

Curs 7:

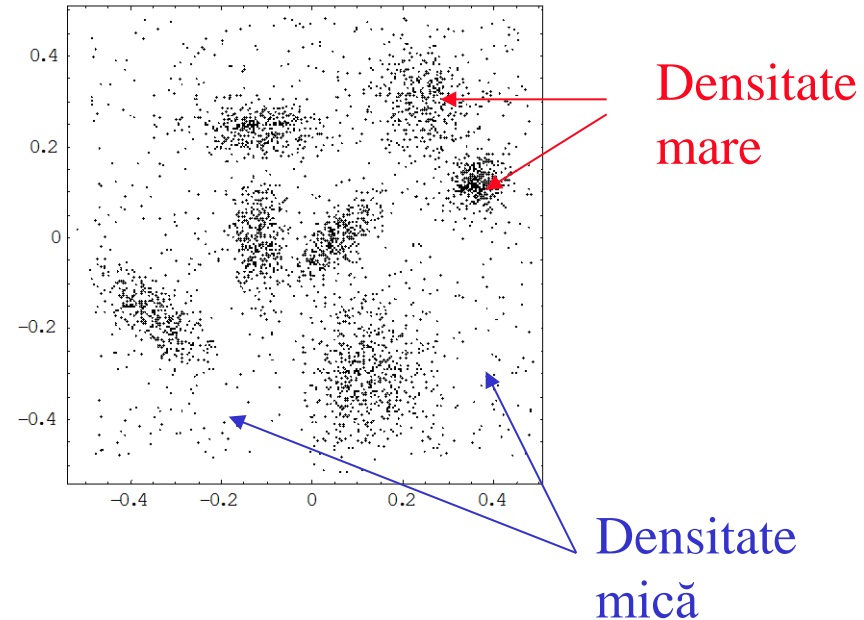
## Gruparea datelor (II)

# Structura

- Metode bazate pe densitate
  - DBSCAN
  - DENCLUE
- Metode probabiliste
  - EM - Expectation Maximization

# Metode bazate pe densitate

- **Cluster** = grup dens de date similare separate de regiuni cu densitate mai mică de date
- **Idee de bază**: estimarea densității locale a datelor
  - se determină numărul de date din vecinătatea punctului analizat (**DBSCAN**)
  - se utilizează funcții de influență pt estimarea densității (**DENCLUE**)
- **Problema principală**:
  - Cum se estimează densitatea?

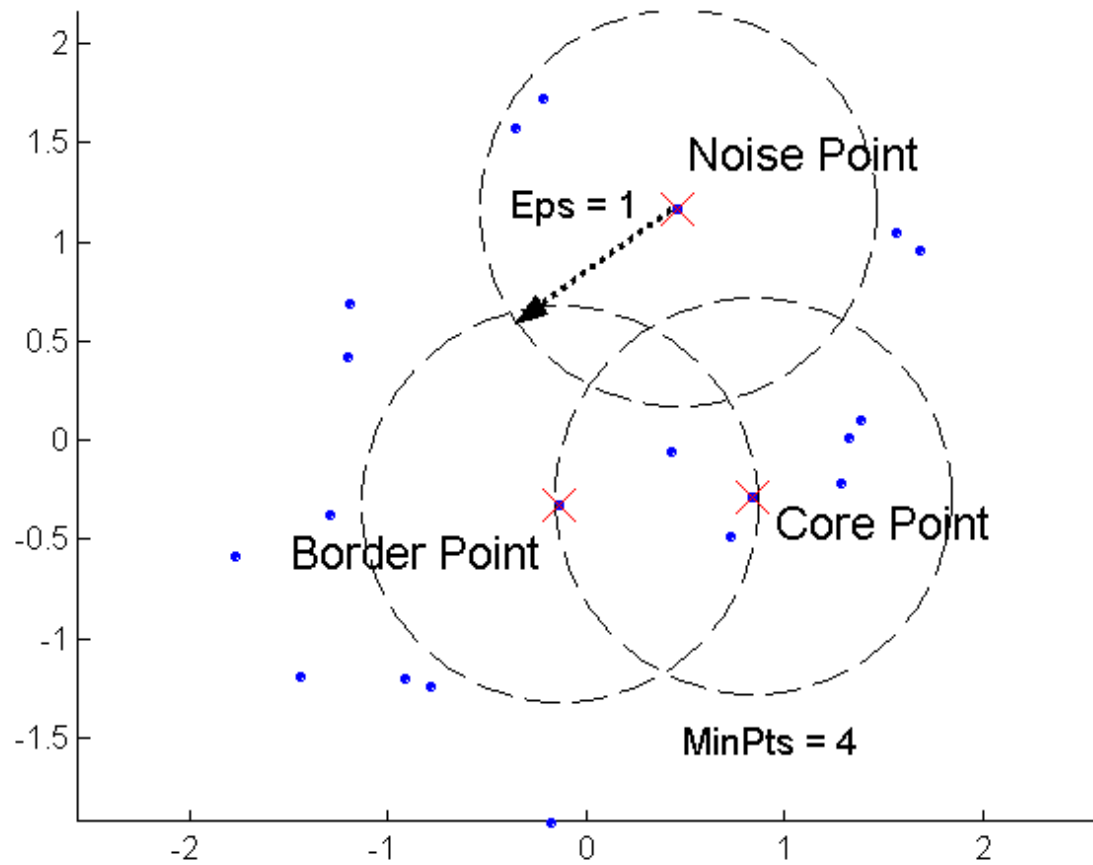


# DBSCAN

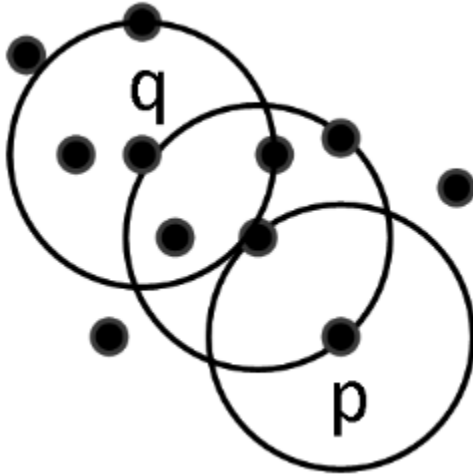
DBSCAN [M.Ester, H Kriegel et al, 1996] este un algoritm de grupare bazat pe următoarele elemente:

- Densitatea estimată într-un punct = numărul de puncte aflate în vecinătatea definită de o anumită rază (**Eps**)
- Un punct este considerat **punct nucleu (core point)** dacă numărul de puncte din vecinătatea sa depășește un prag (**MinPts**); acestea sunt puncte considerate a fi în interiorul clusterului
- Un **punct frontieră (border point)** are un număr de vecini mai mic decât MinPts dar este în vecinătatea unui punct nucleu; Două puncte sunt considerate **conectate** dacă unul este în vecinătatea celuilalt
- Un punct q est accesibil direct (**directly density reachable**) dintr-un punct nucleu p dacă e în vecinătatea lui p; accesibilitatea în sens general este definită ca fiind închiderea tranzitivă a relației de accesibilitate directă (există un lanț de **puncte nucleu** ce începe cu p și se termină cu q cu proprietatea că fiecare e direct accesibil dintr-un punct anterior)
- Un punct de tip zgomot (**noise point**) este orice punct care nu e nici nucleu nici frontieră

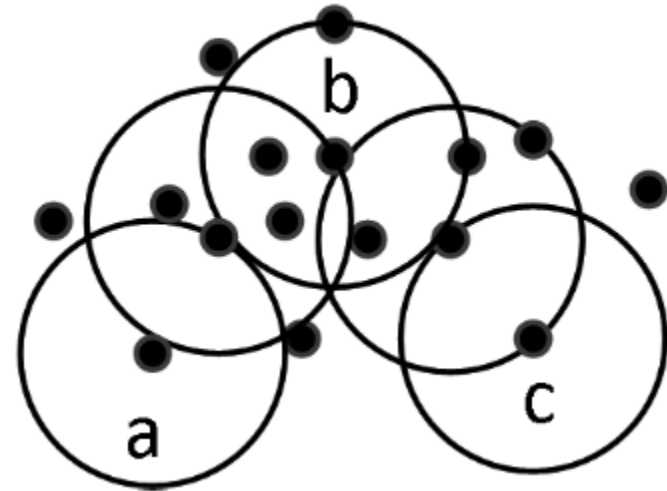
# DBSCAN



# DBSCAN



p este accesibil din q



a este accesibil din b

c este accesibil din b

=> a și c sunt conectate

## Obs:

- Două pcte, a și b, dacă există un punct c astfel încât c este accesibil atât din a cât și din b
- Două puncte conectate ar trebui să aparțină aceluiași cluster => un cluster definit pe baza densității este un set maximal de date conectate

# DBSCAN

```
DBSCAN (SetOfPoints, Eps, MinPts)    // SetOfPoints is UNCLASSIFIED
ClusterId := nextId(NOISE);
FOR i FROM 1 TO SetOfPoints.size DO
    Point := SetOfPoints.get(i);
    IF Point.ClusterId = UNCLASSIFIED THEN
        IF ExpandCluster(SetOfPoints, Point, ClusterId, Eps, MinPts)
            THEN ClusterId := nextId(ClusterId)
        END IF
    END IF
END FOR
END; // DBSCAN
```

[M.Ester, H Kriegel et al, A Density-Based Algorithm for Discovering Clusters in Large Spatial Databases with Noise, 1996]

# DBSCAN

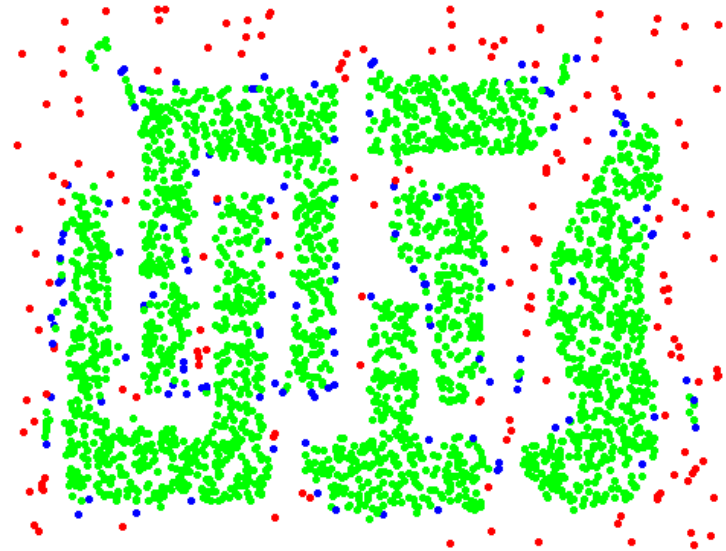
```
ExpandCluster(SetOfPoints, Point, ClId, Eps, MinPts) : Boolean;
  seeds:=SetOfPoints.regionQuery(Point,Eps);
  IF seeds.size<MinPts THEN SetOfPoint.changeClId(Point,NOISE); RETURN False;
  ELSE // all points in seeds are density-reachable from Point
    SetOfPoints.changeClIds(seeds,ClId); seeds.delete(Point);
  WHILE seeds <> Empty DO
    currentP := seeds.first(); result := SetOfPoints.regionQuery(currentP, Eps);
    IF result.size >= MinPts THEN
      FOR i FROM 1 TO result.size DO resultP := result.get(i);
        IF resultP.ClId IN {UNCLASSIFIED, NOISE} THEN
          IF resultP.ClId = UNCLASSIFIED THEN seeds.append(resultP); END IF;
          SetOfPoints.changeClId(resultP,ClId);
        END IF; // UNCLASSIFIED or NOISE
      END FOR; END IF; // result.size >= MinPts
    seeds.delete(currentP); END WHILE; // seeds <> Empty
  RETURN True; END IF; END;
```



# DBSCAN



Date (puncte) de  
prelucrat



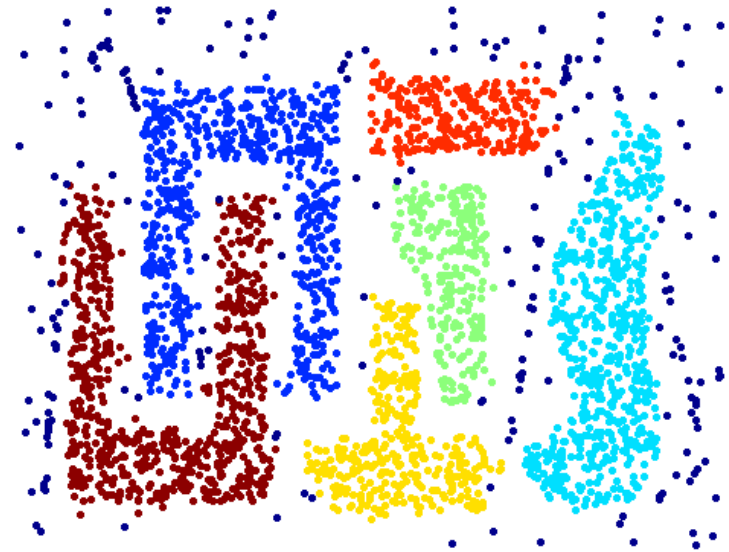
Tipuri de puncte : core,  
border și noise

Eps = 10, MinPts = 4

# DBSCAN



Date

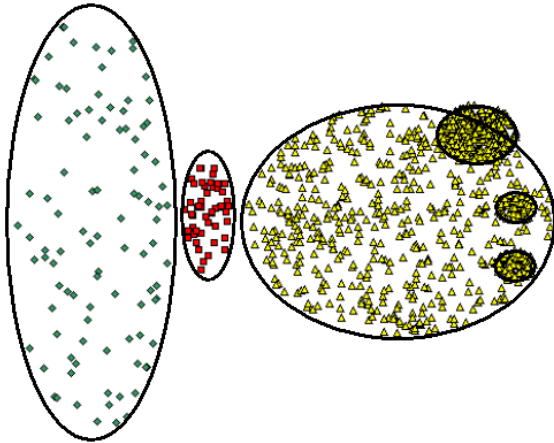


Clustere

Specific:

- Permite identificarea clusterelor de diferite forme
- Robust

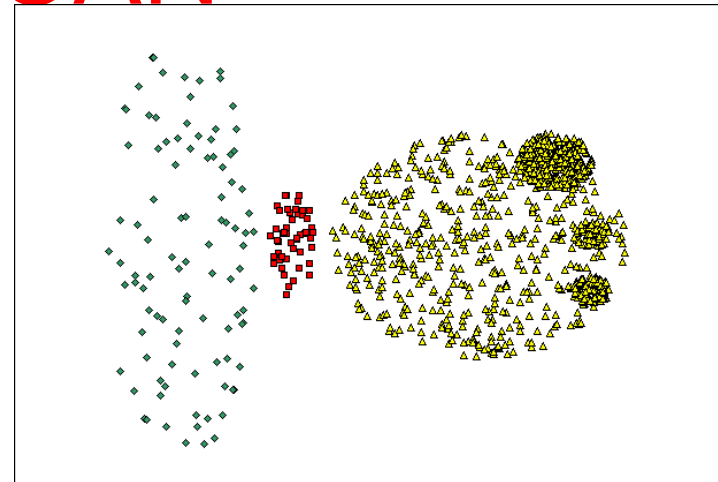
# DBSCAN



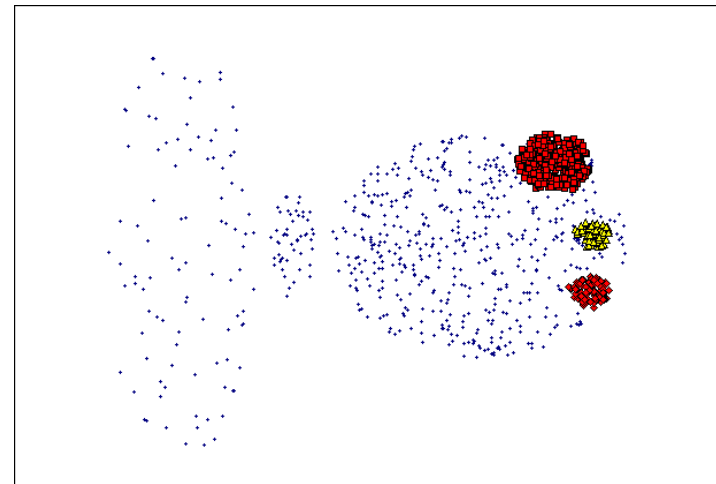
Original Points

Nu este adecvat pt:

- variații în densitatea datelor
- date de dimensiuni mari (cu multe attribute)

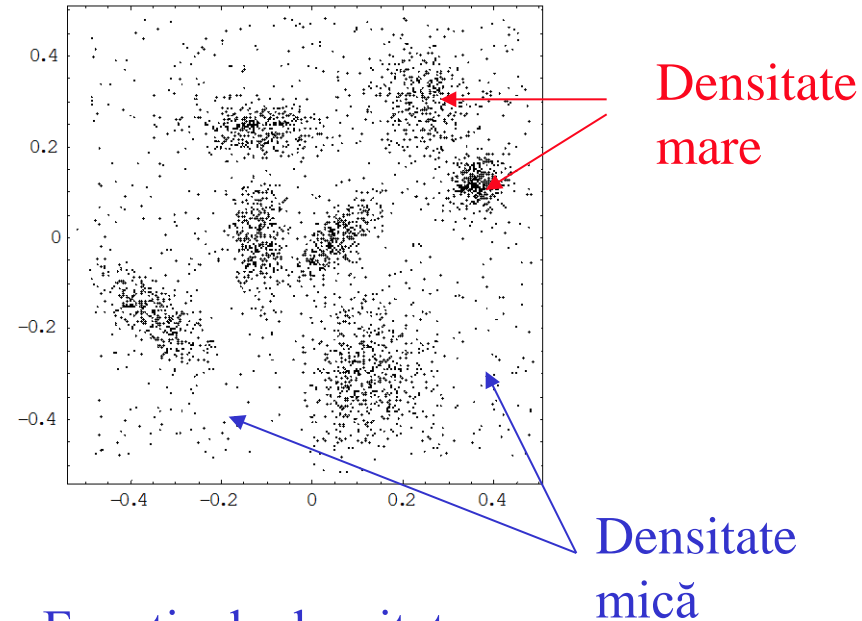


(MinPts=4, Eps=9.75).



# DENCLUE

- **Cluster** = grup dens de date similare separate de regiuni cu densitate mai mică de date
- **Idee de bază:** estimarea densității locale a datelor
  - se utilizează funcții de influență pt estimarea densității



Funcție de influență

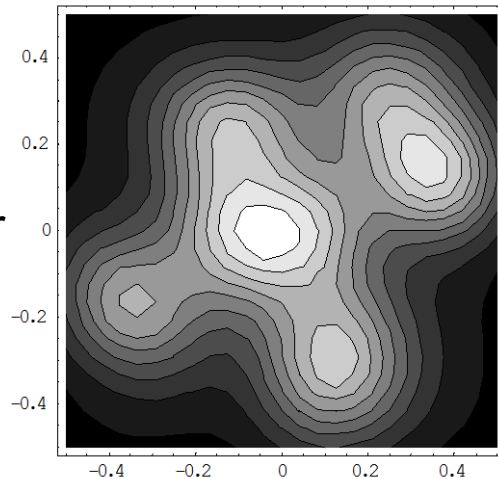
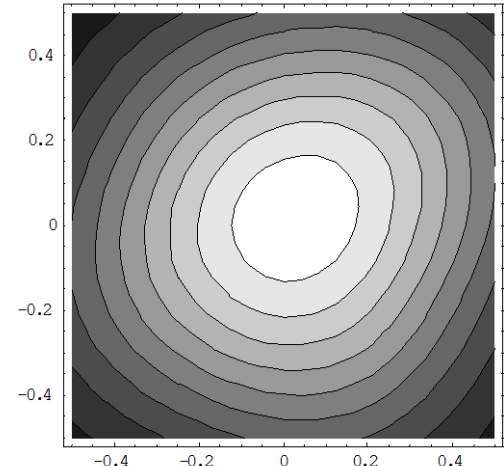
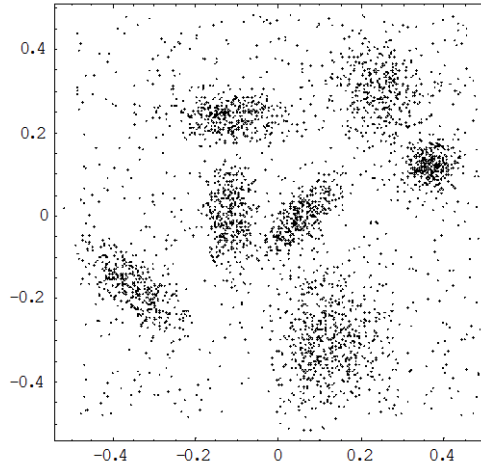
$$I_y(x) = \frac{1}{\sigma^{n/2}} \exp\left(-\frac{\sum_{j=1}^n (x_j - y_j)^2}{2\sigma^2}\right)$$

Funcție de densitate

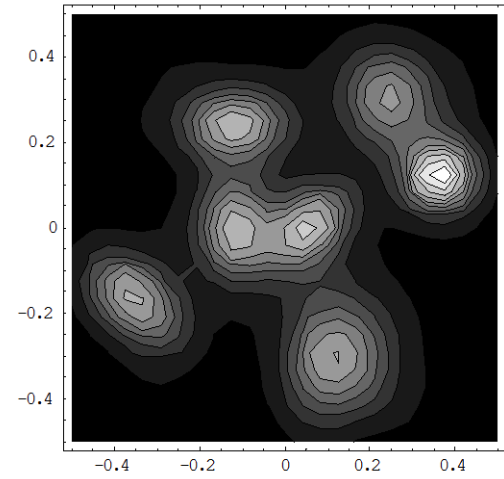
$$f(x) = \frac{1}{N} \sum_{i=1}^N I_{x_i}(x)$$

# DENCLUE

- Forma funcției de densitate depinde de valoarea lui  $\sigma$
- Dacă valoarea lui  $\sigma$  este adecvată, maximele locale ale funcției de densitate corespund reprezentanților clusterilor
- Pt valori mari ale lui  $\sigma$  funcția de densitate are un maxim unic
- Pt valori prea mici ale lui  $\sigma$  maximele locale corespund unor vârfuri izolate și pot fi dificil de detectat



$\sigma=0.1$



$\sigma=0.05$

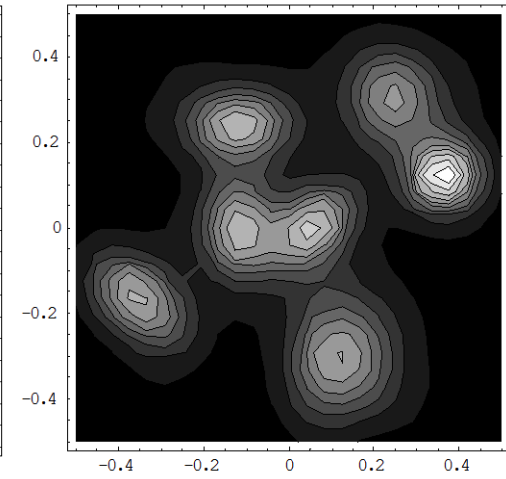
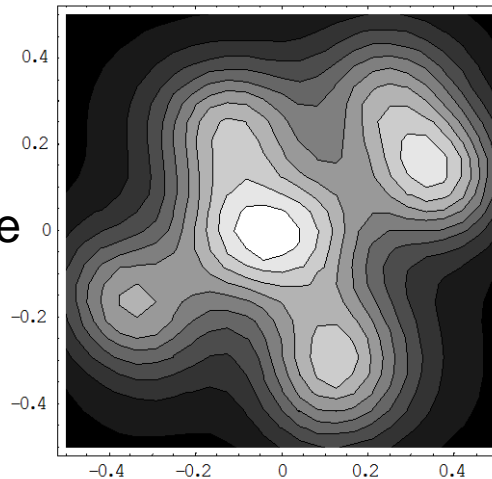
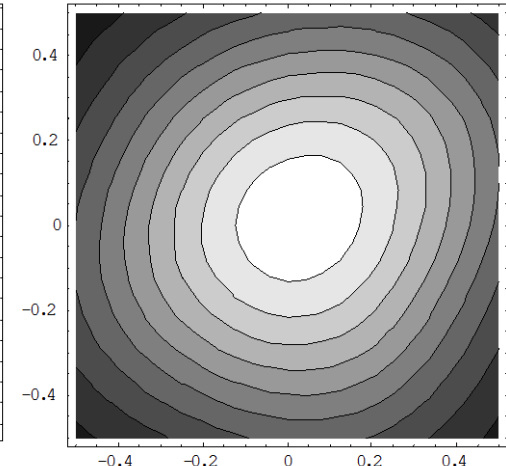
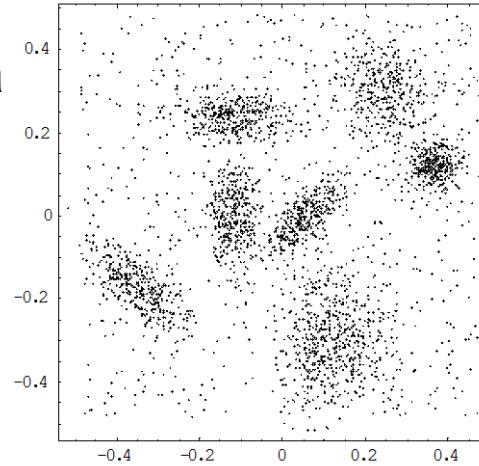
# DENCLUE

Ideea algoritmului DENCLUE

[Hinneburg, Keim – 1998]: se aplică de căutare de tip gradient pornind de la punctele din setul de date cu scopul identificării maximelor locale

Variante:

- Fiecare maxim local corespunde unui cluster (clusterelor detectate vor fi sferice sau elipsoidale)
- Un cluster corespunde unui set de maxime locale “învecinate” (se pot identifica cluster de forma arbitrara)



# Metode probabiliste

## Idee de bay[:

- Datele sunt generate de un proces stohastic (o **mixtură** de distribuții de probabilitate, fiecare dintre ele corespunzând unui cluster)
- Scopul algoritmului de grupare este de a descoperi modelul probabilist, adică de a identifica distribuțiile de probabilitate

## Exemple:

- Algoritmul Expectation–Maximization (EM) – se bazează pe următoarele ipoteze:
  - Fiecare dată este generată de o distribuție de probabilitate (tipul de distribuție depinde de natura datelor)
  - În procesul de generare a datelor fiecare dintre distribuțiile de probabilitate este selectată la rândul ei cu o anumită probabilitate

# Algoritmul EM

- **Input:** set de date  $D=\{x_1, x_2, \dots, x_N\}$ ,  $K$  = număr de clustere
- **Output:** o partiție  $P=\{C_1, C_2, \dots, C_K\}$  a setului  $D$

Inițializarea parametrilor modelelor și a probabilităților de selecție

- **(E-Step)** Se determină probabilitatea de asignare a fiecărei date la clustere (folosind valorile curente ale parametrilor)
- **(M-Step)** Se determină parametrii modelului folosind valorile curente ale probabilităților de asignare a datelor la clustere

Obs:

- In cazul datelor numerice se poate folosi distribuția normală multi-dimensională pentru a modela datele aparținând unui cluster (modelul de generare a datelor este o mixtură de gaussiene). Parametrii care se estimează sunt media și matricea de covarianță.



# EM algorithm

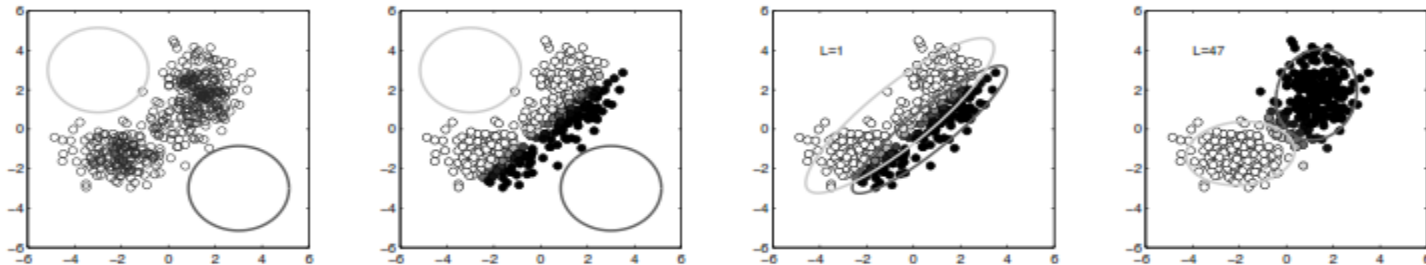
---

**Algorithm 11** EM for Gaussian Mixtures

---

Given a set of data points and a Gaussian mixture model, the goal is to maximize the log-likelihood with respect to the parameters.

- 1: Initialize the means  $\mu_k^0$ , covariances  $\Sigma_k^0$ , and mixing probabilities  $\pi_k^0$ .
  - 2: **E-step**: Calculate the responsibilities  $\gamma(z_{nk})$  using the current parameters based on Equation (3.13).
  - 3: **M-step**: Update the parameters using the current responsibilities. Note that we first update the new means using (3.12), then use these new values to calculate the covariances using (3.14), and finally reestimate the mixing probabilities using (3.15).
  - 4: Compute the log-likelihood using (3.10) and check for convergence of the algorithm. If the convergence criterion is not satisfied, then repeat steps 2–4; otherwise, return the final parameters.
- 



Data Clustering <sup>(a)</sup> Algorithms and applications <sup>(b)</sup> (ed. CC Aggarwal, CK Reddy, <sup>(c)</sup> <sup>(d)</sup> 2014)

# Algoritmul EM

Ecuții care intervin în algoritmul EM

$$p(\mathbf{x}_n|\Theta) = p(\mathbf{x}_n|\pi, \mu, \Sigma) = \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n|\mu_k, \Sigma_k). \quad (3.10)$$

$$\mathcal{N}(\mathbf{x}|\mu, \Sigma) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} \exp\left\{-\frac{1}{2}(\mathbf{x} - \mu)^T \Sigma^{-1}(\mathbf{x} - \mu)\right\}$$

$$\mu_k = \frac{\sum_{n=1}^N \gamma(z_{nk}) \mathbf{x}_n}{\sum_{n=1}^N \gamma(z_{nk})} \quad (3.12)$$

$$\gamma(z_{nk}) = \frac{\pi_k \mathcal{N}(\mathbf{x}_n|\mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}_n|\mu_j, \Sigma_j)}. \quad (3.13)$$

$$\Sigma_k = \frac{\sum_{n=1}^N \gamma(z_{nk}) (\mathbf{x}_n - \mu_k)(\mathbf{x}_n - \mu_k)^T}{\sum_{n=1}^N \gamma(z_{nk})}. \quad (3.14)$$

$$\pi_k = \frac{\sum_{n=1}^N \gamma(z_{nk})}{N} \quad (3.15)$$

Data Clustering Algorithms and applications (ed. CC Aggarwal, CK Reddy, 2014)