

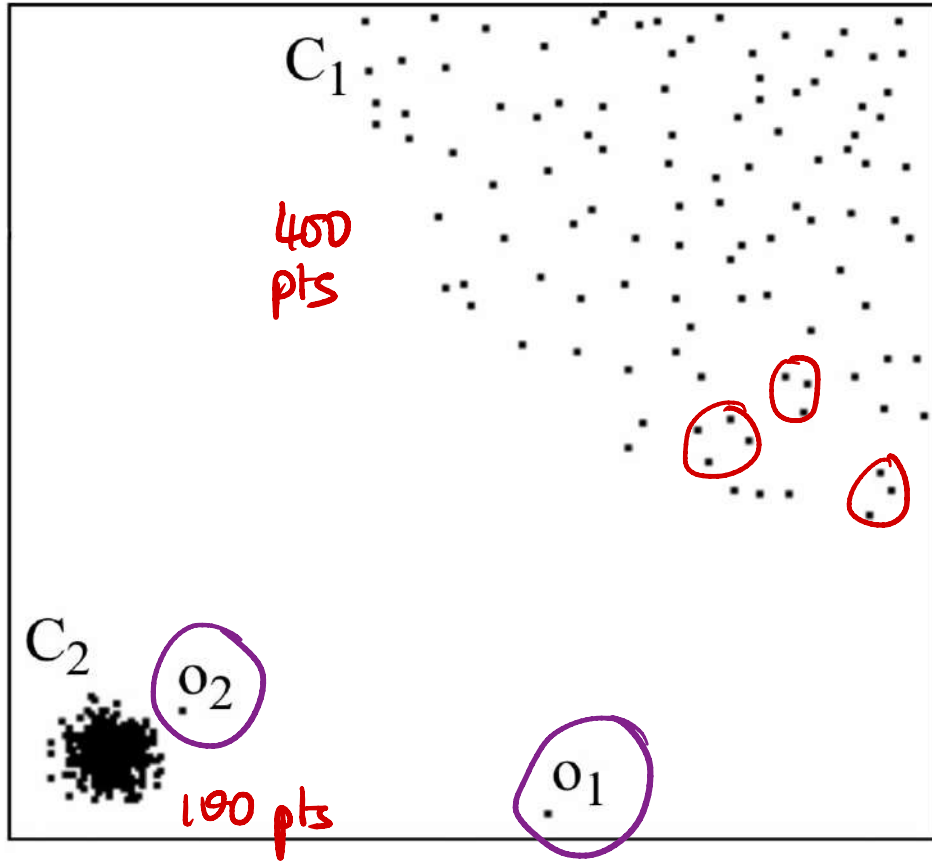
DMML 18 Feb 2020

Clustering - k Means, Hierarchical, Density Based
DBScan

Outliers

Density based - outlier not a core
not reachable from core

Issue - uniform density ?

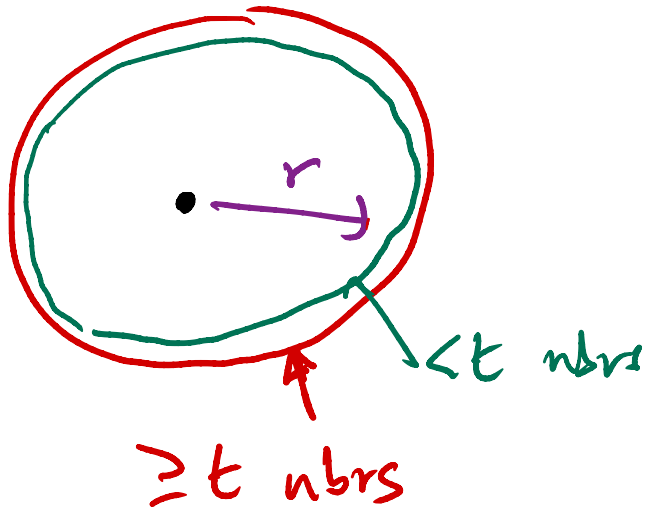


Uniform density
measure will not
work

Local density metric

Previously

- Fix radius r
- Check if t nbrs
within r
- Fix t
- Find r containing t nbrs

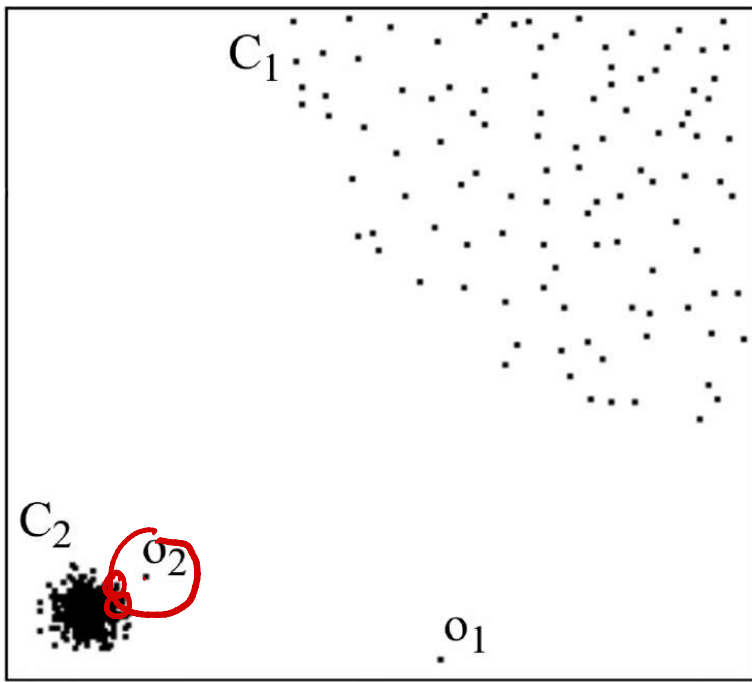


Min-t-radius (p)

Compare this to
 min-t-radius (p') of
 the nearest t nbrs
 \Downarrow
 p

$$\left[\frac{\text{min-radius}(p)}{\text{mean}(\text{min-radius}(p'), p' \text{ t-nbr } \Downarrow p)} \right]$$

"local outlier factor" — larger \Rightarrow outlier



Find outliers

Then cluster

- K Means itself is
good enough

Applications of clustering

Image processing → feature extracting

Colour based clustering of pixels

Original image



10 colors



8 colors



6 colors



4 colors



2 colors



Clustering as preprocessing for classification

MNIST handwritten digits

Raw classification $\sim 97\%$

Cluster images using K-Means $K=2, 3, \dots, 100$

- $K=10$
- each cluster has a centroid
 - assign a digit label to each cluster
 - use distance to centroid as new value

$K=99$ - accuracy improves to about 99%

Mixture Models

Model with parameters \rightarrow observations \rightarrow MLE

2 coins — P_1, P_2

Repeat 100 times $\left[\begin{array}{l} - \text{Choose coin uniformly with probability } 1/2 \\ - \text{Toss it} \end{array} \right.$

H T H T T H H T T H T H H ...

Segregate into red sequence $\rightarrow P_1$

blue sequence $\rightarrow P_2$

What if we don't know the red-blue label?

Iterative approximate

Initial guess \hat{p}_1, \hat{p}_2

H: Com 1: \hat{p}_1
Com 2: \hat{p}_2

T: $1 - \hat{p}_1 = \hat{q}_1$
 $1 - \hat{p}_2 = \hat{q}_2$

Red $\frac{\hat{p}_1}{\hat{p}_1 + \hat{p}_2}$ Blue $\frac{\hat{p}_2}{\hat{p}_1 + \hat{p}_2}$

$\frac{\hat{q}_1}{\hat{q}_1 + \hat{q}_2}$ $\frac{\hat{q}_2}{\hat{q}_1 + \hat{q}_2}$

$$\hat{p}_1 = 0.6 \quad \hat{p}_2 = 0.3$$

	H	T	H	T	T	H	T	T	H	T	
Red	<u>$\frac{2}{3}$</u>	$\frac{4}{11}$	<u>$\frac{2}{3}$</u>	$\frac{4}{11}$	$\frac{4}{11}$	<u>$\frac{2}{3}$</u>	$\frac{4}{11}$	$\frac{4}{11}$	<u>$\frac{2}{3}$</u>	$\frac{4}{11}$	$-\Sigma$
Blue	$\frac{1}{3}$	$\frac{7}{11}$	$\frac{1}{3}$	$\frac{7}{11}$	$\frac{7}{11}$	$\frac{1}{3}$	$\frac{7}{11}$	$\frac{7}{11}$	$\frac{1}{3}$	$\frac{7}{11}$	\downarrow
											total
											red
											coin
											total

$$\text{Revised } \hat{p}_1 = \frac{\frac{2}{3} \times 4}{\frac{2}{3} \times 4 + \frac{4}{11} \times 6}$$

$$\hat{p}_2 = \frac{\frac{1}{3} \times 4}{\frac{1}{3} \times 4 + \frac{7}{11} \times 6}$$

Iterate till convergence

Expectation - Maximization = EM

a Maximum likelihood



5 sets, 10 tosses per set

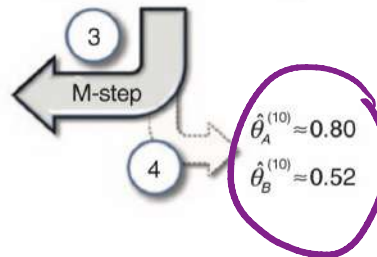
Coin A	Coin B
	5 H, 5 T
9 H, 1 T	
8 H, 2 T	
	4 H, 6 T
7 H, 3 T	
24 H, 6 T	9 H, 11 T



Coin A	Coin B
≈ 2.2 H, 2.2 T	≈ 2.8 H, 2.8 T
≈ 7.2 H, 0.8 T	≈ 1.8 H, 0.2 T
≈ 5.9 H, 1.5 T	≈ 2.1 H, 0.5 T
≈ 1.4 H, 2.1 T	≈ 2.6 H, 3.9 T
≈ 4.5 H, 1.9 T	≈ 2.5 H, 1.1 T
≈ 21.3 H, 8.6 T	≈ 11.7 H, 8.4 T

$$\hat{\theta}_A^{(1)} \approx \frac{21.3}{21.3 + 8.6} \approx 0.71$$

$$\hat{\theta}_B^{(1)} \approx \frac{11.7}{11.7 + 8.4} \approx 0.58$$

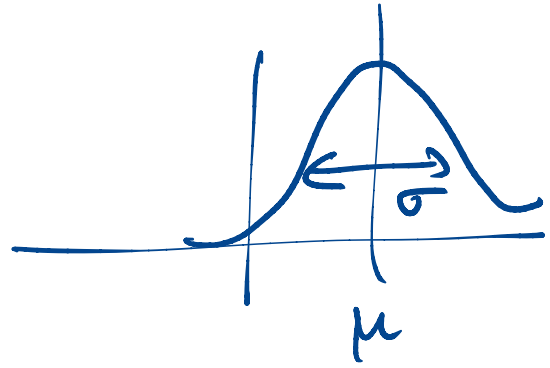


$$\frac{24}{30} = 0.8$$

$$\frac{9}{20} = 0.45$$

Mixture of Gaussian

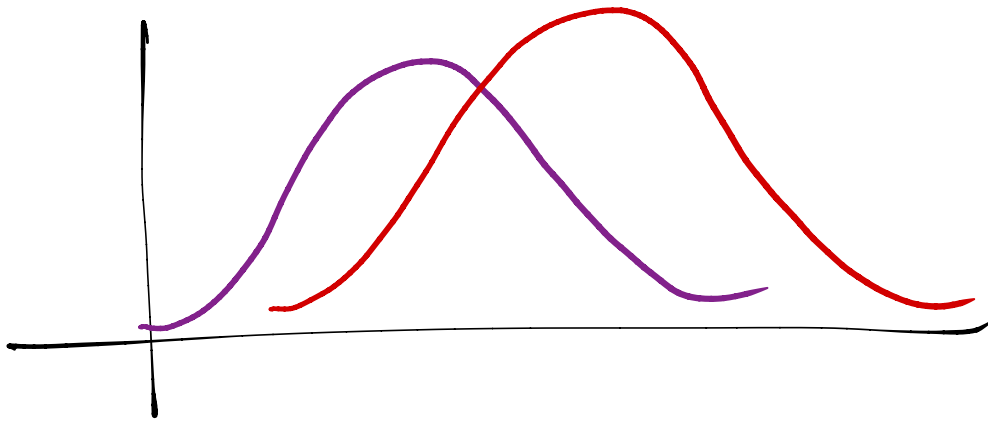
$$\frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$



Single Gaussian

$$x_1, x_2, \dots, x_n = \text{MLE of } \mu = \frac{1}{n} \sum x_i$$

$$\text{MLE of } \sigma = \sqrt{\frac{1}{n} \sum (x_i - \mu)^2}$$



Mixed sequence x_1, x_2, \dots, x_n

$$\hat{\mu}_1, \hat{\sigma}_1$$

$$\hat{\mu}_2, \hat{\sigma}_2$$

$$\begin{array}{c} \downarrow \quad \downarrow \\ \frac{e_1}{e_1 + e_2} \quad \frac{e_2}{e_1 + e_2} \end{array}$$

Mixture of Gaussians clustering

