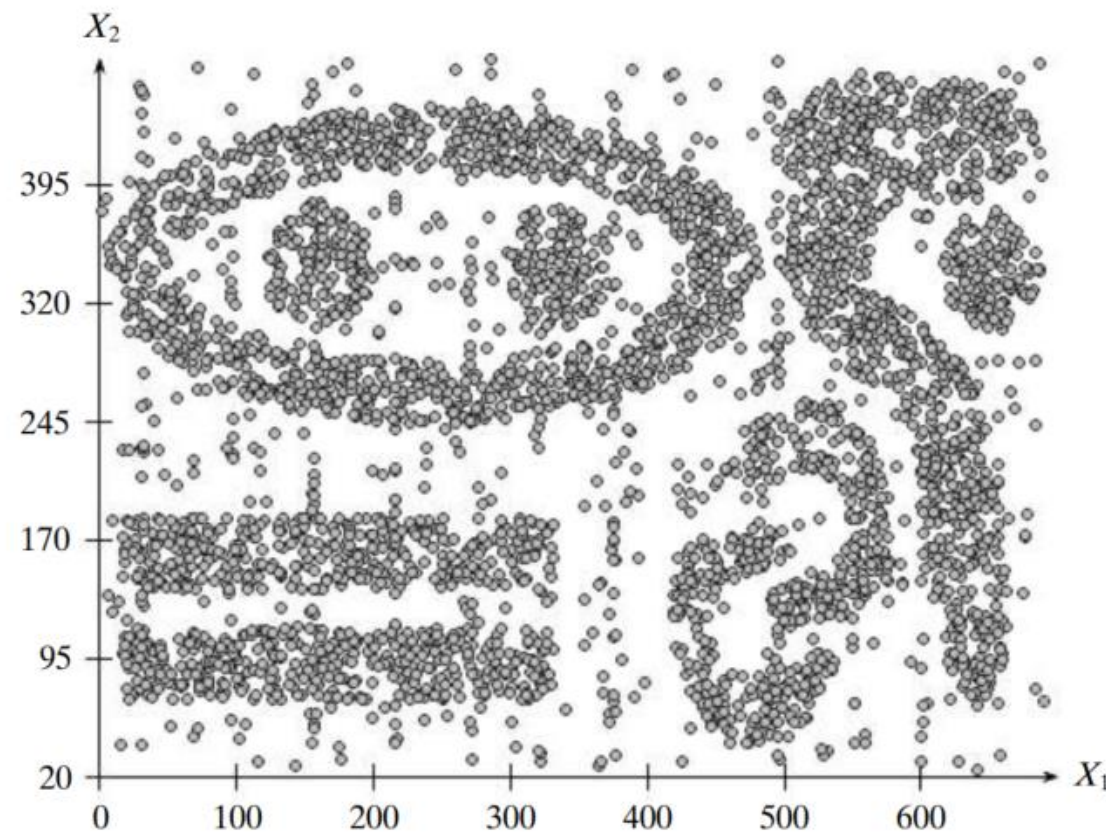


خوشه‌بندی مبتنی بر چگالی

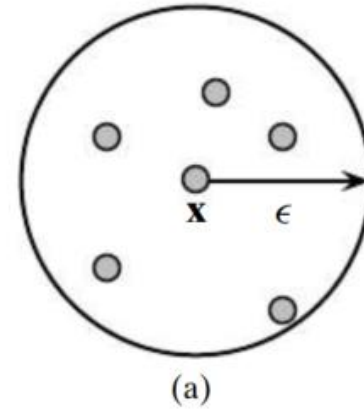
Density-based Clustering

خوشه‌بندی‌های مبتنی بر نماینده برای خوشه‌های غیر محدب (Nonconvex Clusters) مناسب نیستند. زیرا فاصله‌ی دو نقطه از دو خوشه‌ی متفاوت می‌تواند از فاصله‌ی آن‌ها از نماینده‌اشان کوتاه‌تر باشد.



Density-based dataset.

ϵ -همسایگی (ϵ -neighborhood) x توپی به شعاع ϵ حول نقطه‌ی $x \in \mathbb{R}^d$ است.



(a)

$$N_\epsilon(x) = B_d(x, \epsilon) = \{y | \delta(x, y) \leq \epsilon\}$$

که معمولا از فاصله‌ی اقلیدسی استفاده می‌شود.
 $\delta(x, y) = \|x - y\|_2$

(a) Neighborhood of a point.

اگر $x \in N_\epsilon(y)$ باشد و y یک نقطه‌ی اصلی باشد،
 آنگاه x مستقیما دسترسی چگالی به y دارد (Directly Density Reachable).

اگر یک زنجیره‌ای از نقاط اصلی x_0, x_1, \dots, x_l وجود داشته باشد که
 هر دو نقطه‌ی دنبال هم، مستقیما دسترسی چگالی به هم داشته باشند،
 آنگاه دو نقطه‌ی ابتدایی و انتهایی x_l و x_0 نیز دسترسی چگالی به هم دارند.

اگر نقطه‌ی اصلی z وجود داشته باشد که دو نقطه‌ی x و y به نقطه‌ی z
 دسترسی چگالی داشته باشند،
 متصل چگالی (Density Directed) می‌باشند.

یک خوشه‌ی مبتنی بر چگالی، مجموعه بیشینه‌ی نقاط متصل چگالی است.

نقطه‌ی اصلی (Core Point): $|N_\epsilon(x)| \geq \text{minpts}$

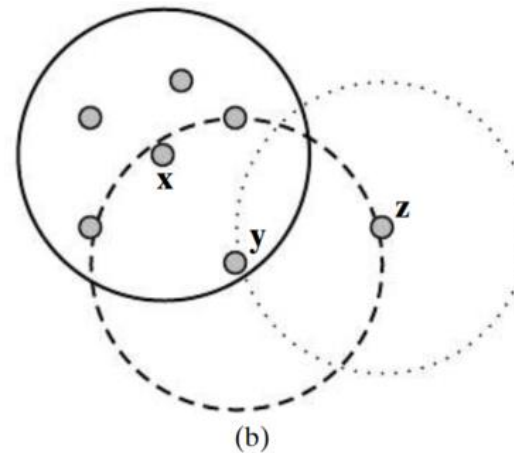
نقطه‌ی مرزی (Border Point):

$$|N_\epsilon(x)| < \text{minpts}$$

اما در همسایگی یک نقطه‌ی اصلی قرار دارد.

نقطه‌ی یرت یا نویز (Noise Point):

هر نقطه‌ای که نه اصلی و نه مرزی باشد.



(b)

(b) Core, border, and noise points.

Algorithm 15.1: Density-based Clustering Algorithm

DBSCAN ($\mathbf{D}, \epsilon, minpts$):

```

1  $Core \leftarrow \emptyset$ 
2 foreach  $\mathbf{x}_i \in \mathbf{D}$  do // Find the core points
3   Compute  $N_\epsilon(\mathbf{x}_i)$ 
4    $id(\mathbf{x}_i) \leftarrow \emptyset$  // cluster id for  $\mathbf{x}_i$ 
5   if  $N_\epsilon(\mathbf{x}_i) \geq minpts$  then  $Core \leftarrow Core \cup \{\mathbf{x}_i\}$ 
6  $k \leftarrow 0$  // cluster id
7 foreach  $\mathbf{x}_i \in Core$ , such that  $id(\mathbf{x}_i) = \emptyset$  do
8    $k \leftarrow k + 1$ 
9    $id(\mathbf{x}_i) \leftarrow k$  // assign  $\mathbf{x}_i$  to cluster id  $k$ 
10  DENSITYCONNECTED ( $\mathbf{x}_i, k$ )
11  $\mathcal{C} \leftarrow \{C_i\}_{i=1}^k$ , where  $C_i \leftarrow \{\mathbf{x} \in \mathbf{D} \mid id(\mathbf{x}) = i\}$ 
12  $Noise \leftarrow \{\mathbf{x} \in \mathbf{D} \mid id(\mathbf{x}) = \emptyset\}$ 
13  $Border \leftarrow \mathbf{D} \setminus \{Core \cup Noise\}$ 
14 return  $\mathcal{C}, Core, Border, Noise$ 

```

DENSITYCONNECTED (\mathbf{x}, k):

```

15 foreach  $\mathbf{y} \in N_\epsilon(\mathbf{x})$  do
16    $id(\mathbf{y}) \leftarrow k$  // assign  $\mathbf{y}$  to cluster id  $k$ 
17   if  $\mathbf{y} \in Core$  then DENSITYCONNECTED ( $\mathbf{y}, k$ )

```

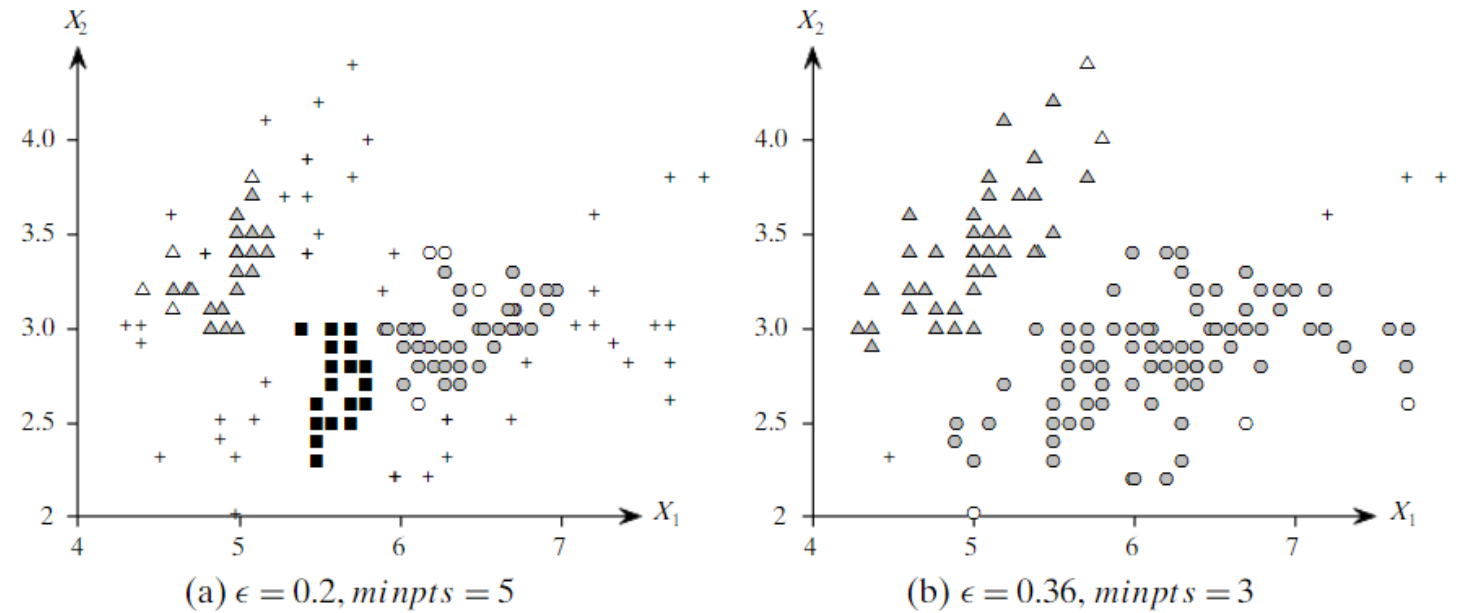


Figure 15.4. DBSCAN clustering: Iris dataset.

برآورد تابع توزیع تجمعی (Cumulative Distribution Function) بطور مستقیم از داده

فرض کنید X متغیر تصادفی پیوسته با تابع چگالی (نامشخص برای ما) $f(x)$ باشد. هدف برآورد $f(x)$ به کمک n نمونه x_1, \dots, x_n می باشد.

برآوردگر هسته‌ای (Kernel Estimator)

خواص تابع هسته‌ای چگالی (Density Kernel Function)

$$K(x) \geq 0, \quad \int K(x)dx = 1, \quad K(x) = K(-x)$$

$$\hat{f}(x) = \frac{1}{nh} \sum_{i=1}^n K\left(\frac{x - x_i}{h}\right),$$

$$\hat{F}(x) = \frac{1}{n} \sum_{i=1}^n I(x_i \leq x)$$

$$\hat{f}(x) = \frac{\hat{F}\left(x + \frac{h}{2}\right) - \hat{F}\left(x - \frac{h}{2}\right)}{h} = \frac{k/n}{h} = \frac{k}{nh}$$

k تعداد نقاط درون پنجره‌ای به پهنای h و مرکز x است و تابع چگالی برآورد شده می‌شود، نسبت تعداد نقاط درون پنجره حول x ، به تعداد کل نقاط به حجم پنجره. دقت الگوریتم به پارامتر h بستگی دارد.

هسته‌ی گاوسی (Gaussian Kernel)

$$K(x) = \frac{1}{\sqrt{2\pi}} \exp\left\{-\frac{z^2}{2}\right\}$$

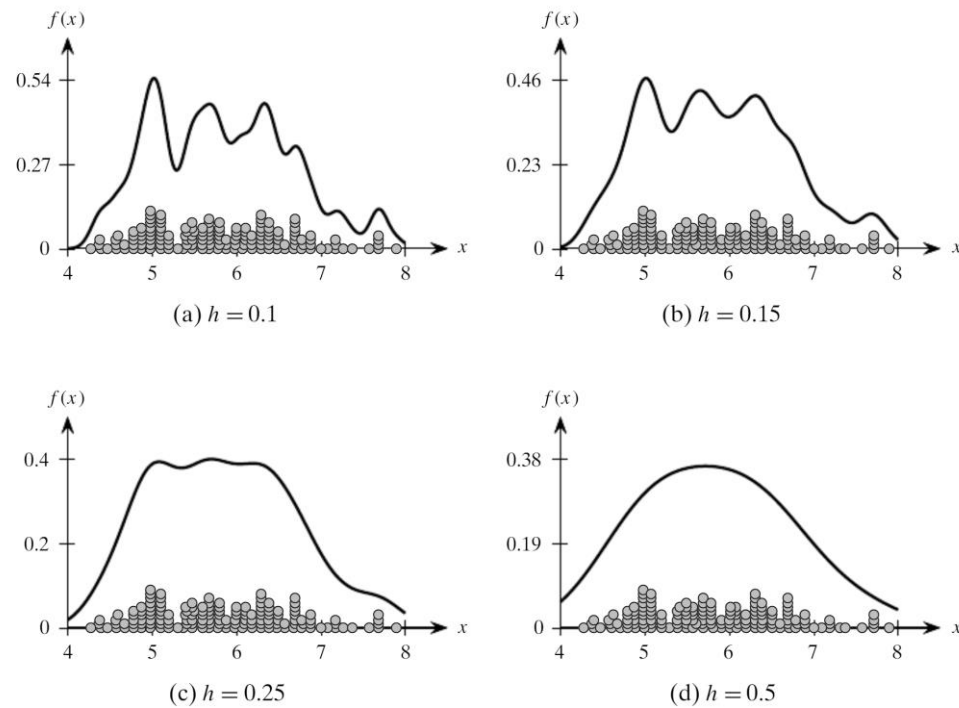


Figure 15.6. Kernel density estimation: Gaussian kernel (varying h).

هسته‌ی گسسته (Discrete Kernel)

$$K(x) = \begin{cases} 1 & |z| \leq 1/2 \\ 0 & \text{otherwise} \end{cases}$$

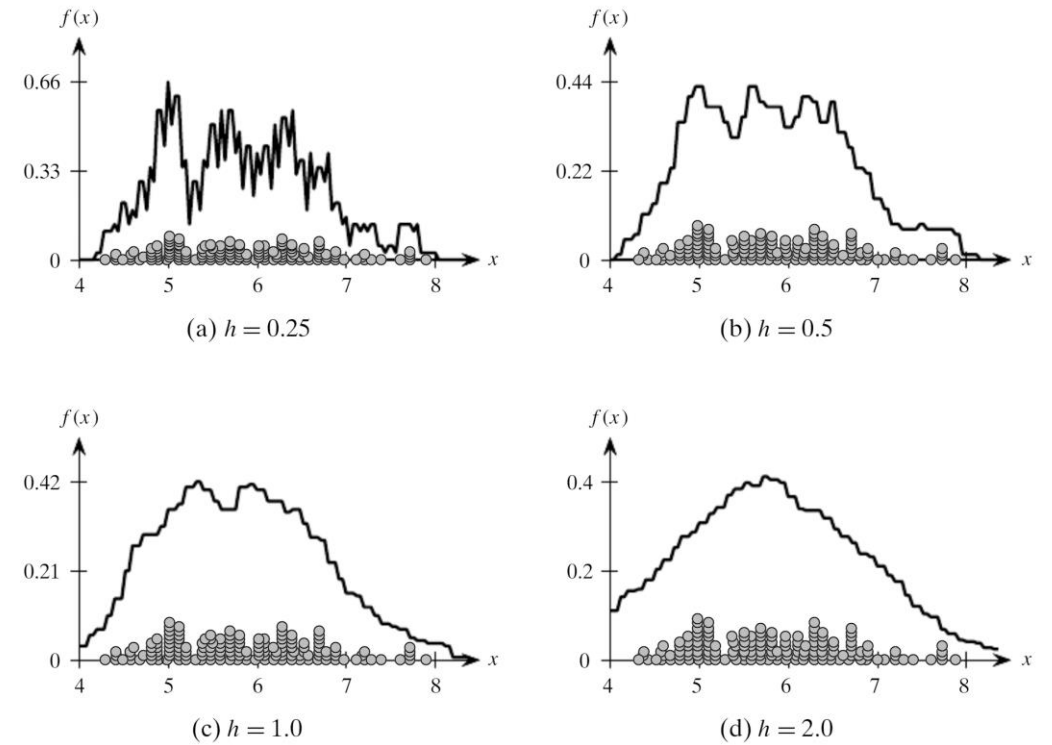


Figure 15.5. Kernel density estimation: discrete kernel (varying h).

برآورد چگالی چندمتغیره (Multivariate Density Estimation)

$$\mathbf{x} = (x_1, \dots, x_d)^T$$

$$K(\mathbf{x}) \geq 0, \quad \int K(\mathbf{x}) d\mathbf{x} = 1, \quad K(\mathbf{x}) = K(-\mathbf{x})$$

$$\hat{f}(\mathbf{x}) = \frac{1}{nh^d} \sum_{i=1}^n K\left(\frac{\mathbf{x} - \mathbf{x}_i}{h}\right),$$

هسته‌ی گسسته (Discrete Kernel)

$$K(\mathbf{z}) = \begin{cases} 1 & \text{If } |z_j| \leq \frac{1}{2}, \text{ for all dimensions } j = 1, \dots, d \\ 0 & \text{Otherwise} \end{cases}$$

هسته‌ی گاوسی (Gaussian Kernel)

$$K(\mathbf{z}) = \frac{1}{(2\pi)^{d/2}} \exp \left\{ -\frac{\mathbf{z}^T \mathbf{z}}{2} \right\} \quad \Sigma = \mathbf{I}_d.$$

$$K \left(\frac{\mathbf{x} - \mathbf{x}_i}{h} \right) = \frac{1}{(2\pi)^{d/2}} \exp \left\{ -\frac{(\mathbf{x} - \mathbf{x}_i)^T (\mathbf{x} - \mathbf{x}_i)}{2h^2} \right\}$$

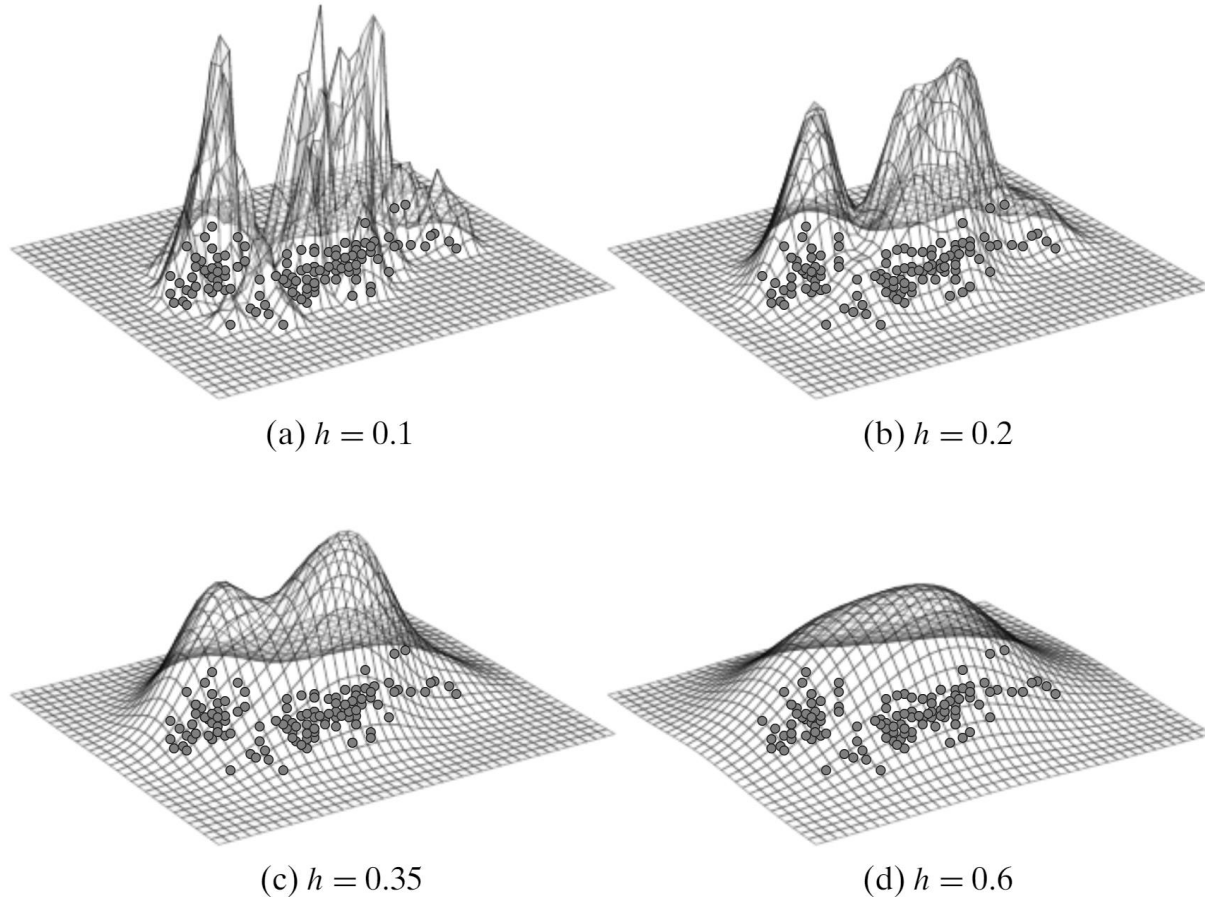


Figure 15.7. Density estimation: 2D Iris dataset (varying h).

$$\hat{f}(\mathbf{x}) = \frac{k}{n \text{vol}(S_d(h_{\mathbf{x}}))}$$

$$\text{vol}(S_d(r)) = K_d \cdot r^d = \left(\frac{\pi^{\frac{d}{2}}}{\Gamma(\frac{d}{2} + 1)} \right) r^d$$

$$\Gamma\left(\frac{d}{2} + 1\right) = \begin{cases} \left(\frac{d}{2}\right)! & \text{if } d \text{ is even} \\ \sqrt{\pi} \left(\frac{d!!}{2^{(d+1)/2}}\right) & \text{if } d \text{ is odd} \end{cases}$$

در این رابطه h_x فاصله‌ی k امین نزدیکترین همسایه‌ی x و $\text{vol}(S_d(h_x))$ حجم کره‌ی d بعدی به مرکز x و شعاع h_x است. مقدار k از ابتدا مشخص می‌شود.

$$\Gamma(\alpha) = \int_0^{\infty} x^{\alpha-1} e^{-x} dx$$

$$\Gamma(\alpha) = (\alpha - 1)\Gamma(\alpha - 1)$$

$$\Gamma(1) = 1 \quad \text{and} \quad \Gamma\left(\frac{1}{2}\right) = \sqrt{\pi}$$

هر نقطه‌ی x جذب یک x^* (Density Attractor) که چگالی بیشتر از آستانه دارد، می‌شود.

$$\mathbf{x}_{t+1} = \mathbf{x}_t + \eta \cdot \nabla \hat{f}(\mathbf{x}_t)$$

$$\nabla \hat{f}(\mathbf{x}) = \frac{1}{nh^{d+2}} \sum_{i=1}^n K\left(\frac{\mathbf{x} - \mathbf{x}_i}{h}\right) \cdot (\mathbf{x}_i - \mathbf{x})$$

$$\mathbf{x}_{t+1} = \frac{\sum_{i=1}^n K\left(\frac{\mathbf{x}_t - \mathbf{x}_i}{h}\right) \mathbf{x}_i}{\sum_{i=1}^n K\left(\frac{\mathbf{x}_t - \mathbf{x}_i}{h}\right)}$$

$$\mathbf{x}_{t+1} = \frac{\sum_{\mathbf{x}_i \in B_d(\mathbf{x}_t, r)} K\left(\frac{\mathbf{x}_t - \mathbf{x}_i}{h}\right) \mathbf{x}_i}{\sum_{\mathbf{x}_i \in B_d(\mathbf{x}_t, r)} K\left(\frac{\mathbf{x}_t - \mathbf{x}_i}{h}\right)}$$

برای افزایش سرعت می‌توان جمع را روی همسایه‌ها با فاصله‌ی zh انجام داد. $r = zh$

$$\hat{f}(\mathbf{x}^*) = \frac{1}{nh^d} \sum_{i=1}^n K\left(\frac{\mathbf{x}^* - \mathbf{x}_i}{h}\right) \geq \xi$$

ξ یک مقدار از پیش مشخص شده برای آستانه چگال است.

Algorithm 15.2: DENCLUE Algorithm

DENCLUE ($\mathbf{D}, h, \xi, \epsilon$):

- 1 $\mathcal{A} \leftarrow \emptyset$
- 2 **foreach** $\mathbf{x} \in \mathbf{D}$ **do** // find density attractors
- 3 $\mathbf{x}^* \leftarrow \text{FINDATTRACTOR}(\mathbf{x}, \mathbf{D}, h, \epsilon)$
- 4 **if** $\hat{f}(\mathbf{x}^*) \geq \xi$ **then**
- 5 $\mathcal{A} \leftarrow \mathcal{A} \cup \{\mathbf{x}^*\}$
- 6 $R(\mathbf{x}^*) \leftarrow R(\mathbf{x}^*) \cup \{\mathbf{x}\}$
- 7 $\mathcal{C} \leftarrow \{\text{maximal } C \subseteq \mathcal{A} \mid \forall \mathbf{x}_i^*, \mathbf{x}_j^* \in C, \mathbf{x}_i^* \text{ and } \mathbf{x}_j^* \text{ are density reachable}\}$
- 8 **foreach** $C \in \mathcal{C}$ **do** // density-based clusters
- 9 **foreach** $\mathbf{x}^* \in C$ **do** $C \leftarrow C \cup R(\mathbf{x}^*)$
- 10 **return** \mathcal{C}

FINDATTRACTOR ($\mathbf{x}, \mathbf{D}, h, \epsilon$):

- 11 $t \leftarrow 0$
- 12 $\mathbf{x}_t \leftarrow \mathbf{x}$
- 13 **repeat**
- 14 $\mathbf{x}_{t+1} \leftarrow \frac{\sum_{i=1}^n K\left(\frac{\mathbf{x}_t - \mathbf{x}_i}{h}\right) \cdot \mathbf{x}_i}{\sum_{i=1}^n K\left(\frac{\mathbf{x}_t - \mathbf{x}_i}{h}\right)}$
- 15 $t \leftarrow t + 1$
- 16 **until** $\|\mathbf{x}_t - \mathbf{x}_{t-1}\| \leq \epsilon$
- 17 **return** \mathbf{x}_t

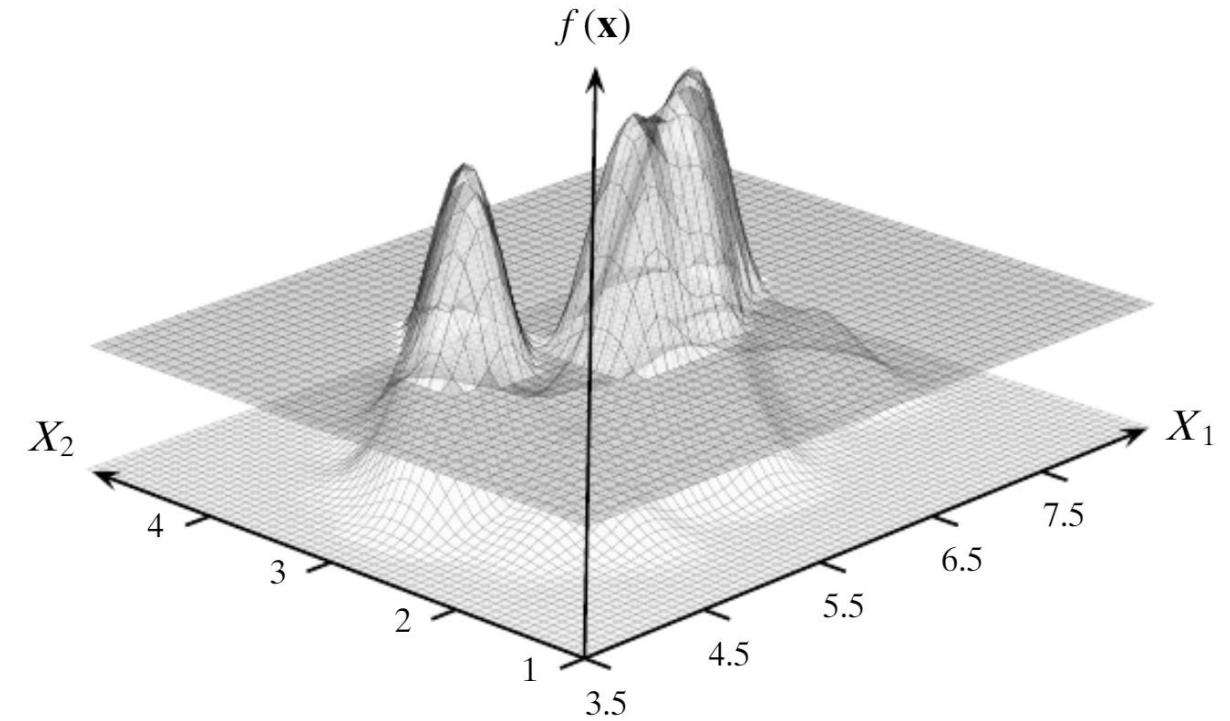


Figure 15.10. DENCLUE: Iris 2D dataset.

می‌توان نشان داد که DBSCAN حالت خاصی از DENCLUE (برآورد چگالی هسته) با انتخاب $h = \epsilon$ و $\xi = \text{minpts}$ است.

دو جاذب x_i^* و x_j^* از طریق چگالی به هم در دسترس می‌باشند اگر بتوان مسیری بین آن دو پیدا کرد که برای همه نقاط \mathbf{y} روی مسیر $\hat{f}(\mathbf{y}) \geq \xi$