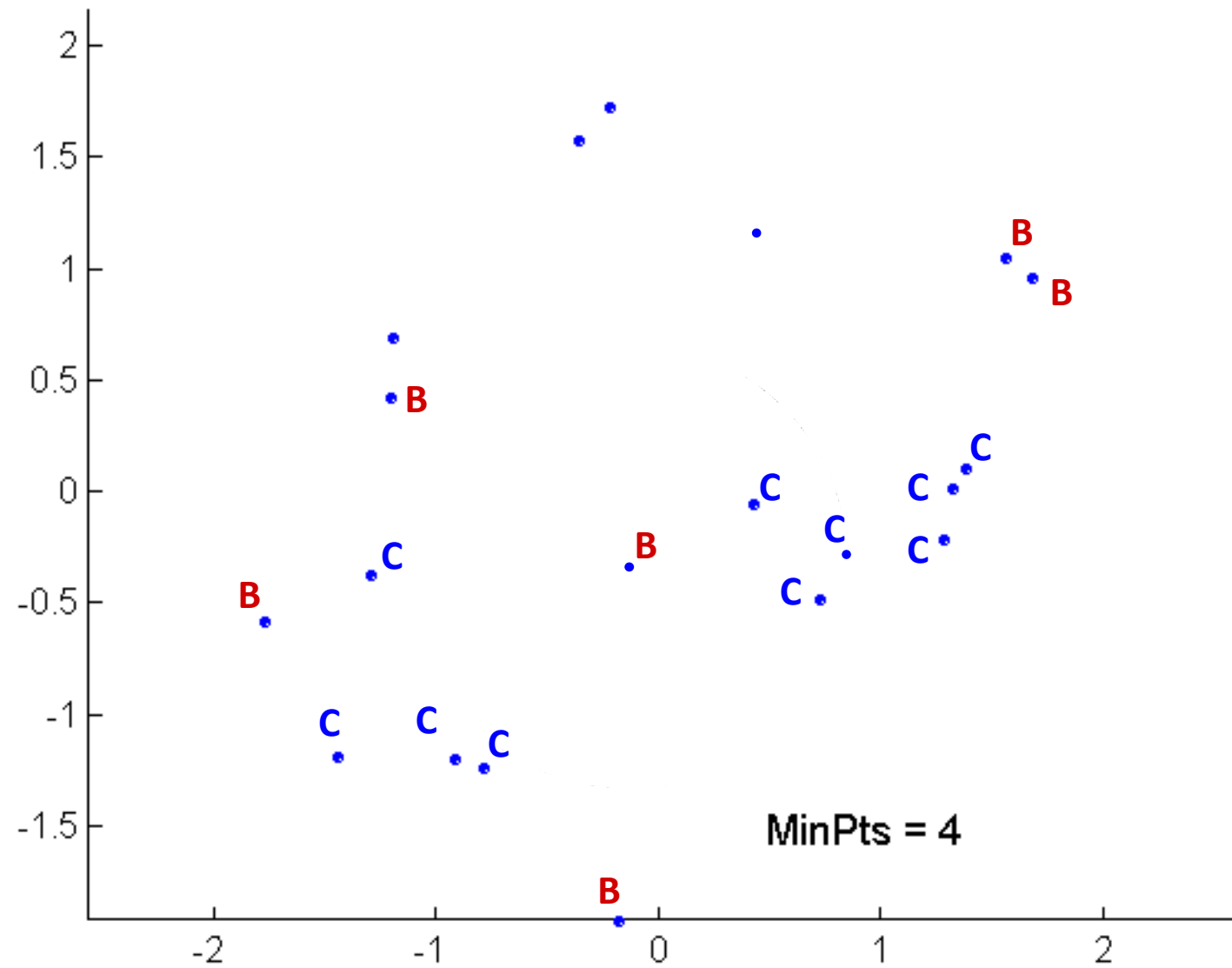


# Ejemplo

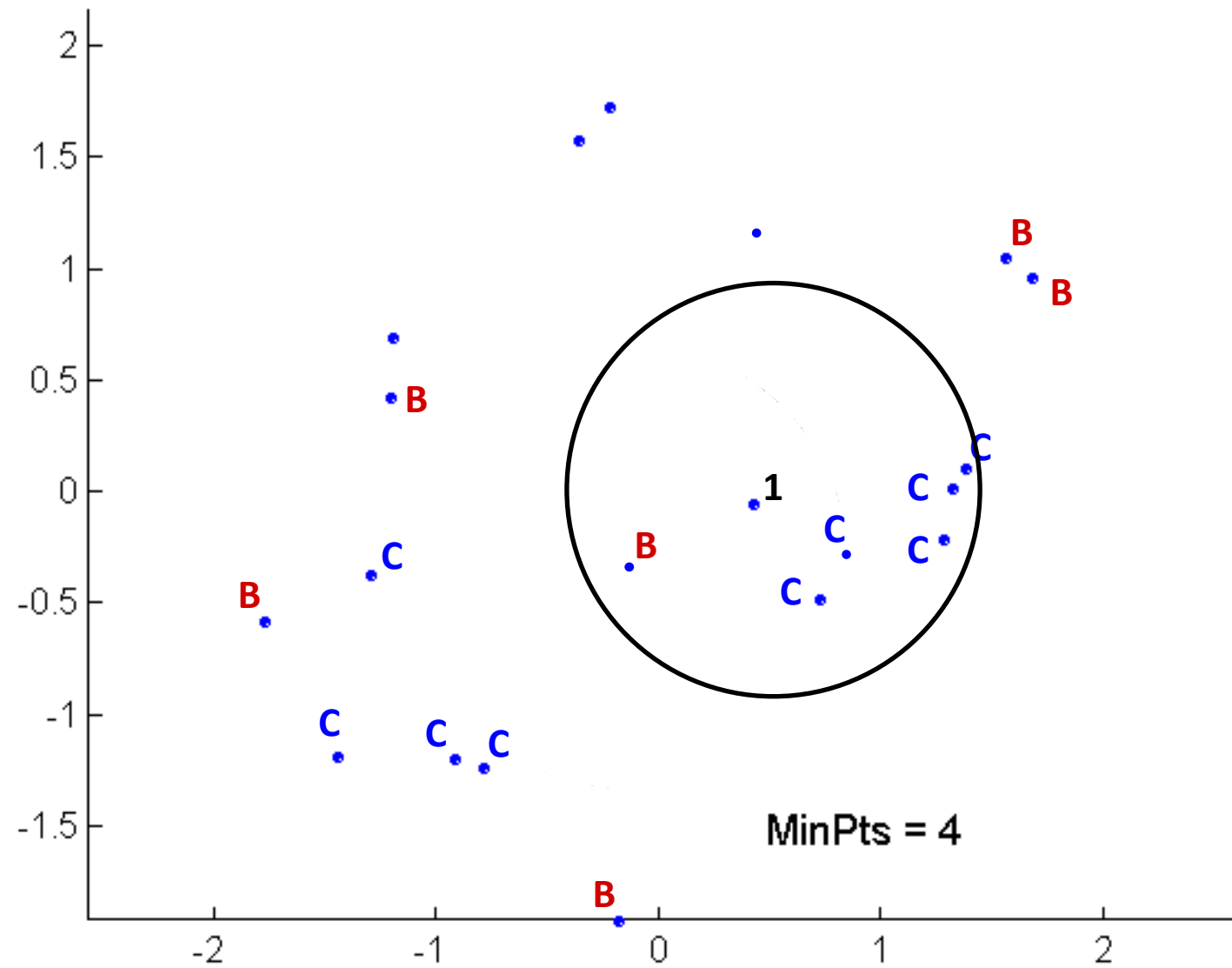




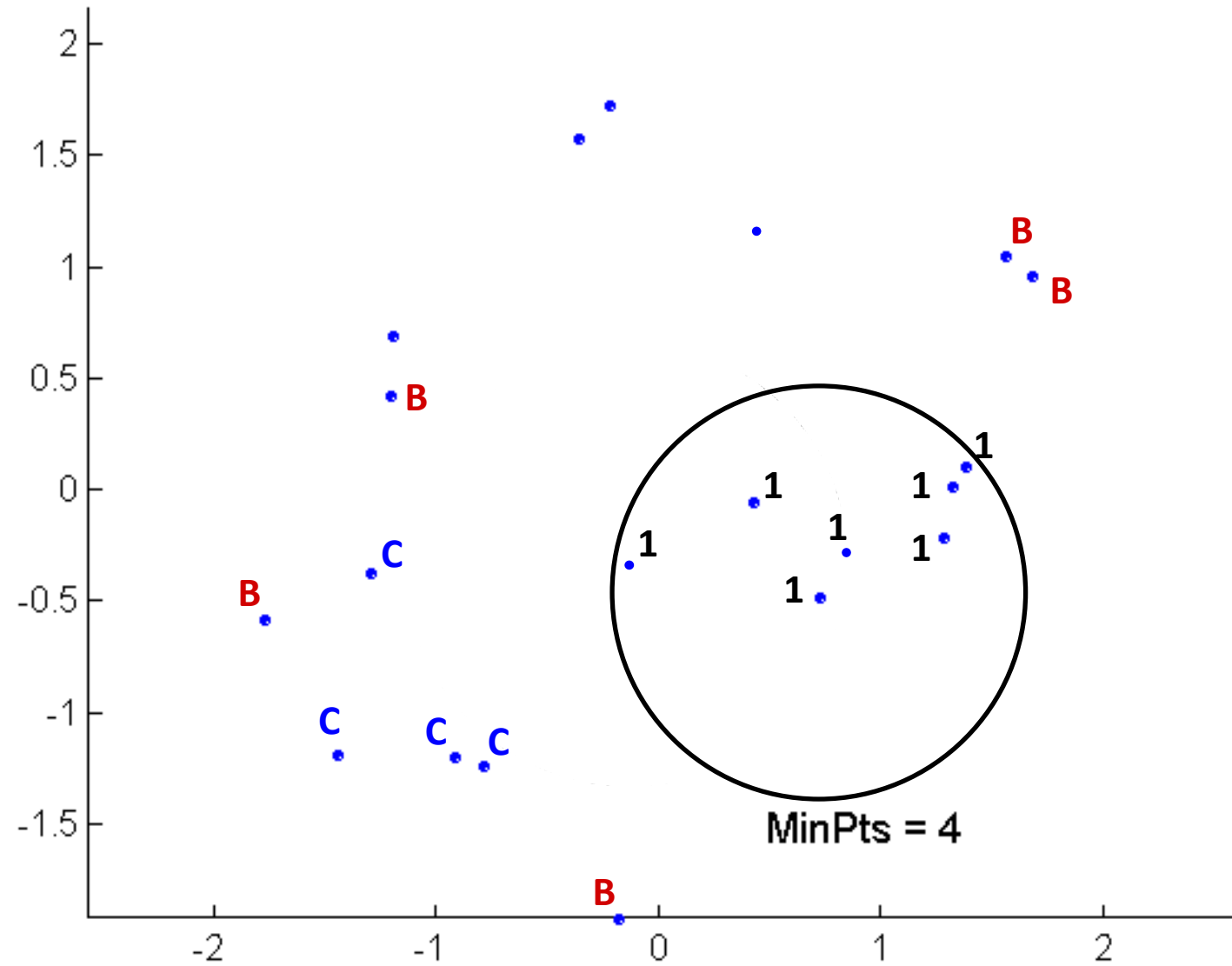
# Ejemplo



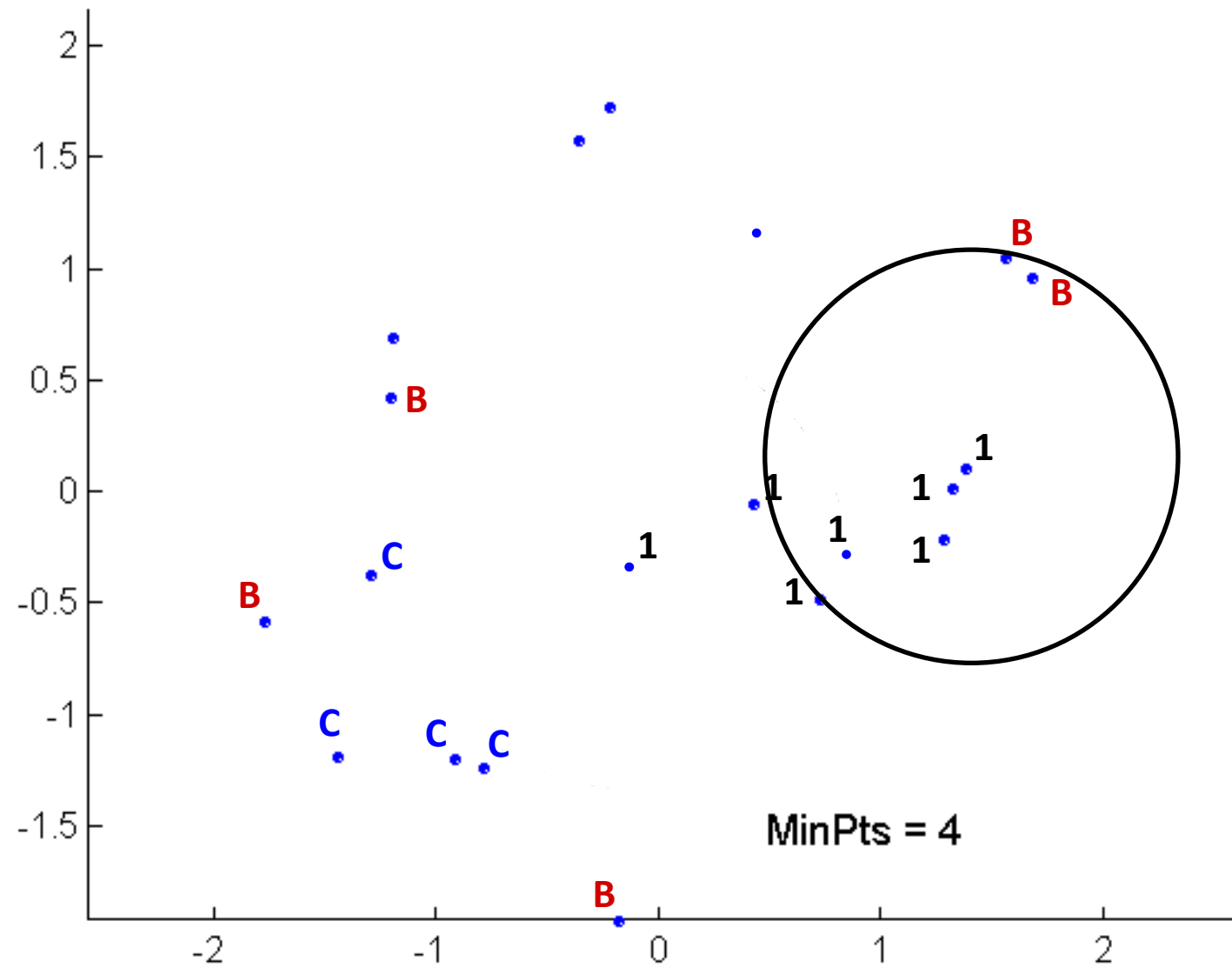
# Ejemplo



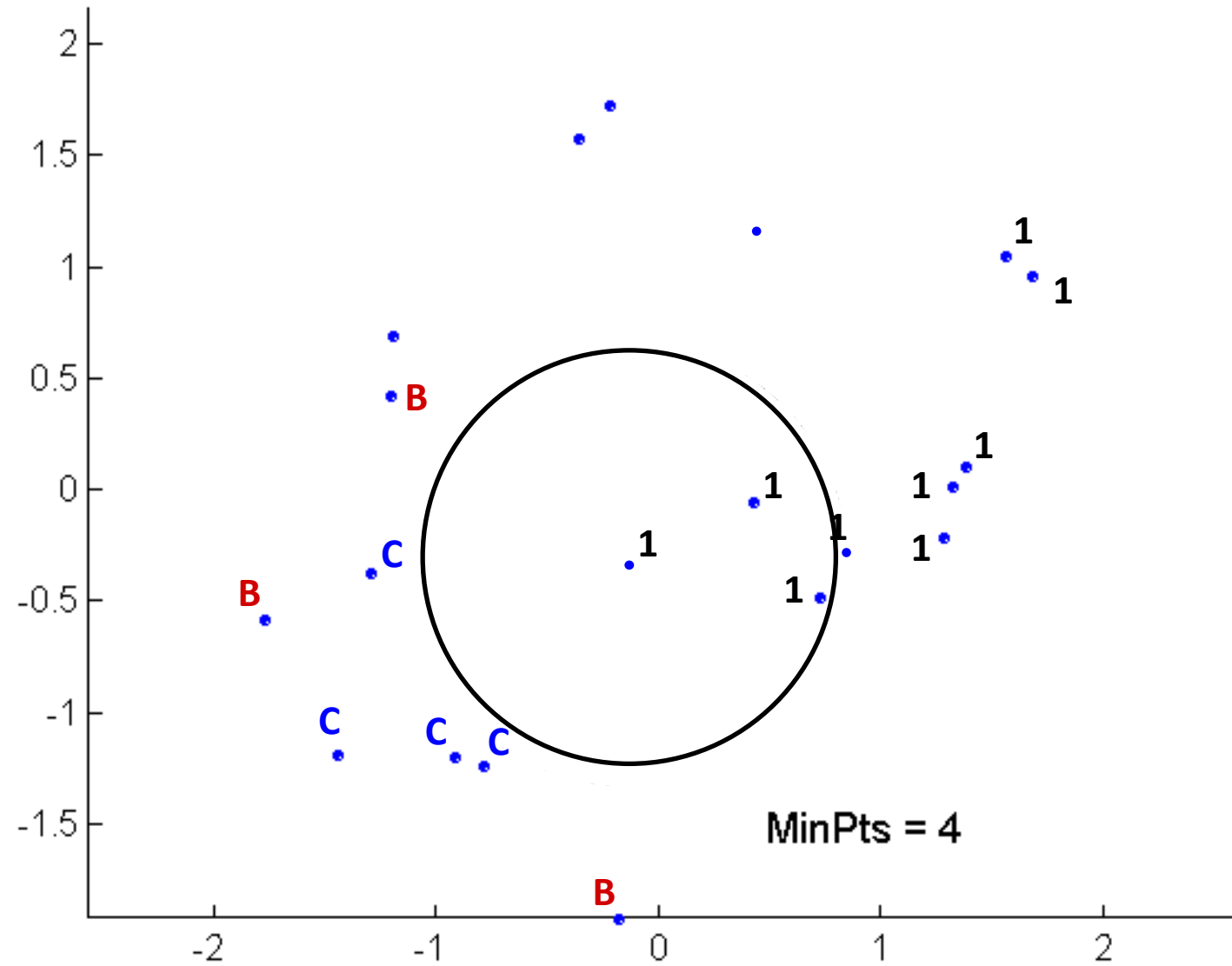
# Ejemplo



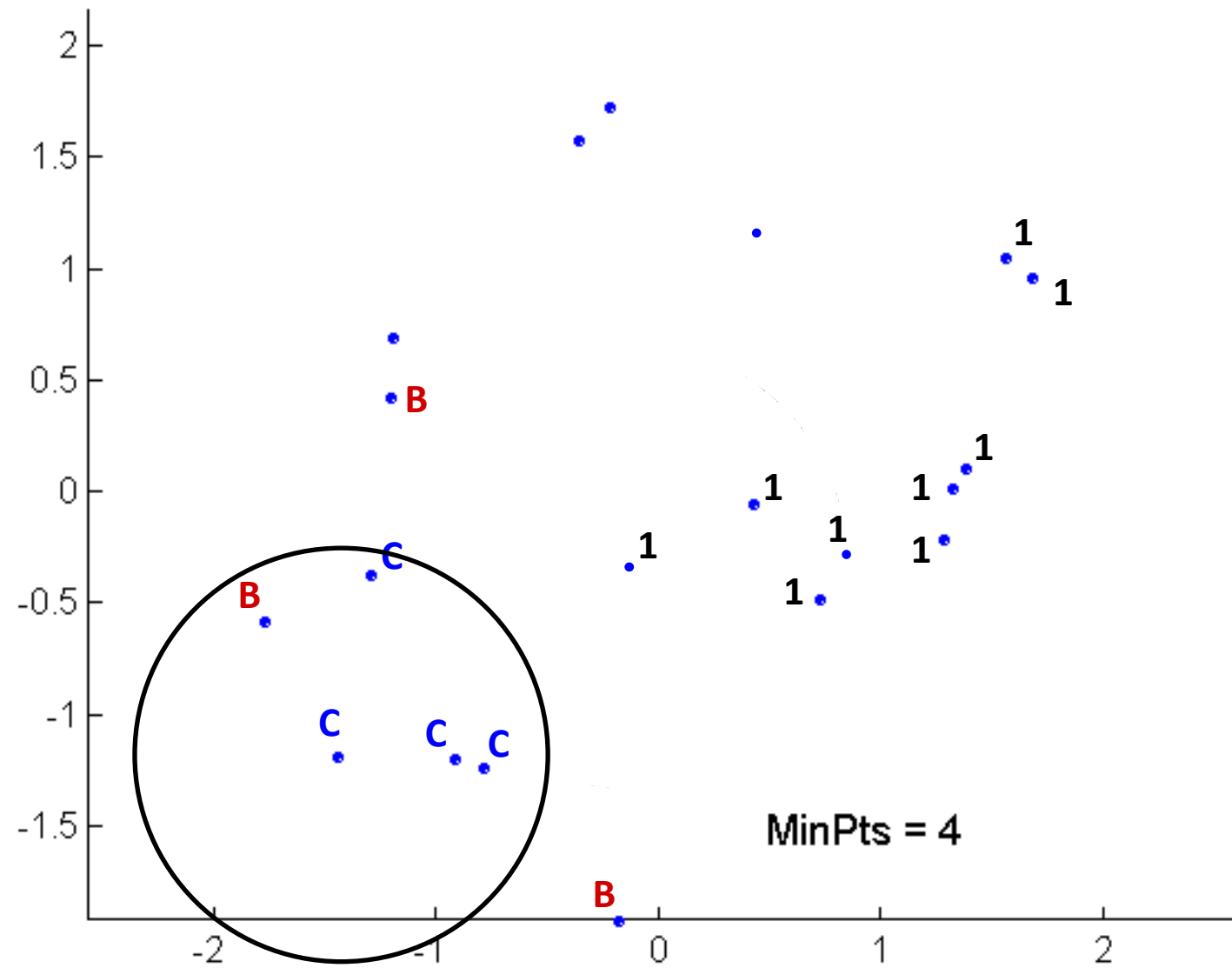
# Ejemplo



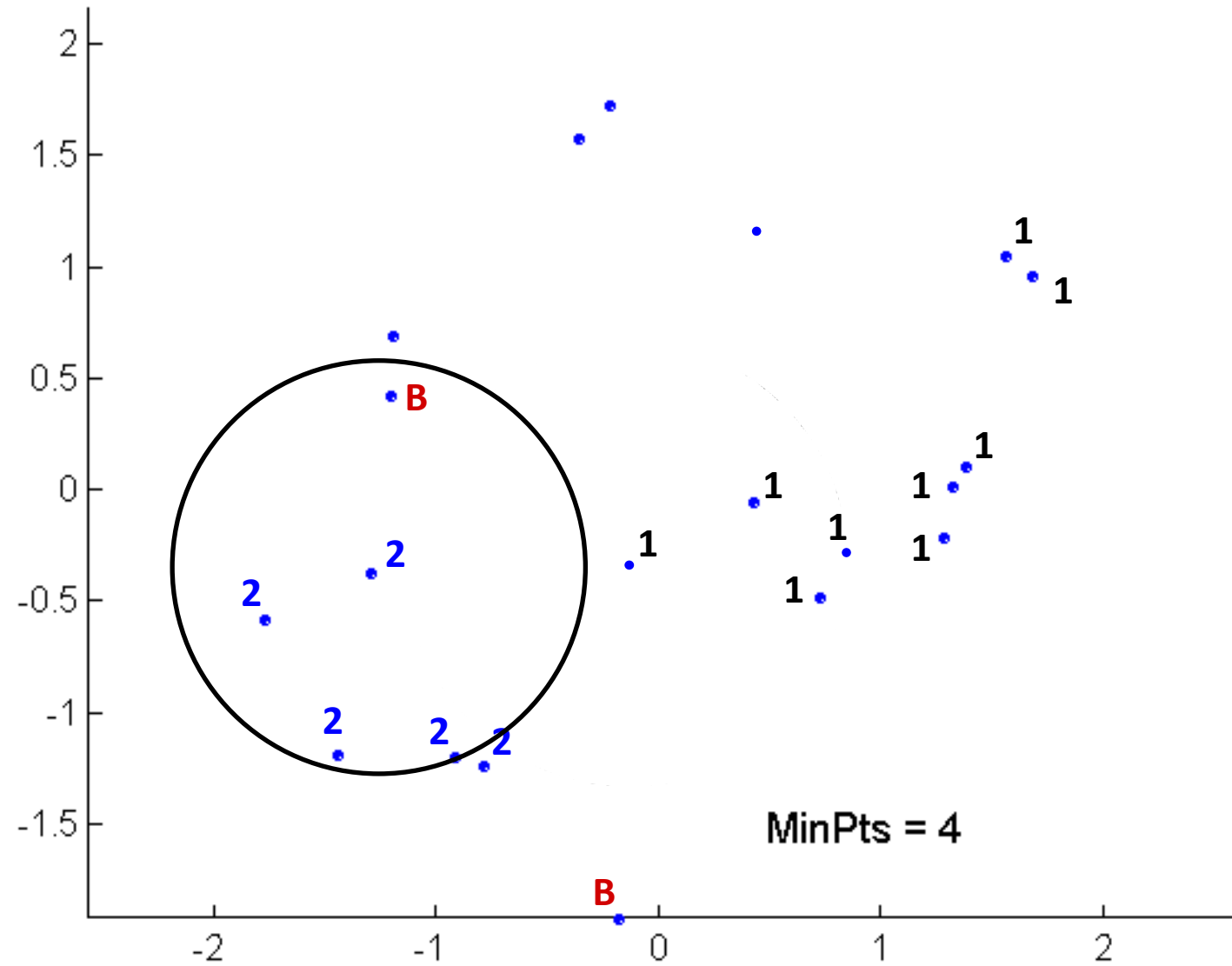
# Ejemplo



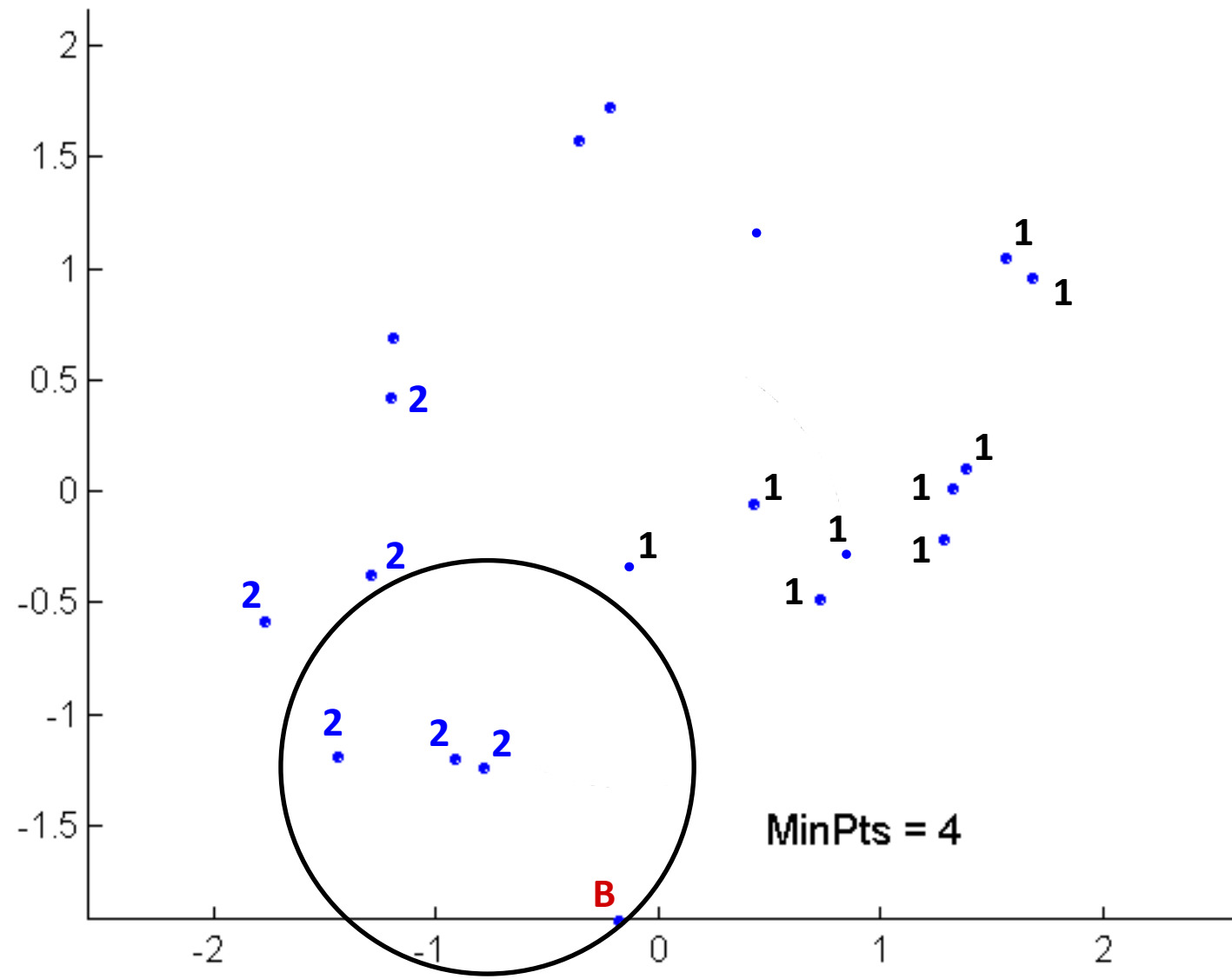
# Ejemplo



# Ejemplo

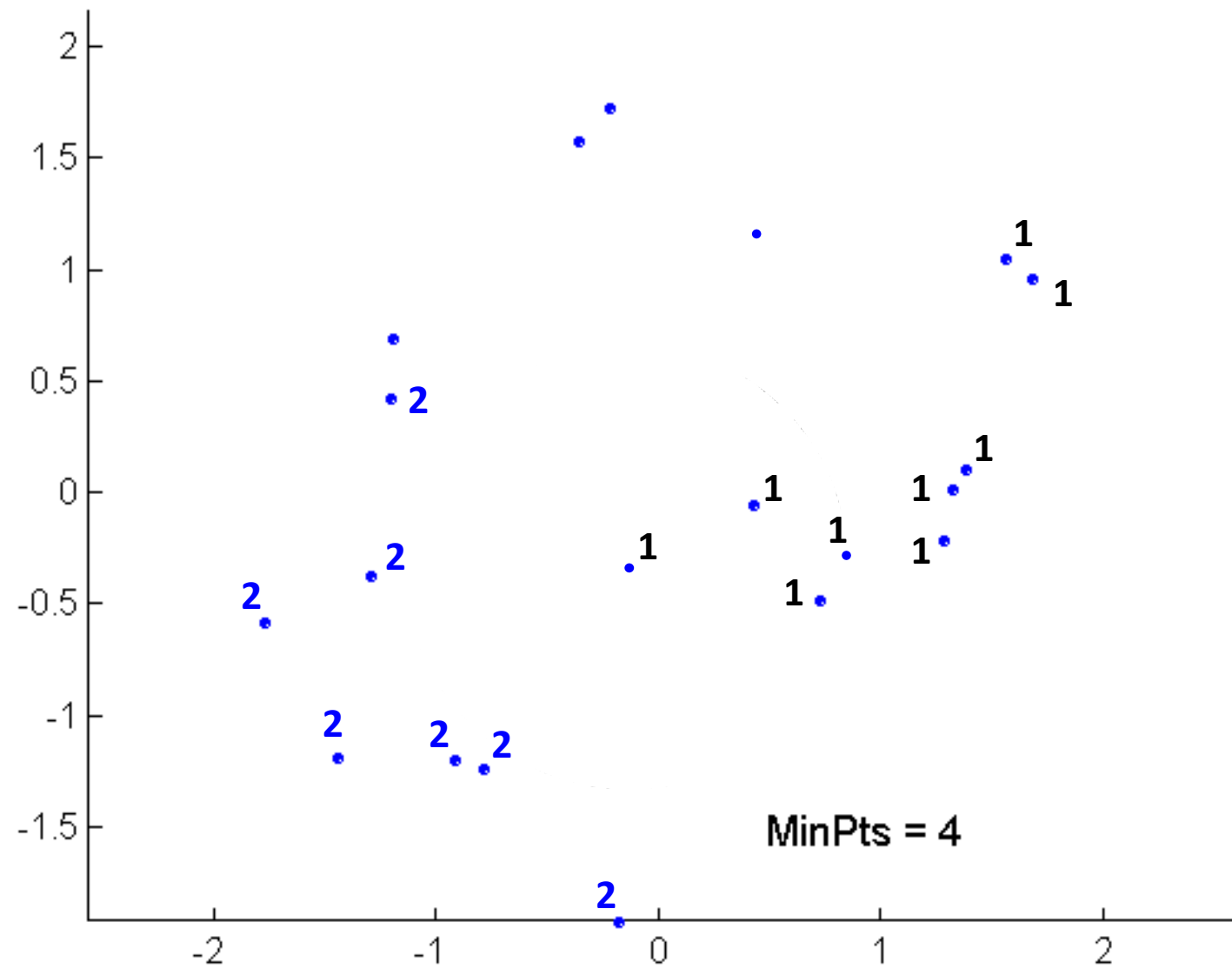


# Ejemplo





# Ejemplo

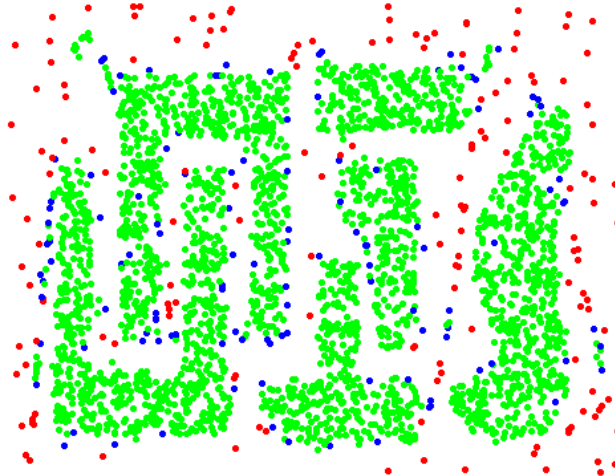


# Ejemplo

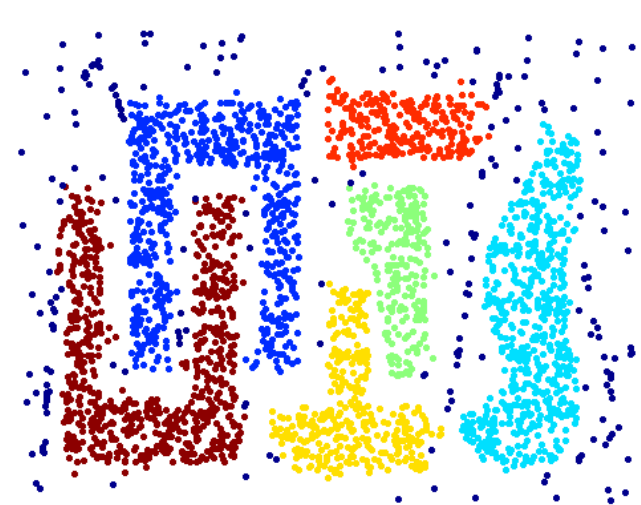
Eps = 10, MinPts = 4



Puntos originales



Tipos de puntos:  
centro, frontera y  
ruido



Clusters

# Fortalezas y debilidades

La complejidad del algoritmo es  $O(n^2)$ , debido a la distancia entre todos los puntos.

## Fortalezas:

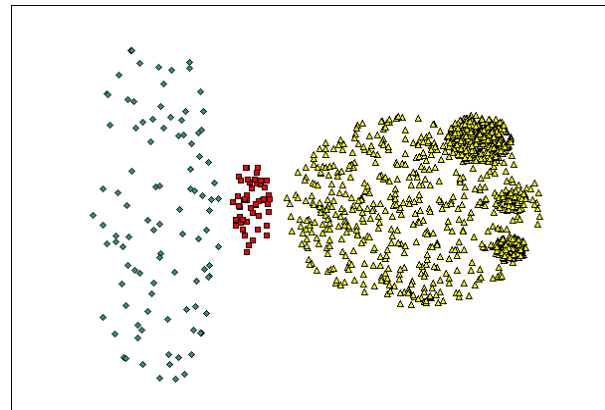
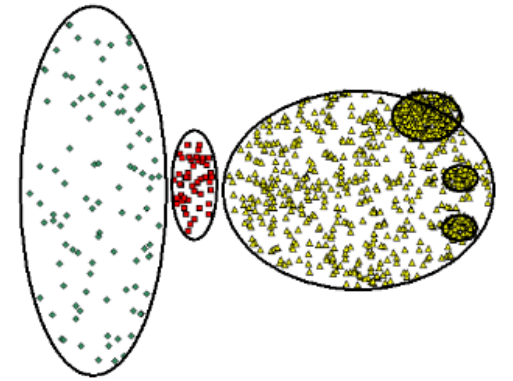
Resistente al ruido

Puede manejar grupos de diferentes formas y tamaños.

## Debilidades:

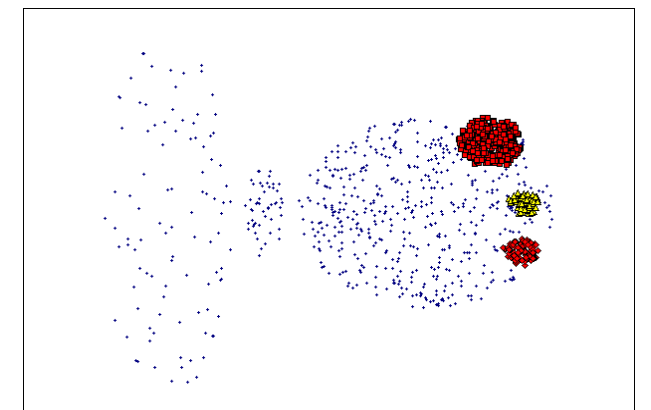
Varianza de densidades

Datos de alta dimensión



MinPts = 4

Eps = 9,75



MinPts = 4

Eps = 9,92

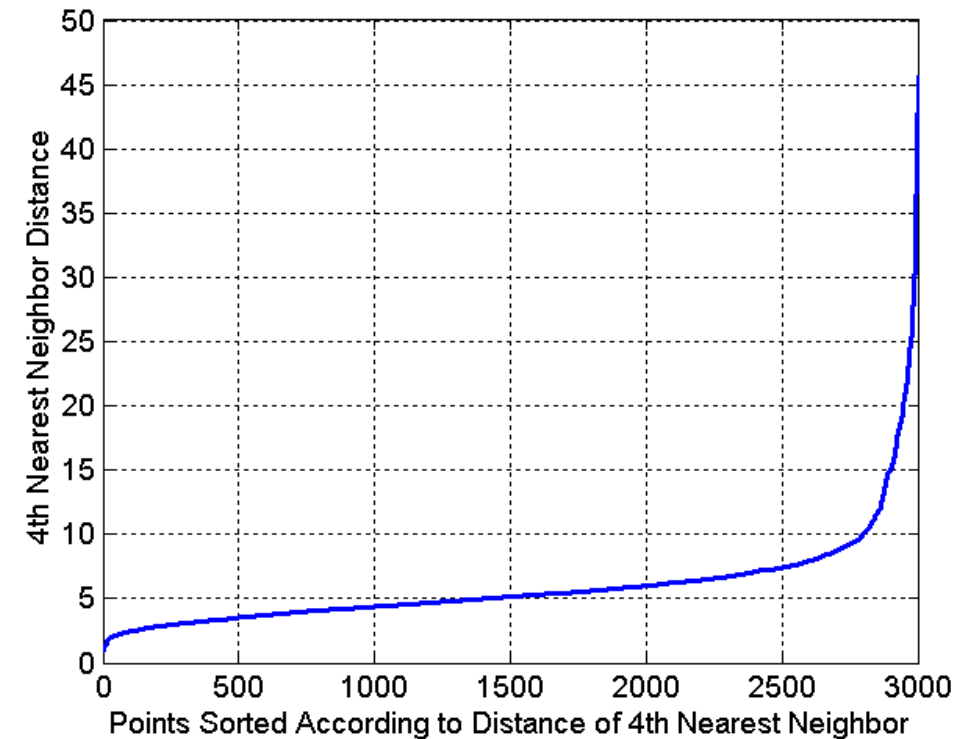
# Ajustes de parámetros

Cómo seleccionar **EPS** y **MinPts**?

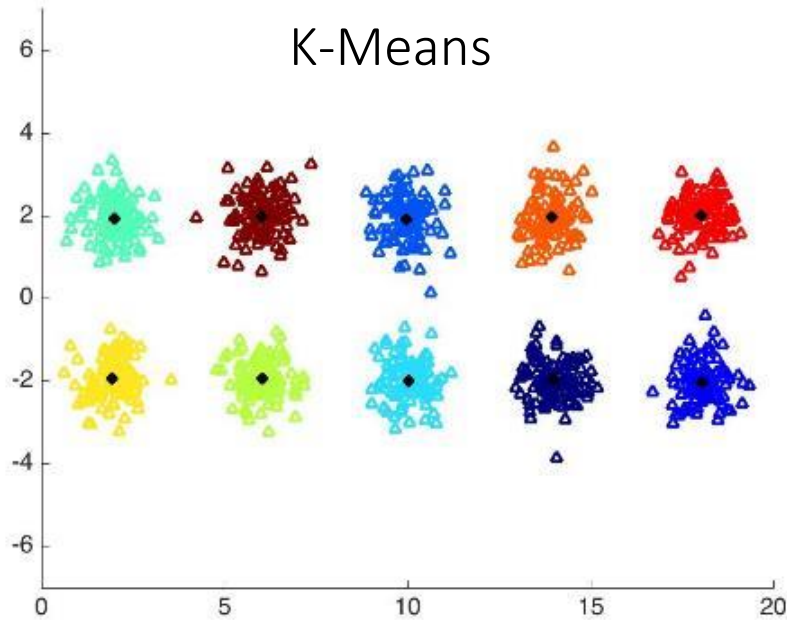
La idea es que para los puntos de un grupo, sus  $k$  vecinos más cercanos están aproximadamente a la misma distancia.

Los puntos de ruido tienen al  $k$  vecino más cercano a mayor distancia.

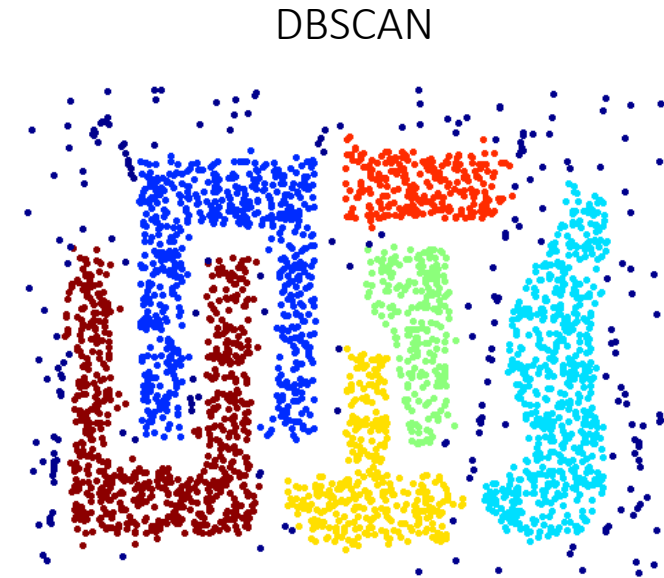
Trace la distancia ordenada de cada punto a su  $k$  vecino más cercano, y seleccione **EPS** cerca del crecimiento exponencial.



# K-medias y DBSCAN



- Eficiente  $O(K*n*i)$
- Encuentra clusters esféricos
- Sensible a condiciones iniciales y  $k$
- Solo se puede usar cuando el “mean” esta definido (variables continuas)
- Susceptible a los outliers
- Existe variaciones (k-modes, k-medoids)
- No es bueno con clusters de diferente tamaño / densidad

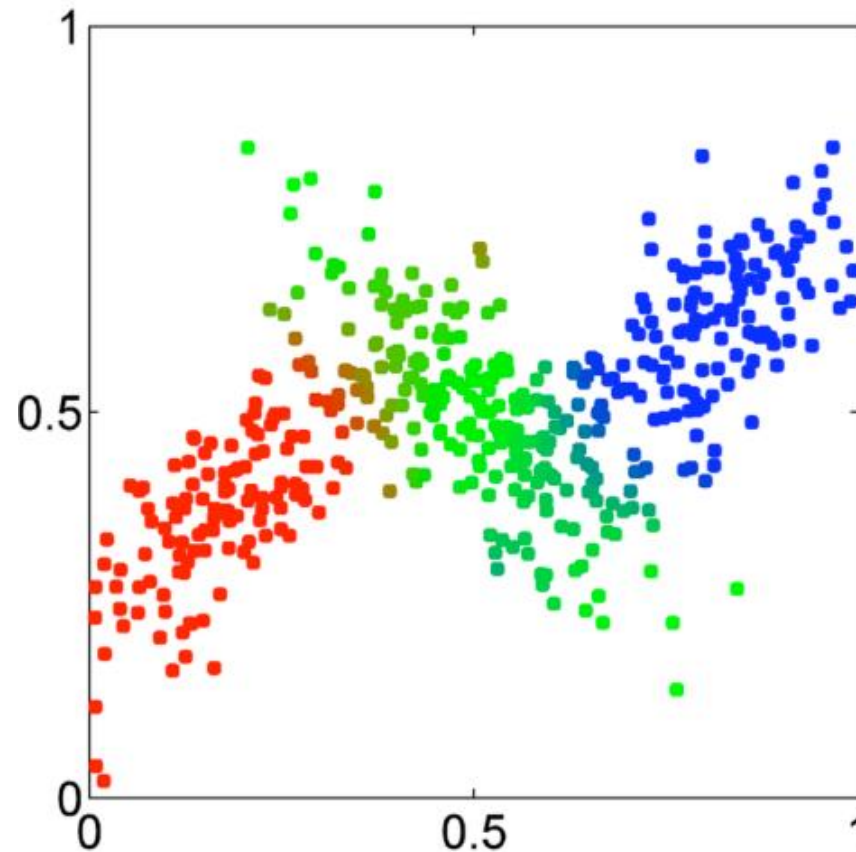


- Menos eficiente  $O(n^2)$
- Resistente al ruido
- Puede buscar clusters de distinto tamaño y forma
- No es bueno en espacios con mucha variación de densidad entre puntos
- No es tan bueno con altas dimensionalidades

# **Soft clustering**

# Definición

En los algoritmos de clústeres difusos, cada punto de datos puede tener una pertenencia o probabilidad de asignación distinto de cero a muchos (normalmente todos) clústeres.

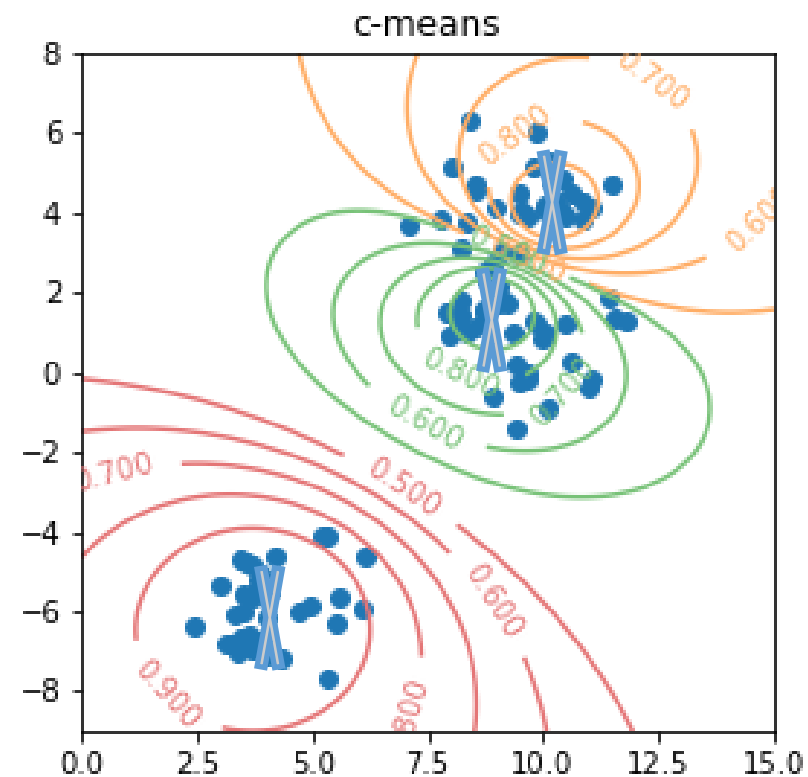
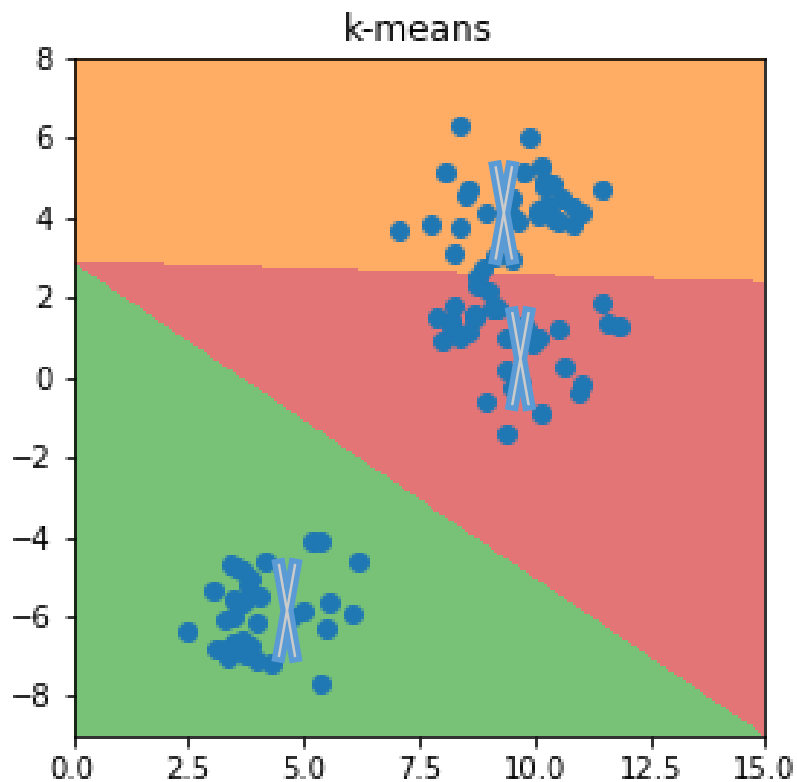


# Fuzzy C-Means

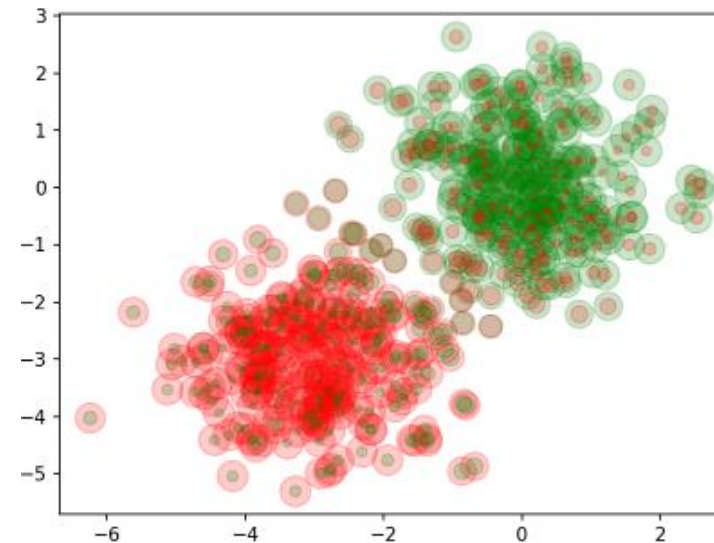
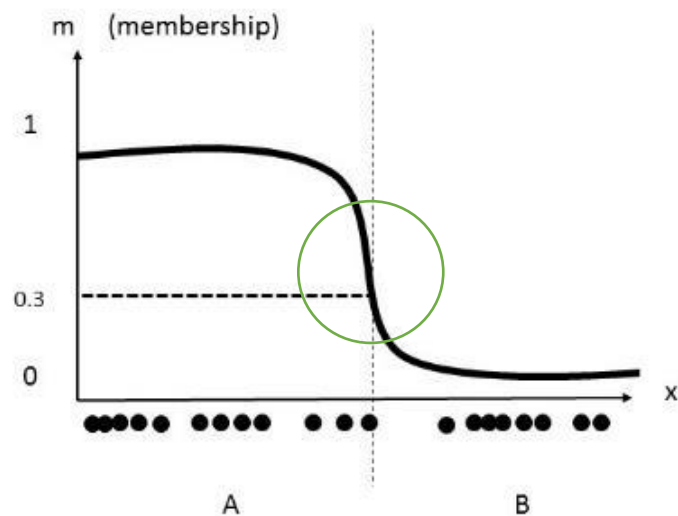
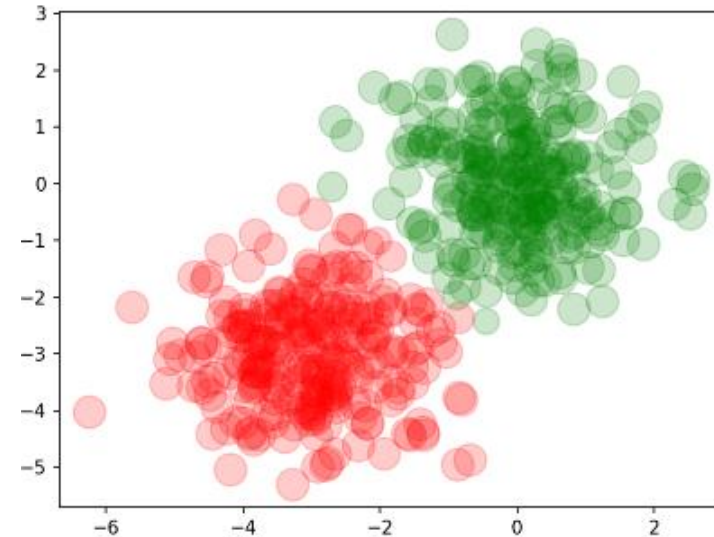
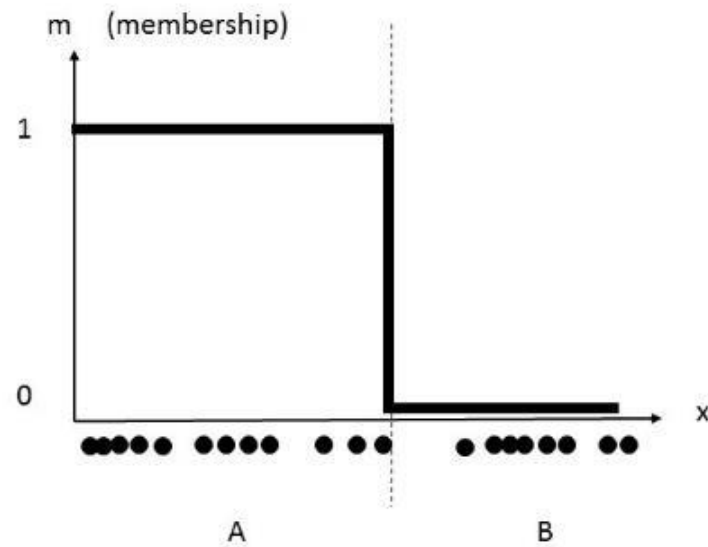


# Fuzzy C-means

Permite que cada punto de datos pertenezca a más de un clúster. Dado el número de clúster  $K$  (determinado por el usuario), cada clúster está asociado a un centroide y cada punto se asigna a uno o varios clústeres en función de una medida de similitud.



# Pertenencia difusa



# Algoritmo

- Paso 1: Inicialización aleatoria de U de tal manera que:

Número de clústeres, definido por el usuario

$$\sum_{j=1}^K u_{ij} = 1, \quad \forall i = 1, \dots, n$$

Número de entidades

Pertenencia de punto i a cluster j

- Paso 2: Seleccione K centroides basado en:

$$r_j = \frac{\sum_{i=1}^n u_{ij}^m x(i)}{\sum_{i=1}^n u_{ij}^m}$$

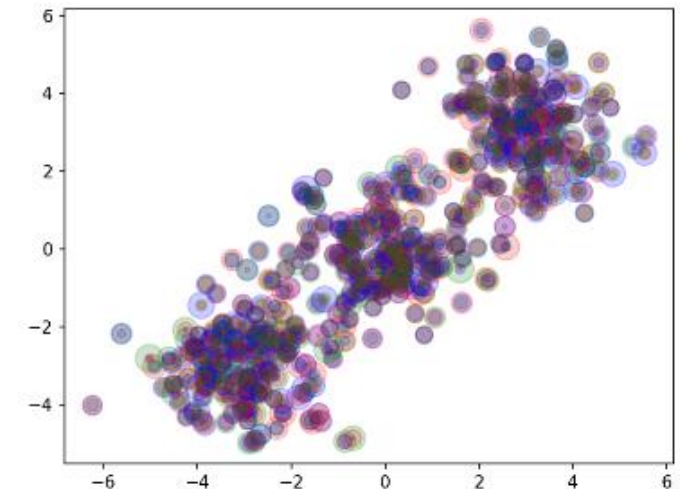
Centro de cluster j

m: Fuzzifier [1, ∞)

X(i): Punto de datos i

Matriz de pertenencia

K clusters X points	K j=1	K j=2	K j=3
X i=1	0,5	0,2	0,3
X i=2	1	0	0
X i=3	0,3	0,3	0,4
X i=4	0,1	0,1	0,8
X i=5	0	0,5	0,5



- Paso 3: Calcule la función de costo  $wc(C)$ .

Si la función de costo obtiene un valor más alto o está por debajo de un umbral específico, DETENGA EL ALGORITMO.

$$wc(C) = \sum_{k=1}^K wc(C_k) = \sum_{k=1}^K \sum_{x(i) \in C_k} u_{ik}^m d(x(i), r_k)$$

- Paso 4: Vuelva a calcular la matriz de pertenencia utilizando:

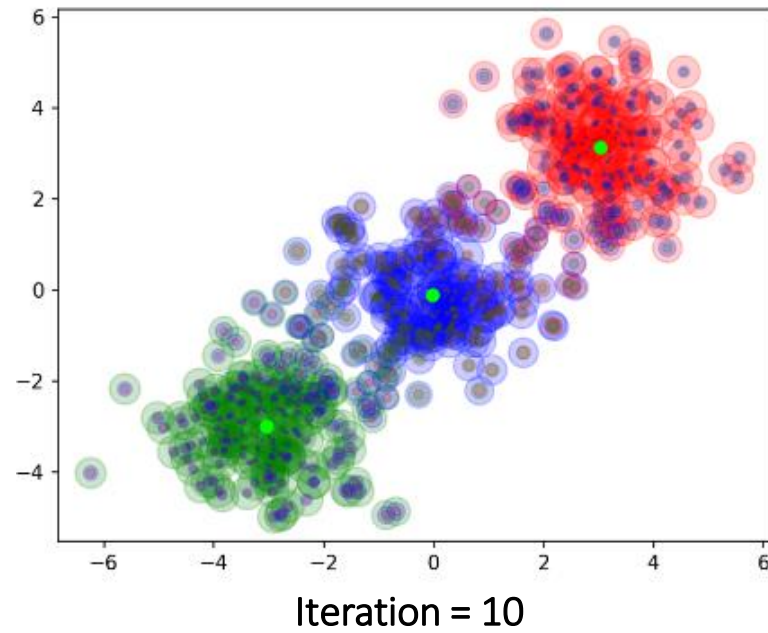
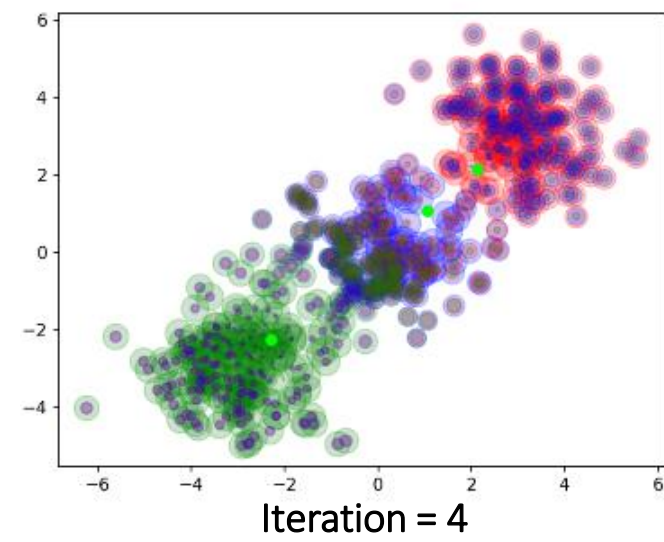
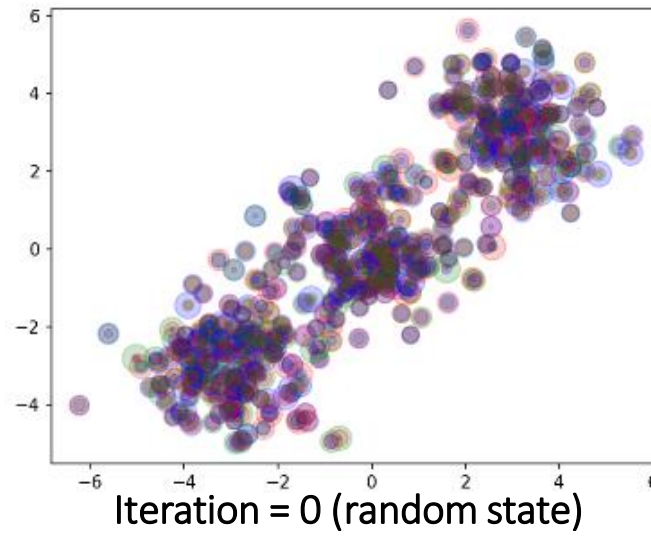
$$u_{ij} = \frac{\frac{1}{d(r_j, x(i))^{\frac{1}{m-1}}}}{\sum_{k=1}^K \left( \frac{1}{d(r_k, x(i))^{\frac{1}{m-1}}} \right)}$$

**Matriz de pertenencia**

K clusters X points	K j=1	K j=2	K j=3
X i=1	0,9	0,2	0,9
X i=2	0,3	0,3	0,4
X i=3	0,3	0,3	0,8
X i=4	0,1	0,1	0,8
X i=5	0	0,5	0,5

- Paso 5: Volver al Paso 2

# FCM: $K=3$ $M=2$



# Ajuste de parámetros

K y M se definen por el usuario

Función objetivo: generalización de la distancia dentro del clúster (cohesión)

$$wc(C) = \sum_{k=1}^K wc(C_k) = \sum_{k=1}^K \sum_{x(i) \in C_k} u_{ik}^m d(x(i), r_k)$$

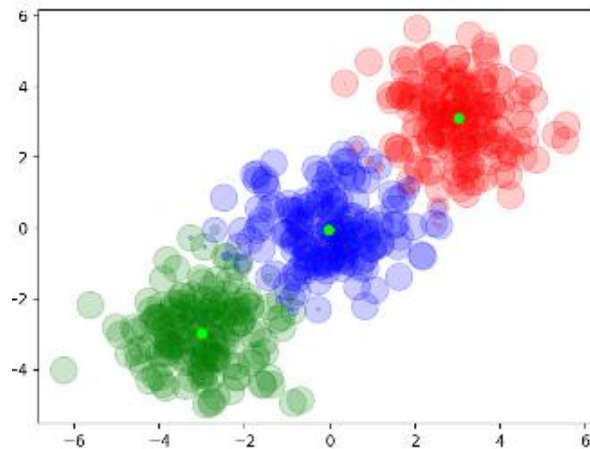
- U es una matriz de pertenencia / peso con valores entre 0 y 1, de tal forma que

$$\sum_{j=1}^K u_{ij} = 1, \quad \forall i = 1, \dots, n$$

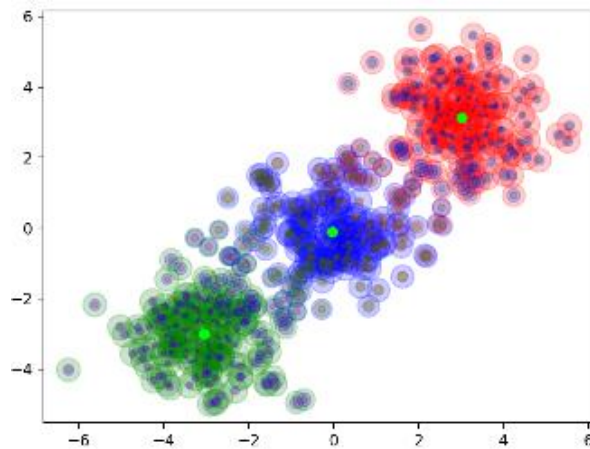
- $m \in [1, \infty)$  es el fuzzificador que determina el nivel de difusividad del cluster.
- Un m grande da como resultado valores de pertenencia  $u_{ij}$  más pequeños, y por lo tanto, clústeres mas difusos. En el límite  $m=1$ , todos los  $u_{ij}$  convergen un 0 o 1, lo que implica un particionamiento dur equivalente a k-medias
- El estado inicial es aleatorio

# Fuzzificador M

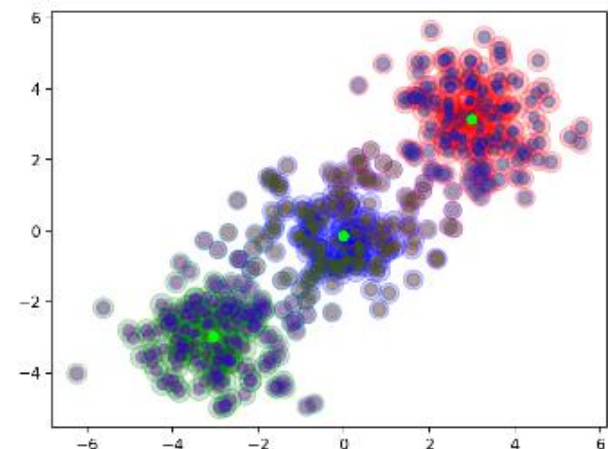
FCM:  $K=3$



$M=1.1$



$M=2$



$M=3$

# Coeficiente de partición difusa (FPC)

El coeficiente de partición difuso evalúa la variabilidad de las asignaciones.

$$F(U) = \frac{\text{tr}(U * U^T)}{n}$$

donde  $U$  es la matriz de pertenencia,  $*$  es el operador de multiplicación entre matrices, y  $\text{tr}()$  es la traza de la matriz, es decir, la suma de los valores diagonales.

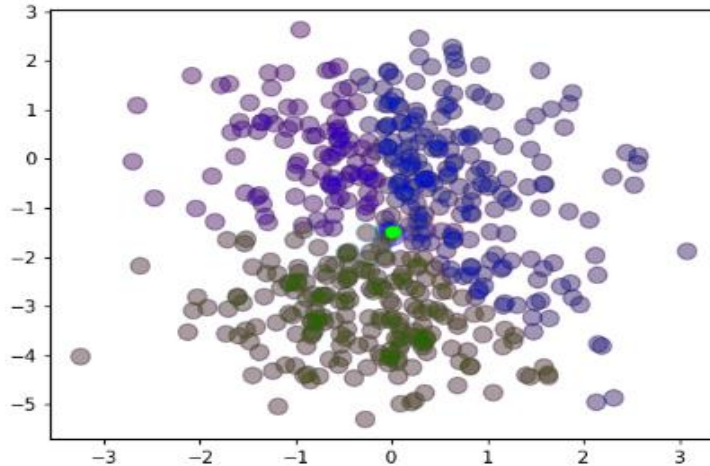
El coeficiente de partición difusa varía entre 0 y 1 donde un valor cercano a 1 implica menor variabilidad en la matriz de pertenencia, que se asocia a una mejor clusterización de los datos.



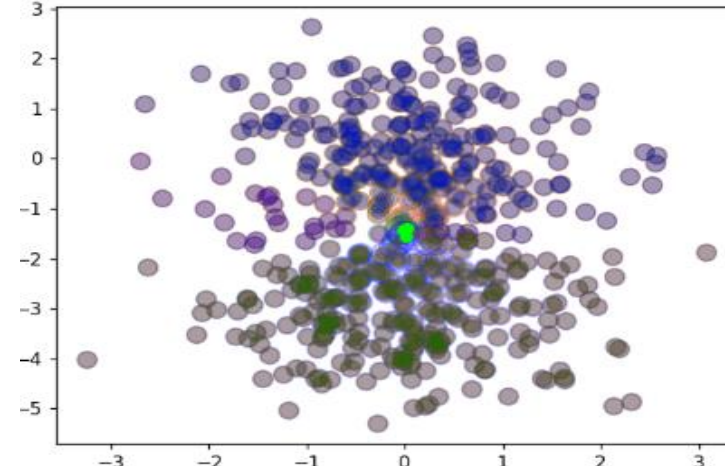
# Coeficiente de partición difusa (FPC)

Puede utilizar FPC para buscar los mejores valores para los parámetros, como en k-medias

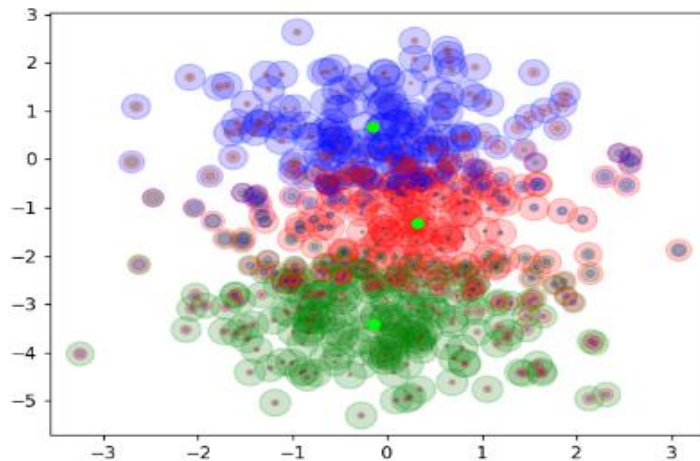
FPC=0.11166



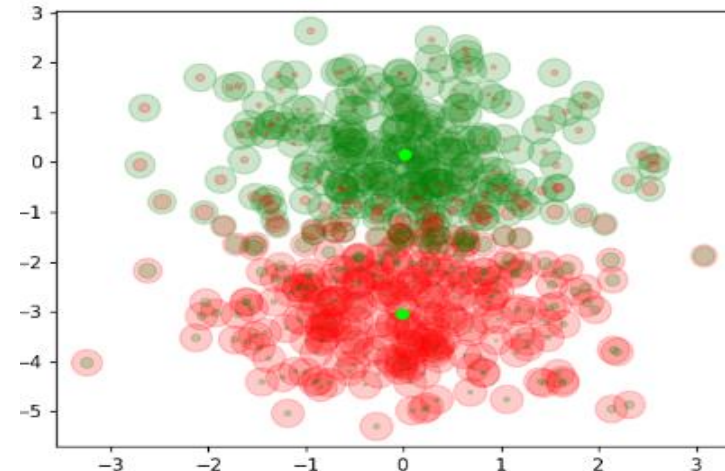
FPC=0.117



FPC=0.27



FPC=0.45

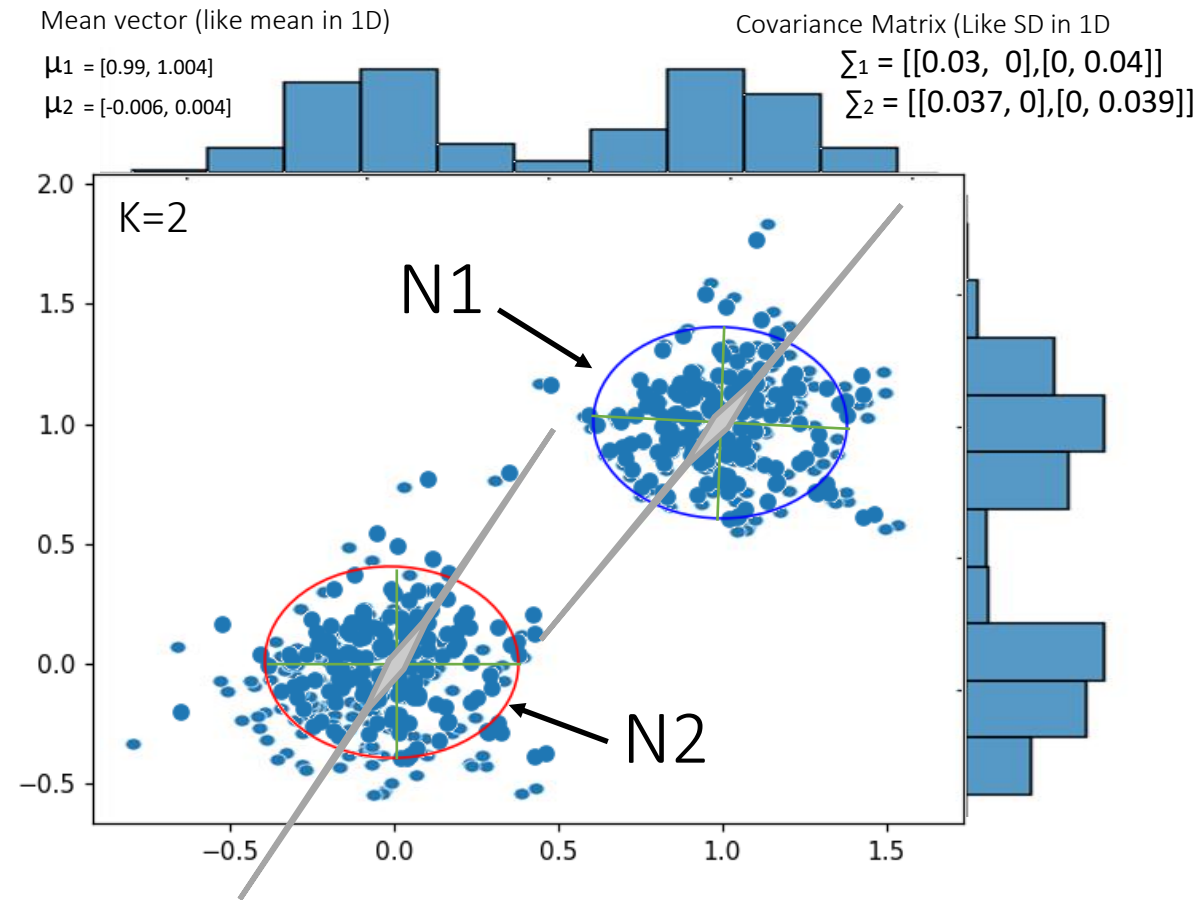


**GMM**

# Definición

Un modelo de mezcla gaussiana (GMM) asume que los datos se generaron a partir de una mezcla de  $K$  gaussianos multidimensionales, donde cada componente tiene parámetros:  $N_k(\mu_k, \Sigma_k)$ .

$K$  es definido por el usuario.



# Gaussiano multivariado

$$X \sim \mathcal{N}(\mu, \sigma^2). \quad \mu = \frac{\sum x}{N}$$
$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2} \quad \sigma = \sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2}$$

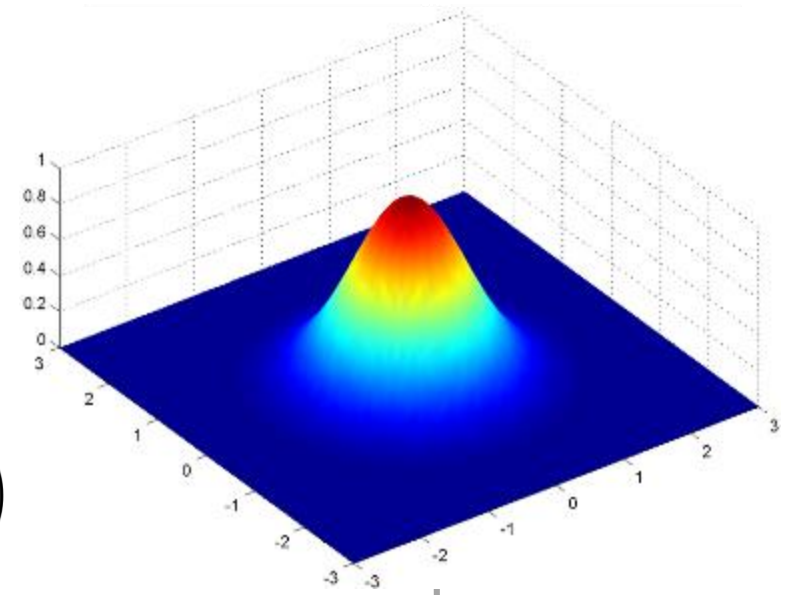
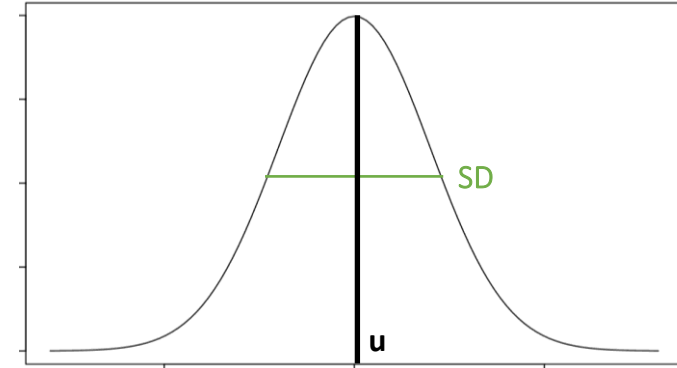
Un Gaussiano multidimensional, para datos con dimensiones  $p$  se especifica de la siguiente manera:

Dónde:

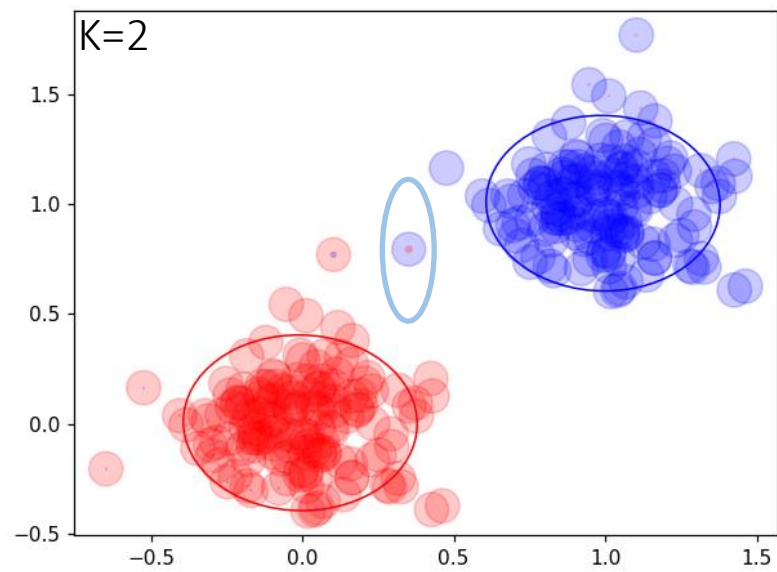
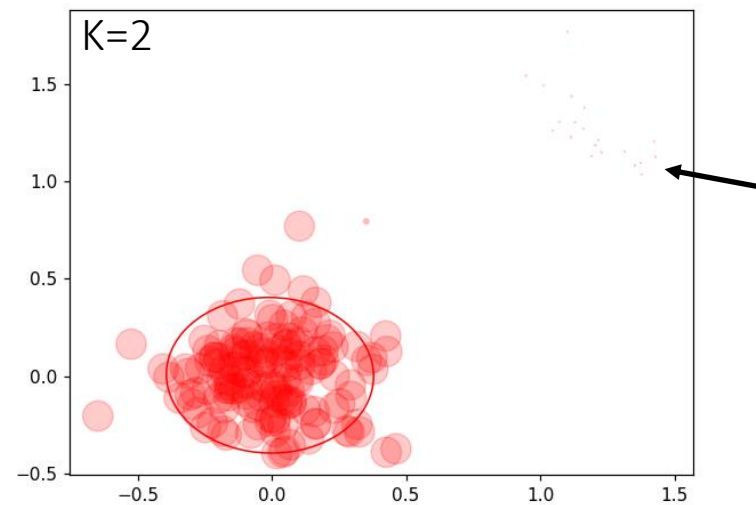
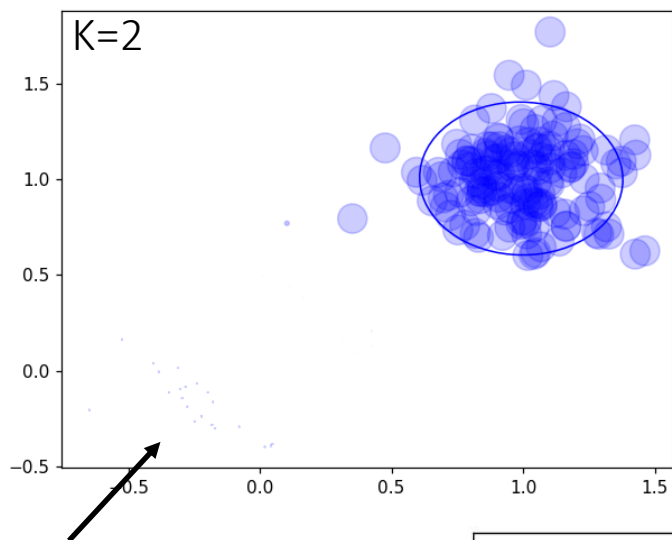
$$x \sim N(\mu, \Sigma)$$

$$\mu = (E[X_1], \dots, E[X_p])$$
$$\Sigma = \begin{bmatrix} \text{Var}(X_1) & \dots & \text{Cov}(X_1, X_p) \\ \dots & \dots & \dots \\ \text{Cov}(X_1, X_p) & \dots & \text{Var}(X_p) \end{bmatrix}$$

$$p(\mathbf{x}) = p(x_1, \dots, x_p) = \frac{1}{\sqrt{(2\pi)^p |\Sigma|}} \exp\left(-\frac{1}{2}(\mathbf{x} - \mu)^T \Sigma^{-1}(\mathbf{x} - \mu)\right)$$



# GMM



# GMM

Prob over total (simple graph)

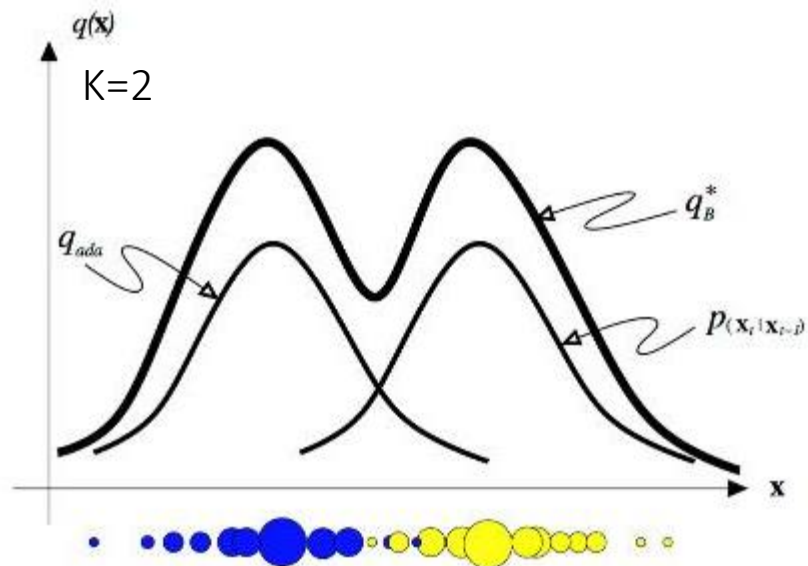
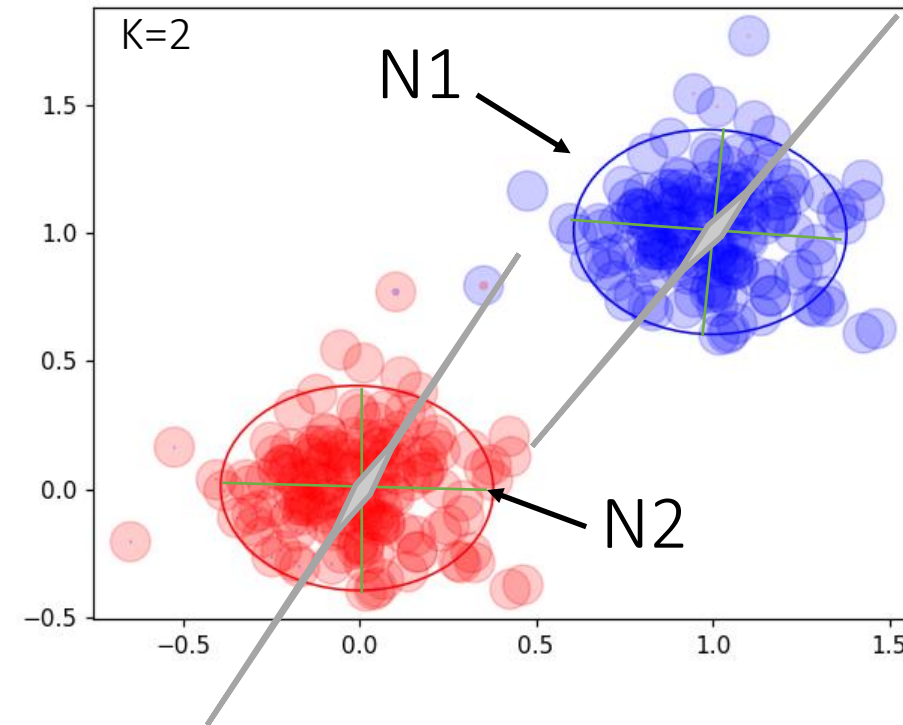
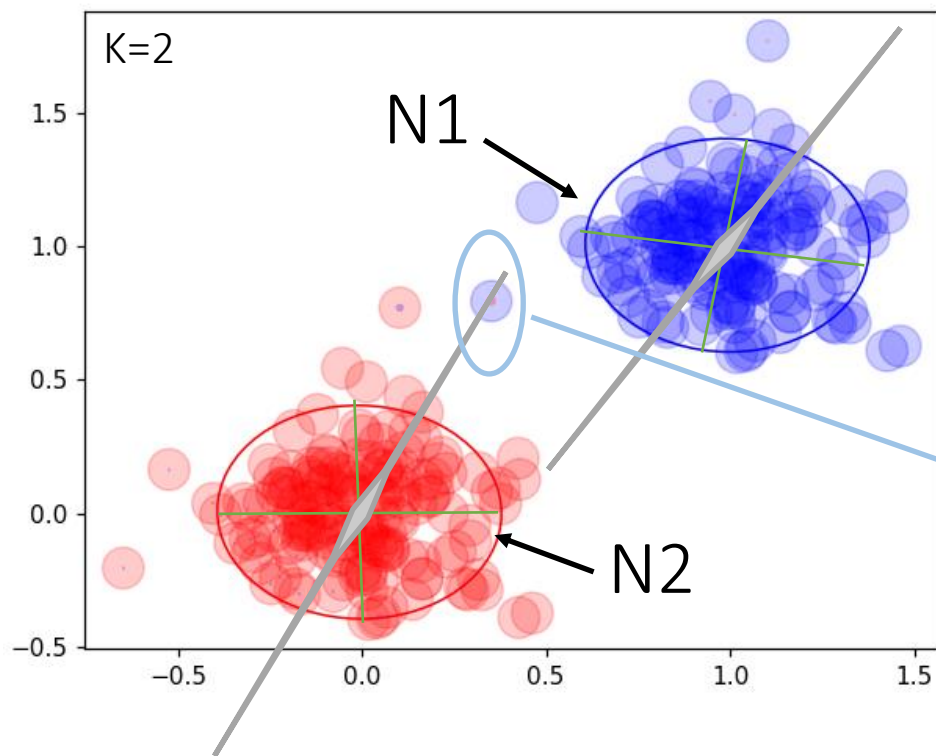


Fig.5. Mixture of Gaussians for the proposal distribution.

Model GMM (K component)



# GMM



$$\mu_1 = [0.99, 1.004]$$

$$\mu_2 = [-0.006, 0.004]$$

$$\pi_1 = 0.49$$

$$\pi_2 = 0.51$$

$$\Sigma_1 = [[0.03, 0], [0, 0.04]]$$

$$\Sigma_2 = [[0.037, 0], [0, 0.039]]$$

Cada punto tiene una probabilidad de pertenecer a N1 y N2

$$p(x) = p(N1)p(x|N1) + p(N2)p(x|N2)$$

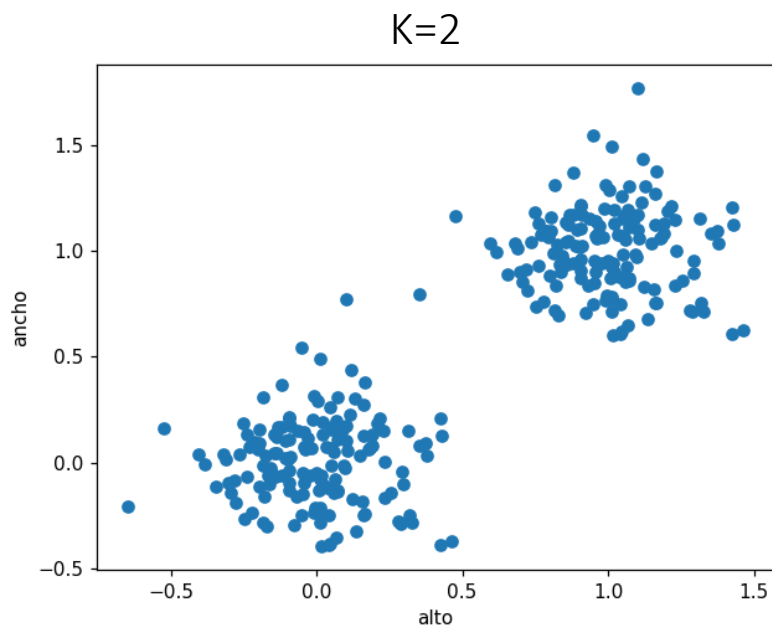
$$p(N1)p(x|N1) \ll p(N2)p(x|N2)$$

$$p(x) = \sum_{k=1}^K \underbrace{p(k)}_{\text{P de k sobre el total de K}} \underbrace{p(x|k)}_{\text{P de x dado k}} = \sum_{k=1}^K p(k) p(x|x \sim N(\mu_k, \Sigma_k)) = \sum_{k=1}^K \underbrace{\pi_k}_{\text{Peso de k}} \underbrace{N(\mu_k, \Sigma_k)}_{\text{Gauss Dist}}$$

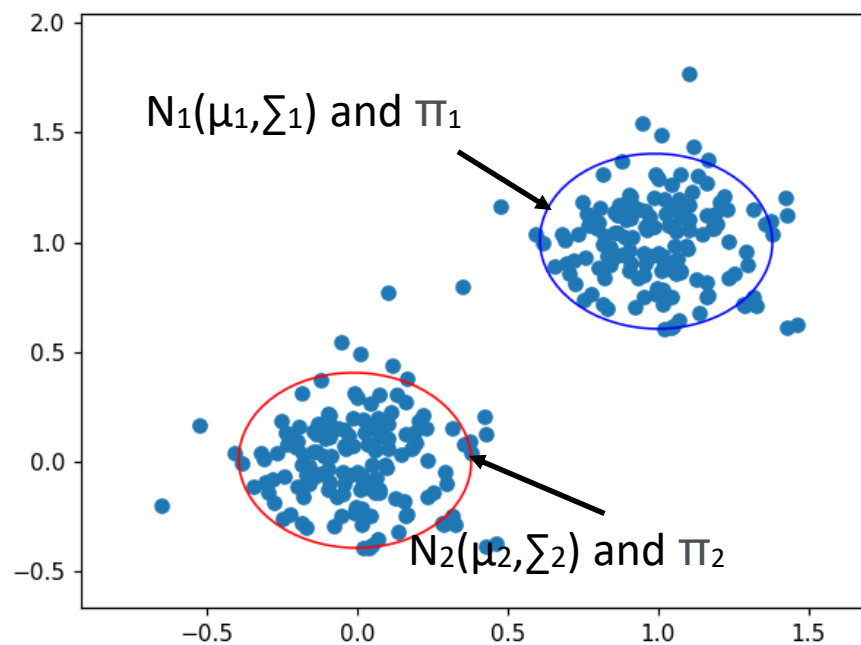
$\mu_k$ : Media de k  
 $\Sigma_k$ : Matriz de covarianza de k

# Algoritmo

Dado un  $K$ , ¿cómo encontramos el mejor valor para cada  $\mu_k$ ,  $\Sigma_k$  y  $\pi_k$ ?



GMM



$\text{MAX}$   
↓

$$p(x) = \sum_{k=1}^K p(k)p(x|k) = \sum_{k=1}^K p(k)p(x|x \sim N(\mu_k, \Sigma_k)) = \sum_{k=1}^K \pi_k N(\mu_k, \Sigma_k)$$



# Algoritmo

- Maximiza la verosimilitud de todo el conjunto de datos:

**MAX LIKELIHOOD**

$$P(X|\pi, \mu, \Sigma) = \prod_{i=1}^N \left[ \sum_{k=1}^K \pi_k N(x_i | \mu_k, \Sigma_k) \right]$$

Diagram annotations:

- $P(X|\pi, \mu, \Sigma)$ : DataSet X with N points
- $\prod_{i=1}^N$ : Each point  $i$
- $\sum_{k=1}^K$ : Each cluster  $k$
- $\pi_k$ : Weight of  $k$
- $N(x_i | \mu_k, \Sigma_k)$ : P of  $i$  in Gauss Dist  $k$

- La suma de todos los pesos tiene que ser 1:

$$\sum_{k=1}^K \pi_k = 1$$

Tenemos que ajustar  $\mu_k$ ,  $\Sigma_k$  y  $\pi_k$  para maximizar la probabilidad = Derivar verosimilitud con respecto a  $\mu_k$ ,  $\Sigma_k$  y  $\pi_k$

# Algoritmo

- Defina algunas cantidades auxiliares:

$$Z_{nk} = \begin{cases} 1 & \text{if } x_n \text{ is in cluster } k \\ 0 & \text{otherwise.} \end{cases}$$

$$\gamma(Z_{nk}) = P(Z_{nk} = 1 | X_n) = \frac{P(Z_k=1)P(X_n|Z_{nk}=1)}{\sum_{j=1}^K P(Z_{nj}=1)P(X_n|Z_{nj}=1)} = \frac{\pi_k N(x_n | \mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j N(x_n | \mu_j, \Sigma_j)}$$

(Bayes thm.)

$$N_k = \sum_{n=1}^N \gamma(Z_{nk})$$

La derivada de:  $P(X|\pi, \mu, \Sigma) = \prod_{i=1}^N \left[ \sum_{k=1}^K \pi_k N(x_i | \mu_k, \Sigma_k) \right]$  respecto con  $\mu_k$ ,  $\Sigma_k$  y  $\pi_k$

$$\mu_k = \frac{1}{N_k} \sum_{n=1}^N \gamma(Z_{nk}) x_n$$

$$\Sigma_k = \frac{1}{N_k} \sum_{n=1}^N \gamma(Z_{nk}) (x_n - \mu_k)(x_n - \mu_k)^T$$

$$\pi_k = \frac{N_k}{N}$$

Gamma => Mu, Sigma, Pi

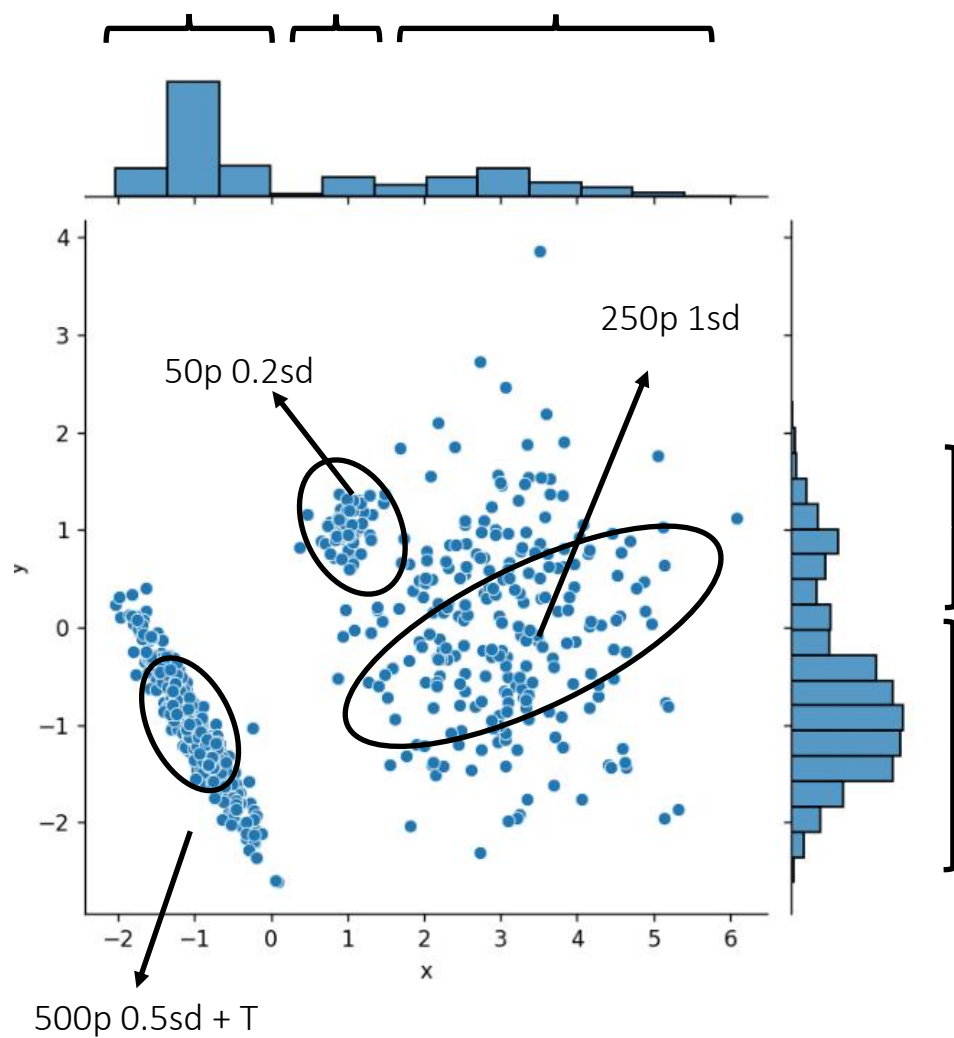
Mu, Sigma, Pi => Gamma

Dependencia circular=> algoritmo

# Resumen algoritmo

- Maximización de verosimilitud (EM)
  - Paso 0) Seleccione K
  - Paso 1) Inicializar  $\mu_k$ ,  $\Sigma_k$  y  $\pi_k$  con valores aleatorios (inicializar Q veces, utilizar una función de puntuación, tomar la mejor condición inicial)
  - Paso 2) Calcular gamma (paso de expectativa)
  - Paso 3) Recalcular  $\mu_k$ ,  $\Sigma_k$  y  $\pi_k$  con Gamma (Paso de Maximización)
  - Paso 4) Compruebe la convergencia de los parámetros, si la convergencia no se satisface vaya al paso 2

# Ejemplo



Datos generados artificialmente

Pistas:

Ocupe 3 blobs con diferente numero de puntos: 500, 250 y 50

Los 3 blobs tienen dista densidad (SD): 0.5, 1 y 0.2

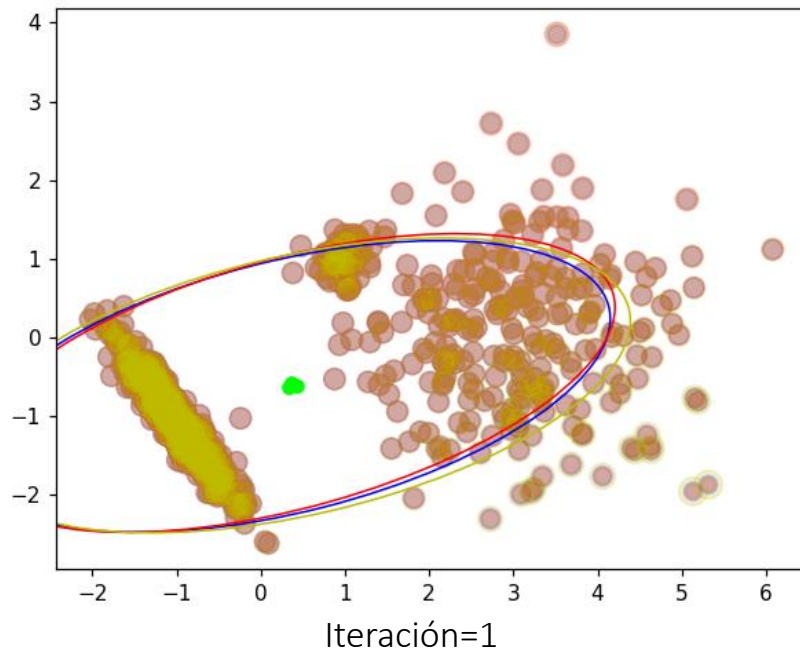
2 de ellos son esféricos, 1 de ellos a sufrido una transformación para alargarlo.

¿Cómo deberían ser los pesos de  $\pi$ ?

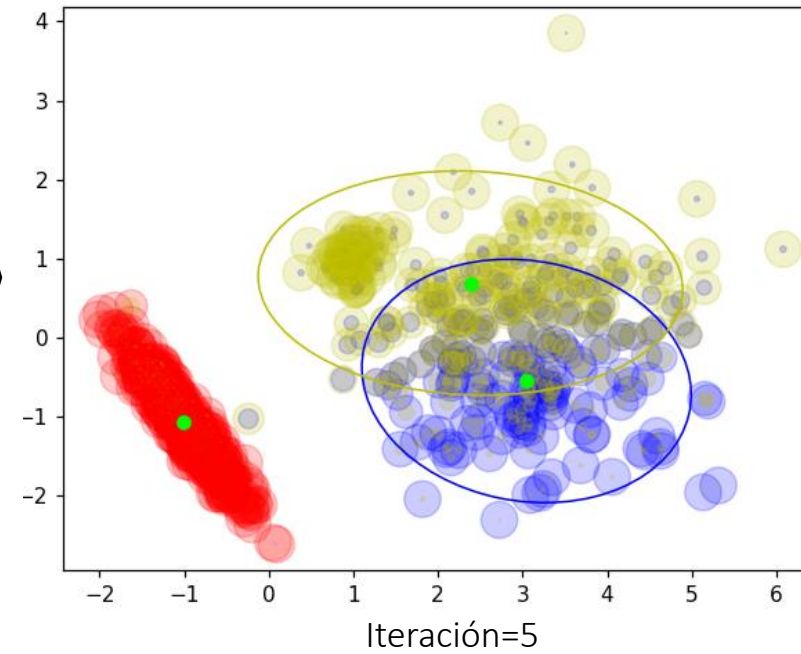
Supongamos que no sabemos nada de lo anterior y usemos GMM

# Ejemplo

K=3, 100 inicios random



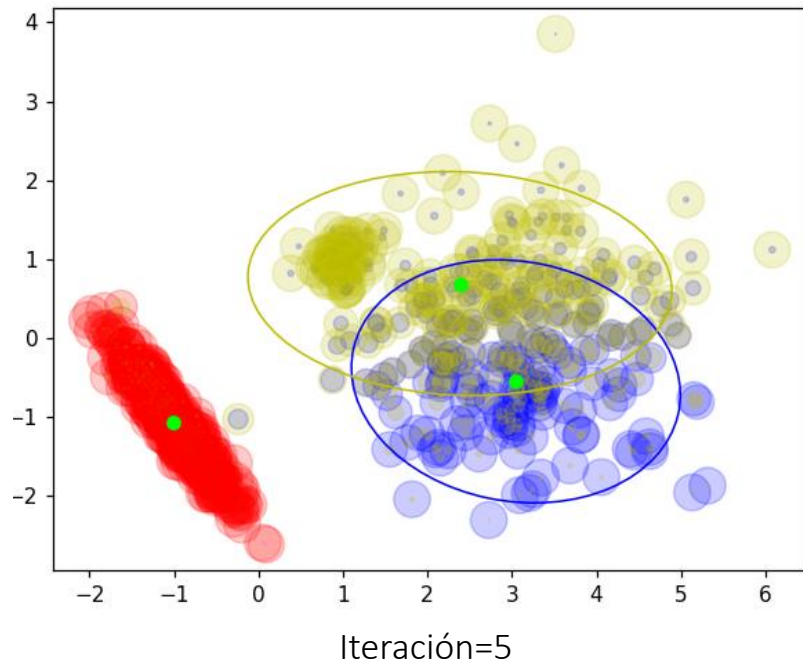
$\pi_1=0.333$   $\pi_2=0.329$   $\pi_3=0.337$



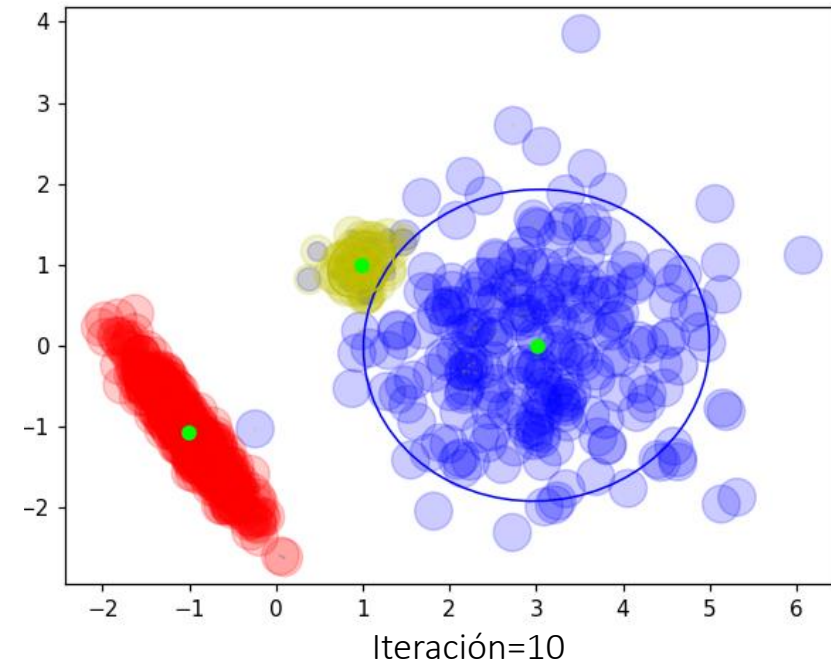
$\pi_1= 0.15$   $\pi_2=0.62$   $\pi_3=0.21$

# Ejemplo

K=3, 100 inicios random



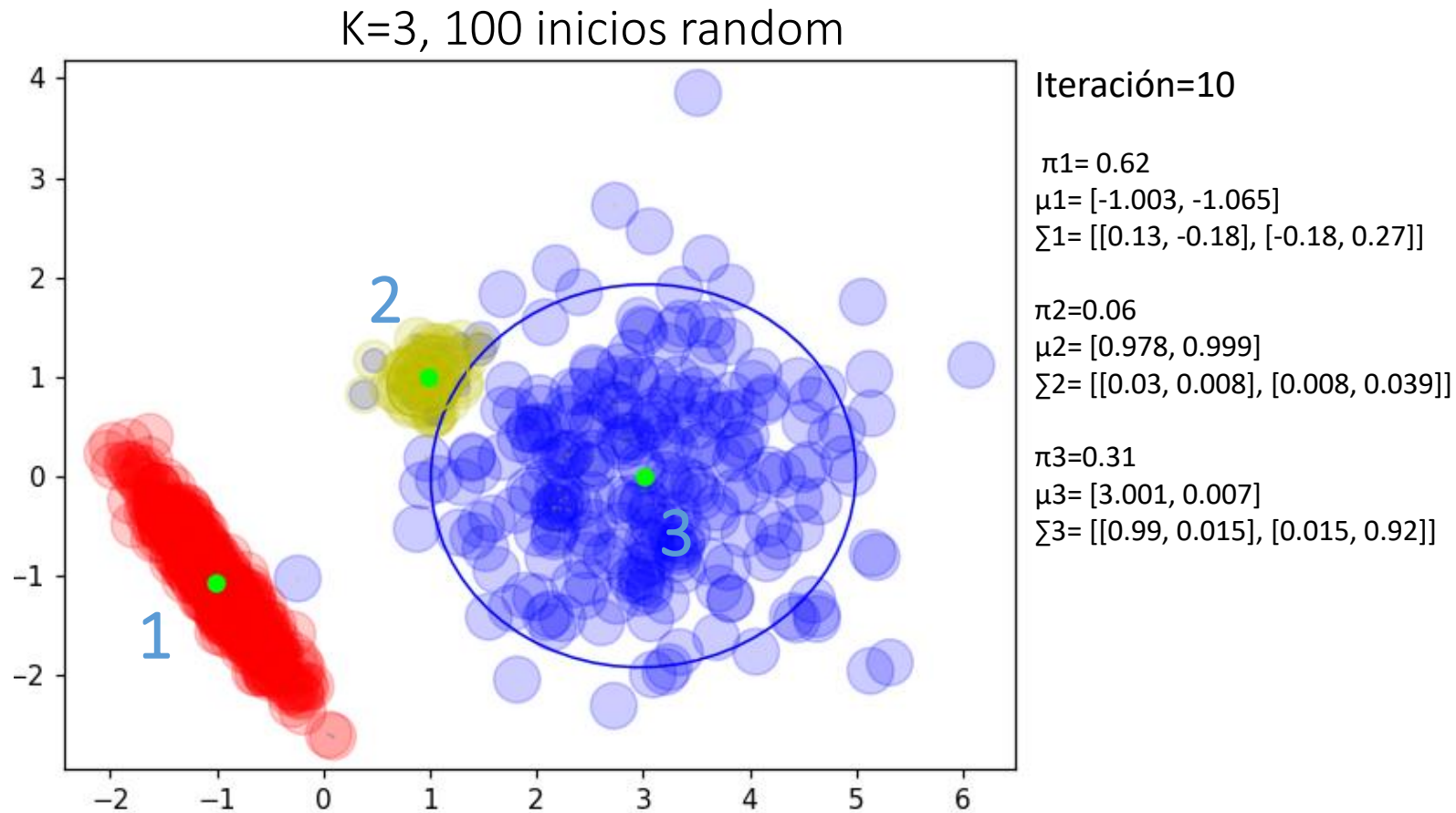
$\pi_1 = 0.15$   $\pi_2 = 0.62$   $\pi_3 = 0.21$



$\pi_1 = 0.31$   $\pi_2 = 0.62$   $\pi_3 = 0.06$

# Ejemplo

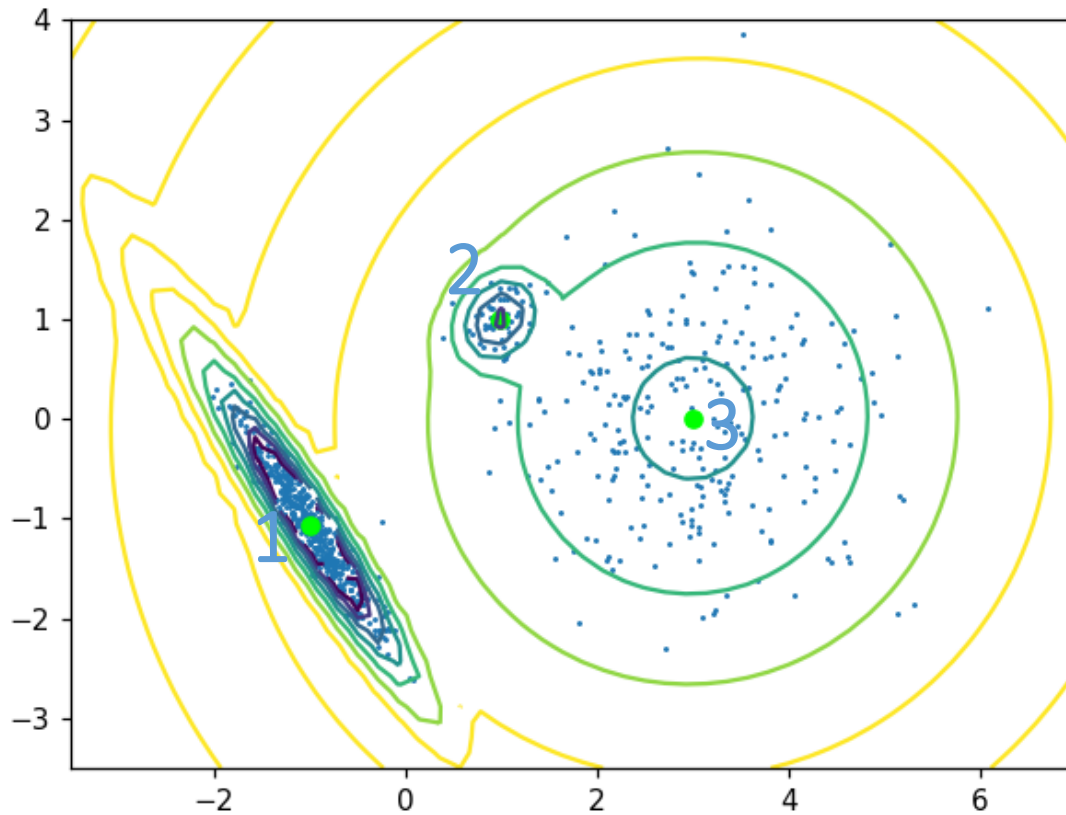
Vista por distribuciones aisladas



# Ejemplo

Vista de la probabilidad total de cada punto

K=3, 100 inicios random



Iteración=10

$\pi_1 = 0.62$

$\mu_1 = [-1.003, -1.065]$

$\Sigma_1 = [[0.13, -0.18], [-0.18, 0.27]]$

$\pi_2 = 0.06$

$\mu_2 = [0.978, 0.999]$

$\Sigma_2 = [[0.03, 0.008], [0.008, 0.039]]$

$\pi_3 = 0.31$

$\mu_3 = [3.001, 0.007]$

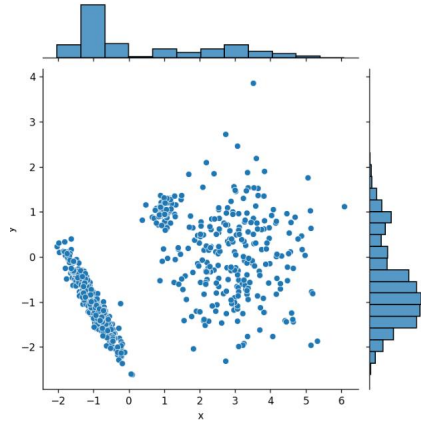
$\Sigma_3 = [[0.99, 0.015], [0.015, 0.92]]$



# Selección de K

- Si elegimos K para maximizar la probabilidad, cuando K aumenta el valor de la probabilidad máxima no puede disminuir.
- Por lo tanto, los modelos más complejos siempre mejorarán la probabilidad.
- Es necesario penalizar la complejidad del modelo.
- Necesitamos un equilibrio entre lo bien que encaja el modelo y los datos y la simplicidad del modelo
  - $\text{Score}(\theta, M) = \text{error}(M) + \text{penalización}(M)$   
La penalización puede depender del número de parámetros del modelo (p) y del número de puntos de datos (n).
  - El error se basa generalmente en la probabilidad de los datos dados el modelo (L).
  - AIC Akaike information criterion  $\text{Score}_{\text{AIC}} = -2 \log L + 2p$
  - BIC: Bayesian information criterion  $\text{Score}_{\text{BIC}} = -2 \log L + p \log n$

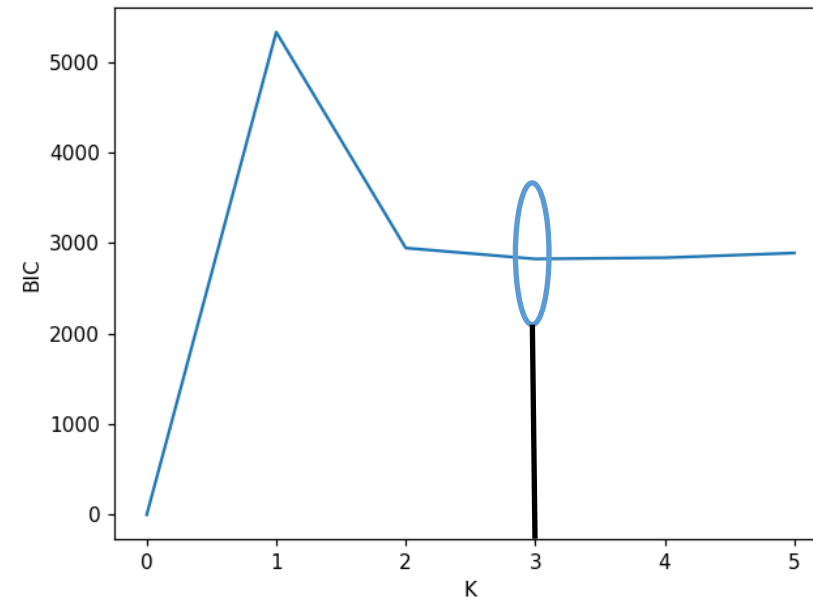
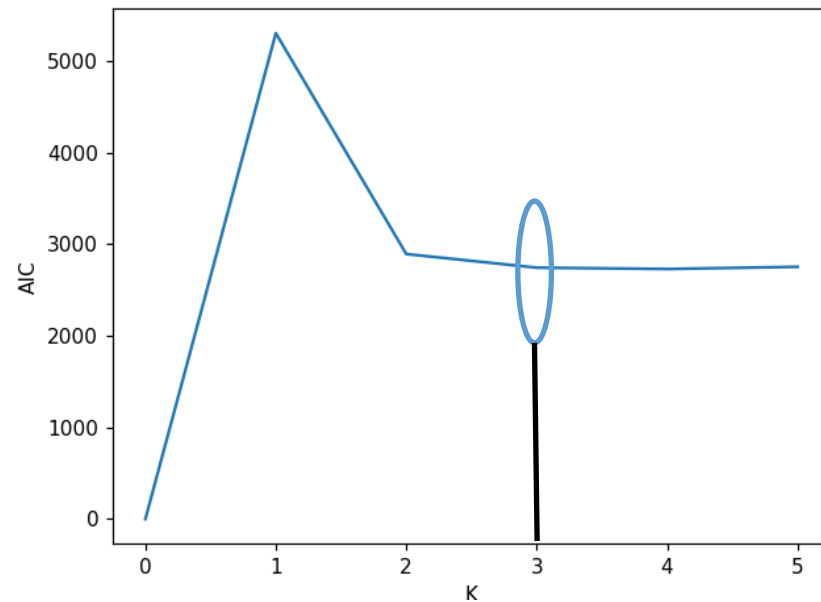
# Selección de K



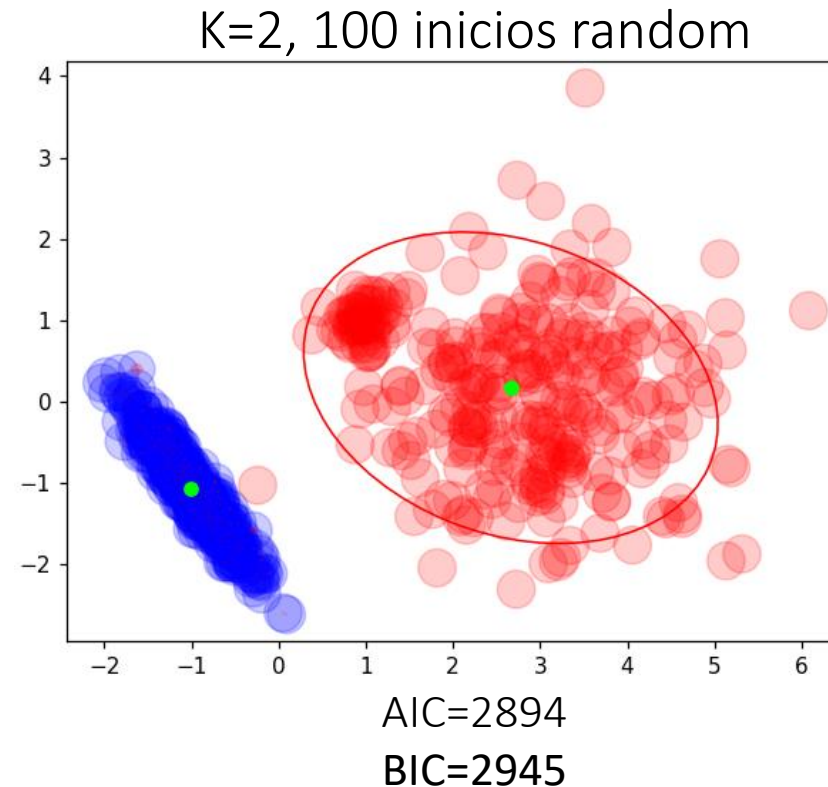
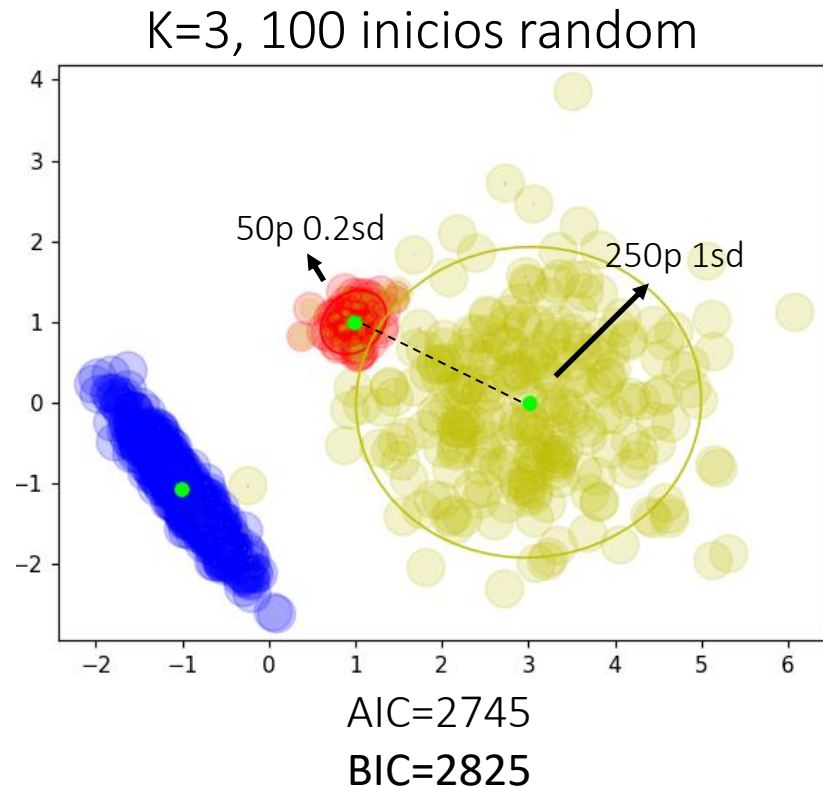
AIC and BIC

100 inicios aleatorios

Hay poca mejora entre  $K=2$  y  $K=3$  debido a la distribución de los datos.



# Ejemplo



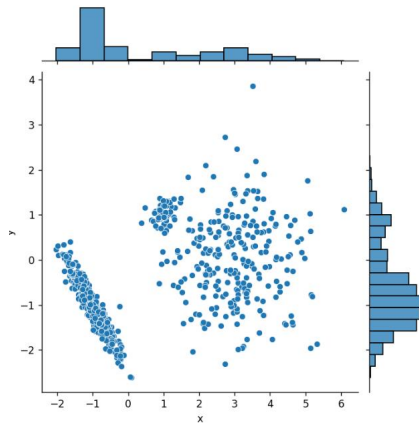
Un cluster puede asimilar a otro muy cercano y que posea pocos datos.  
Hay que tener cuidado al analizar los scores. Como también generar suficientes condiciones iniciales y usar suficientes iteraciones.

# Condiciones iniciales

Opciones para definir la condición inicial de  $\mu_k$ ,  $\Sigma_k$  y  $\pi_k$  para un cierto  $K$

1. El usuario define todos o algunos de ello
2. Utilice otro método, como los medios  $K$ , para buscar en el centroide de los clústeres preliminares ( $\mu_k$ )
3.  $Q$  iteraciones aleatorias iniciales. Puntúa cada uno de ellos y toma el mejor (AIC, BIC o verosimilitud).

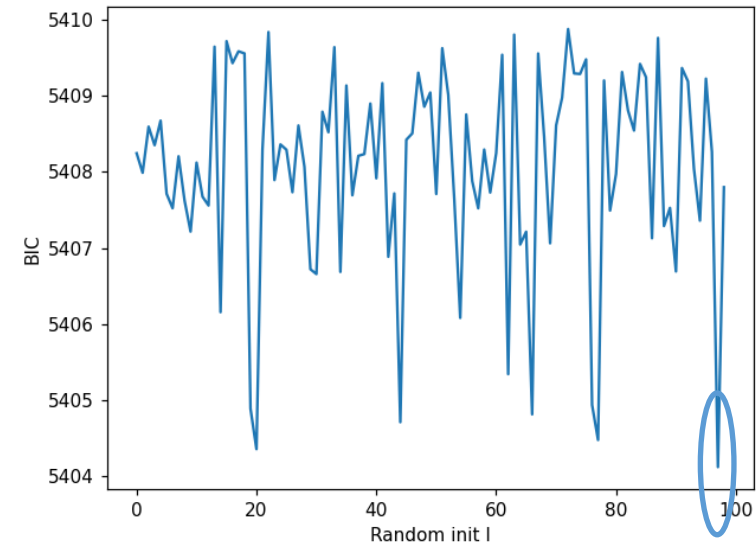
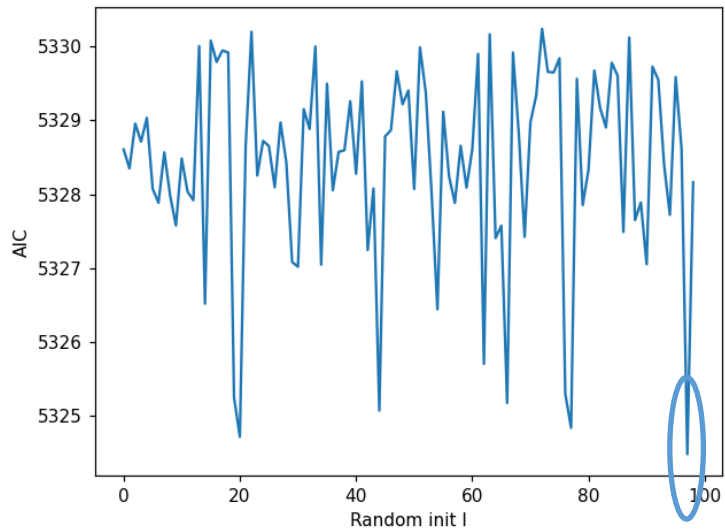
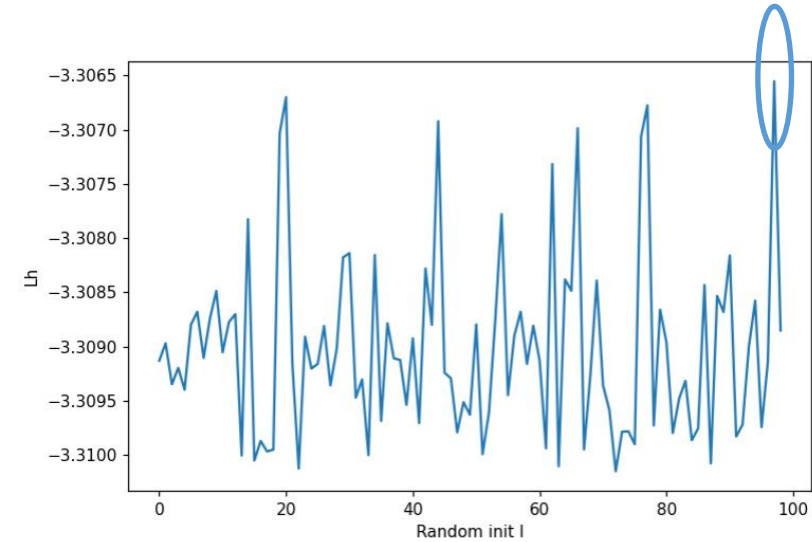
# Condiciones iniciales



AIC, BIC and  
Likelihood Score

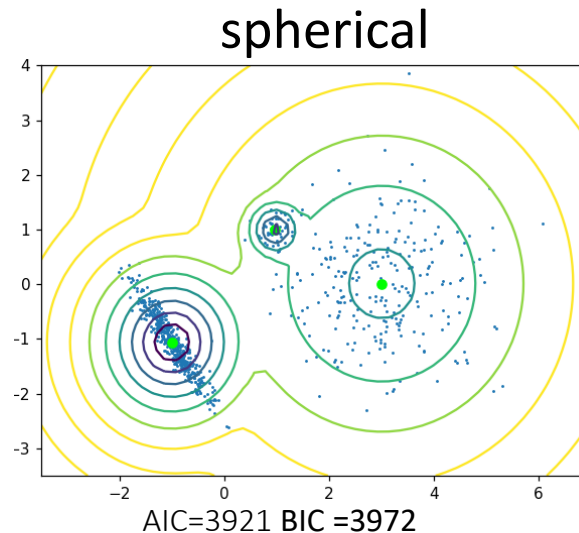
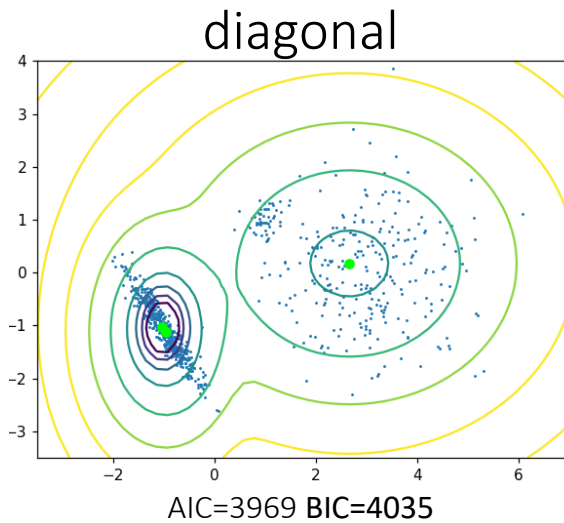
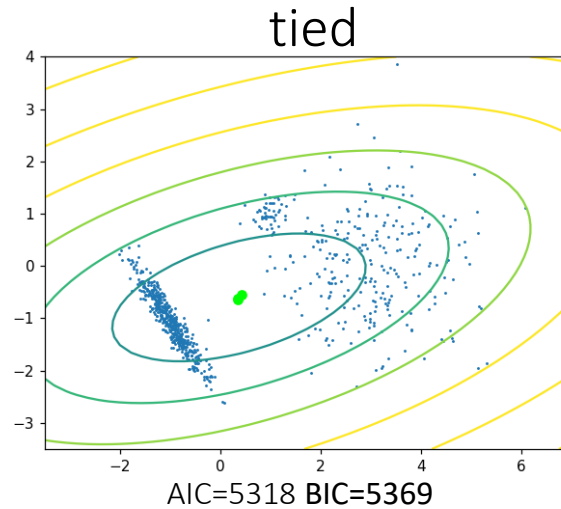
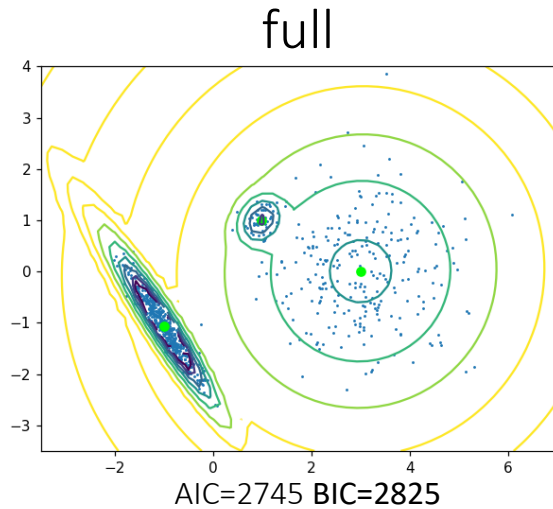
K=3

100 random init



En este caso, la iteración  $i = 97$ ; tiene las mejores condiciones iniciales

# Matrices de covarianzas



- Ell: equal volume, round shape (spherical covariance)
- VII: varying volume, round shape (spherical covariance)
- EEI: equal volume, equal shape, axis parallel orientation (diagonal covariance)
- VEI: varying volume, equal shape, axis parallel orientation (diagonal covariance)
- EVI: equal volume, varying shape, axis parallel orientation (diagonal covariance)
- VVI: varying volume, varying shape, equal orientation (diagonal covariance)
- EEE: equal volume, equal shape, equal orientation (ellipsoidal covariance)
- EEV: equal volume, equal shape, varying orientation (ellipsoidal covariance)
- VEV: varying volume, equal shape, varying orientation (ellipsoidal covariance)
- VVV: varying volume, varying shape, varying orientation (ellipsoidal covariance)

# **Clusters probabilísticos**

Dr. Raimundo Sánchez  
raimundo.sanchez@uai.cl  
@raimun2