

# Clustering

Adam Carter, EPCC, The University of Edinburgh

a.carter@epcc.ed.ac.uk

1 September 2020

[www.archer2.ac.uk](http://www.archer2.ac.uk)



| epcc |

# Reusing this material



This work is licensed under a Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International License.

<https://creativecommons.org/licenses/by-nc-sa/4.0/>

This means you are free to copy and redistribute the material and adapt and build on the material under the following terms: You must give appropriate credit, provide a link to the license and indicate if changes were made. If you adapt or build on the material, you must distribute your work under the same license as the original.

Note that this presentation contains images owned by others. Please seek their permission before reusing these images.

# Partners

| epcc |



Engineering and  
Physical Sciences  
Research Council

Natural  
Environment  
Research Council

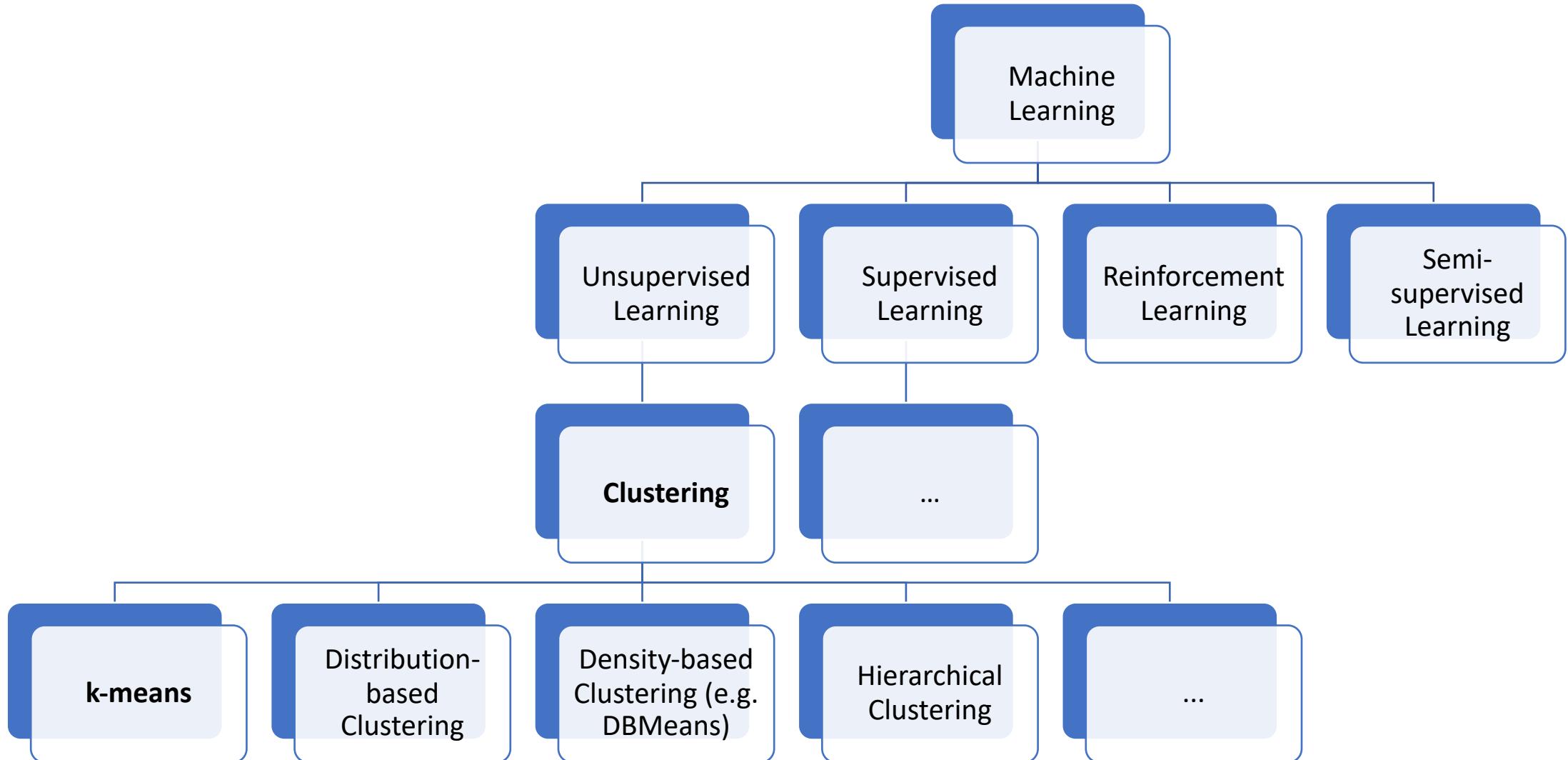


THE UNIVERSITY  
*of* EDINBURGH

| epcc |

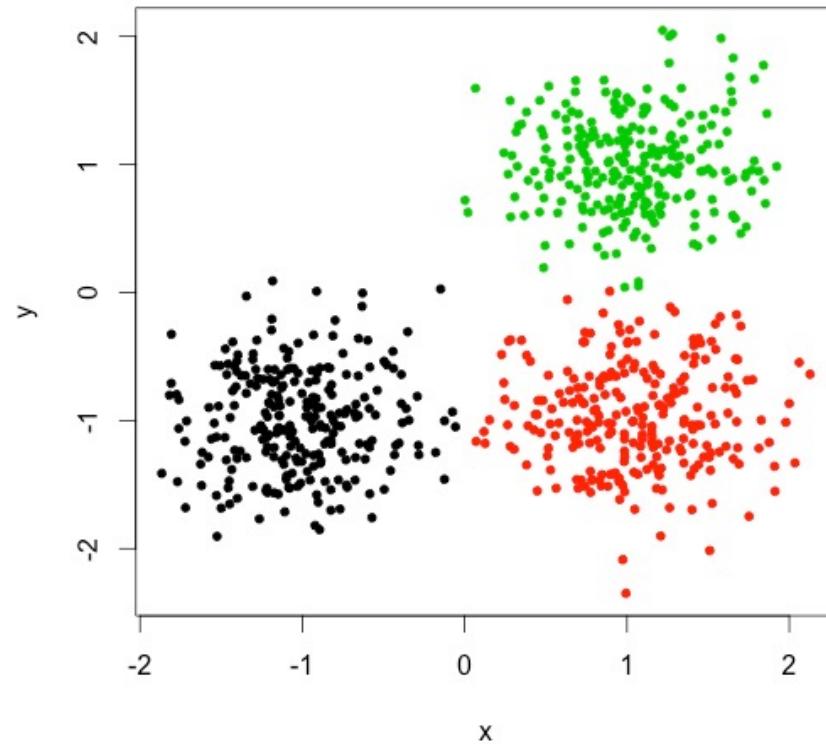
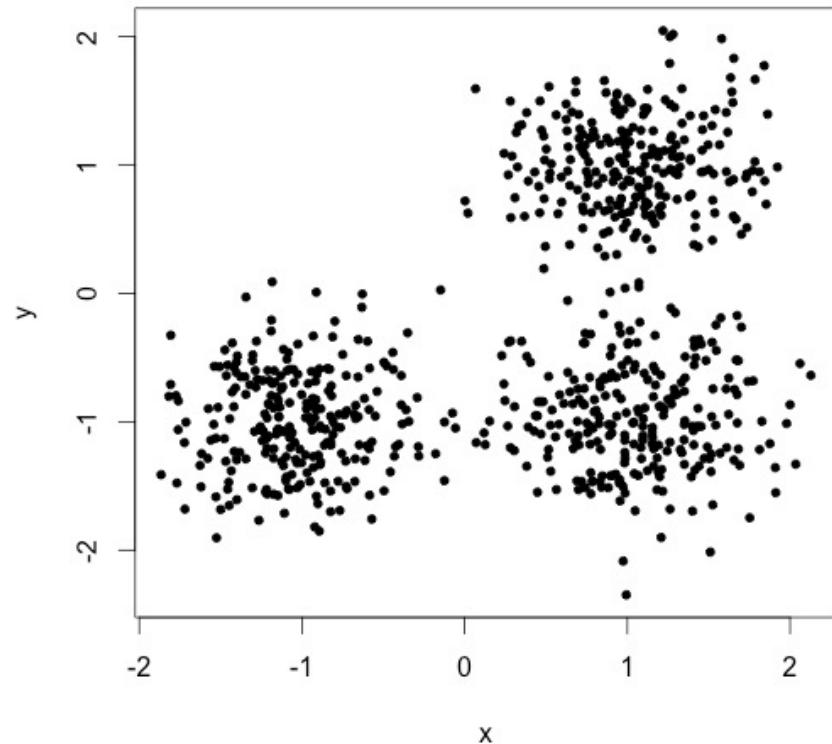


**Hewlett Packard  
Enterprise**



# Clustering – the problem

| epcc |





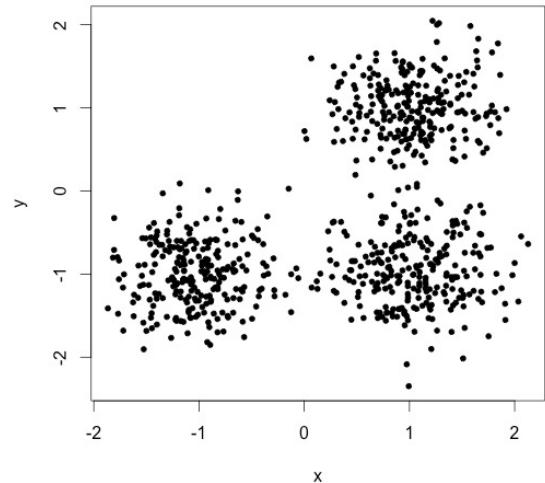
# *k*-Means Clustering



|epcc|

# Brief Aside: Notation (1)

- Unless specified, we'll try to stick to the following notation here:
- $x$  A single instance
  - A single point (as in the graph below). A single member of your population. A single observation. Usually corresponds to a single *row* in a dataset. Usually, the instance has many features or variables. The values of each feature are the values of this vector. So, if there are  $n$  features, then  $x = (x_1, x_2, \dots, x_n)$  or  $x = (x^{(1)}, x^{(2)}, \dots, x^{(n)})$
- $x_1, x_2, \dots, x_m$  A set of  $m$  observations
  - Usually corresponds to a dataset with  $m$  rows and  $n$  columns



# Brief Aside: Notation (2)



- $x$  A single instance
- $x_1, x_2, \dots, x_m$  A set of  $m$  observations
- $x_a, x_b$  Two observations  $a$  and  $b$
- $\|x_b - x_a\|$  The “distance” between the two points  $a$  and  $b$

- Usually (c.f. “square on the hypotenuse”):

$$\|x_b - x_a\| = \sqrt[2]{(x_b^{(1)} - x_a^{(1)})^2 + (x_b^{(2)} - x_a^{(2)})^2 + \dots + (x_b^{(n)} - x_a^{(n)})^2}$$

- Note that  $\|x_b - x_a\| = \|x_a - x_b\|$

# $k$ -means clustering



- Know in advance that there are  $k$  clusters
- Goal:
  - Given observation vectors:  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m$
  - Group them in  $k$  distinct sets:  $S_1, S_2, \dots, S_k$
  - Minimise the within-cluster sum of squares:
  - $\sum_{c=1}^k \sum_{x \in S_c} \|\mathbf{x} - \boldsymbol{\mu}_c\|^2$
  - where  $\boldsymbol{\mu}_c$  is the mean of the points in the set  $S_c$
  -

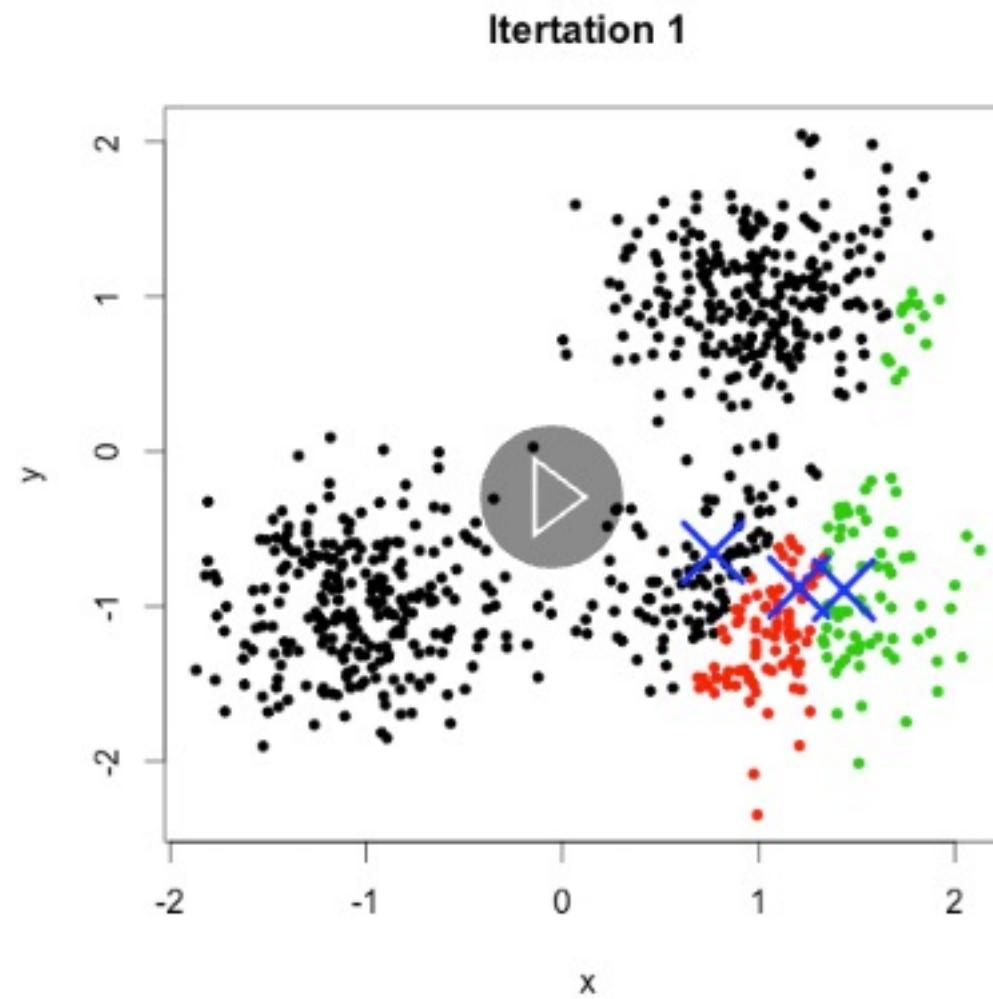
# *k*-means algorithm



- Define each cluster by the centroid (mean) of the values in the cluster
- Algorithm:
  - Choose  $k$  random centroids (possibly  $k$  observations)
  - Repeat:
    - Assign each observation to nearest centroid
    - Update centroids to be mean for those observations assigned to it.
  - Until there is no, or very little, change in the centroids

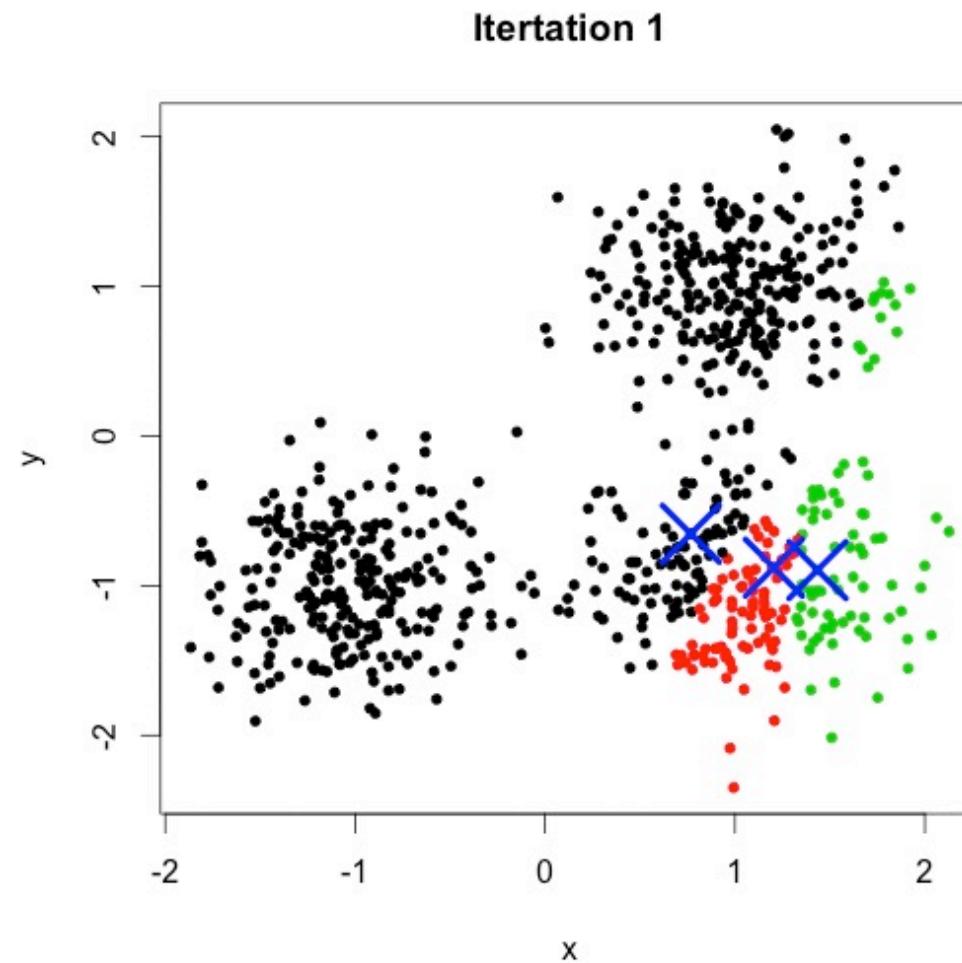
# *k*-means in action

| epcc |



# *k*-means in action

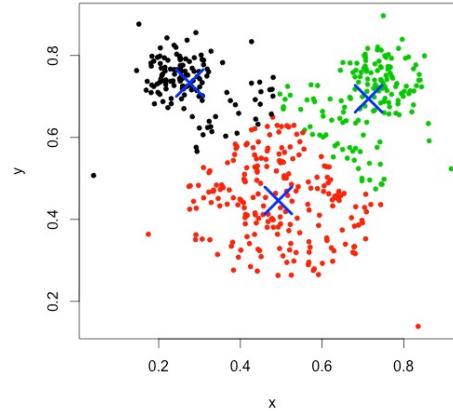
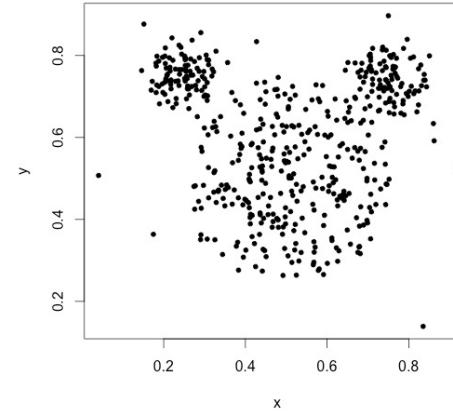
| epcc |



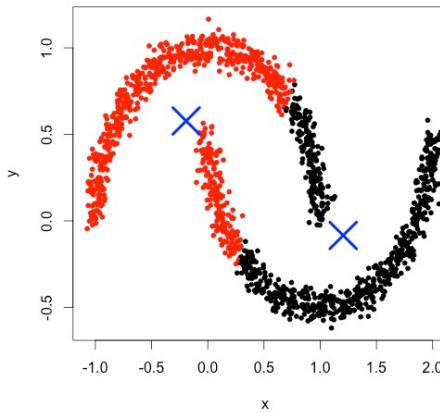
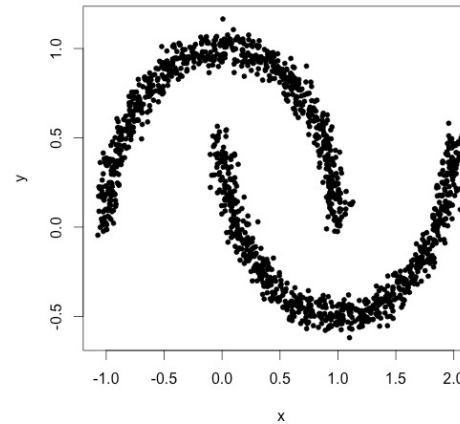
- Must normalise the features so that the distances are not biased to a particular dimension
- Need to think carefully about the features you wish to include as algorithm will give each feature equal weight
  - Unlike certain supervised learning approaches like linear regression where the weight may be zero, or Naïve Bayes where a meaningless feature will have no real impact
- Can be hard to interpret
  - Sometimes the clusters seem meaningless
- Need to choose  $k$ 
  - Sometimes you know there are  $k$  processes generating the data
  - Trial and error
  - Look for 'knee' in plot of cost against  $k$

# Limitations of $k$ -means

- Clusters assumed to be the same size



- Clusters on density not so good





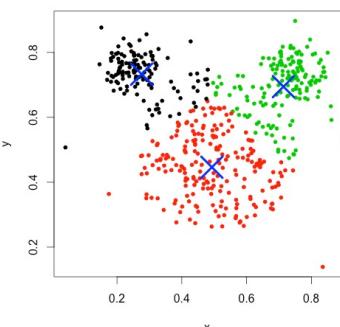
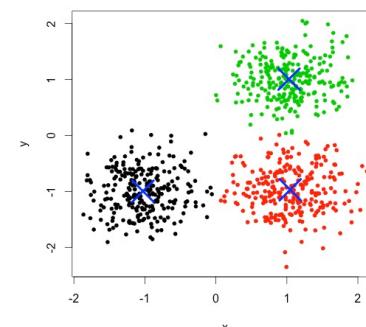
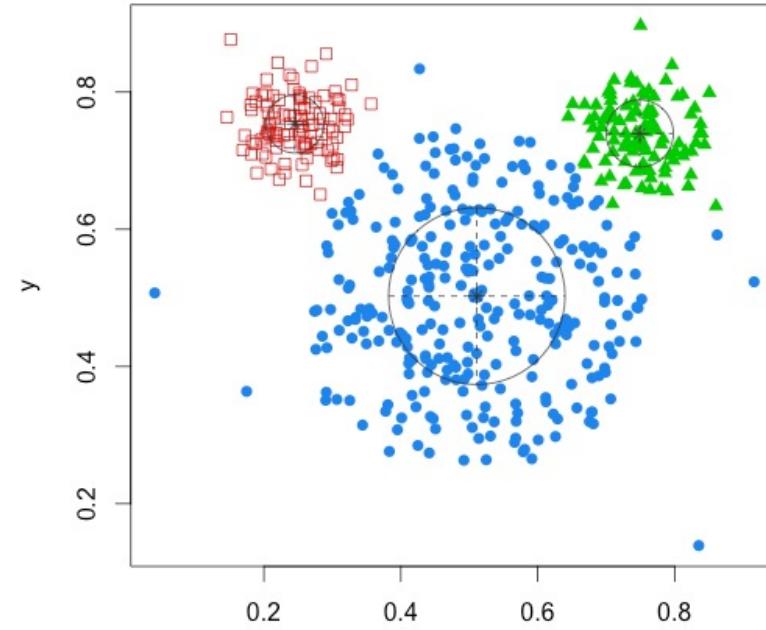
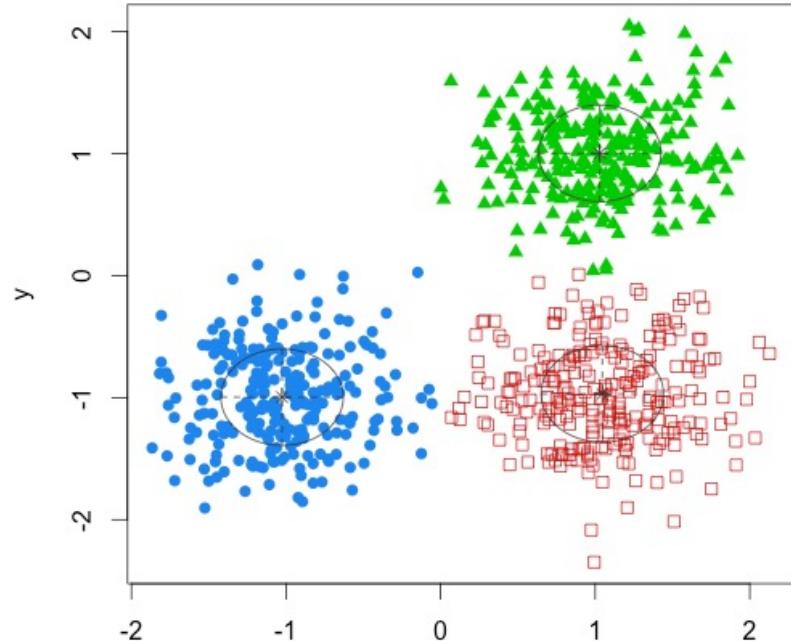
## Some Other Clustering Techniques



| epcc |

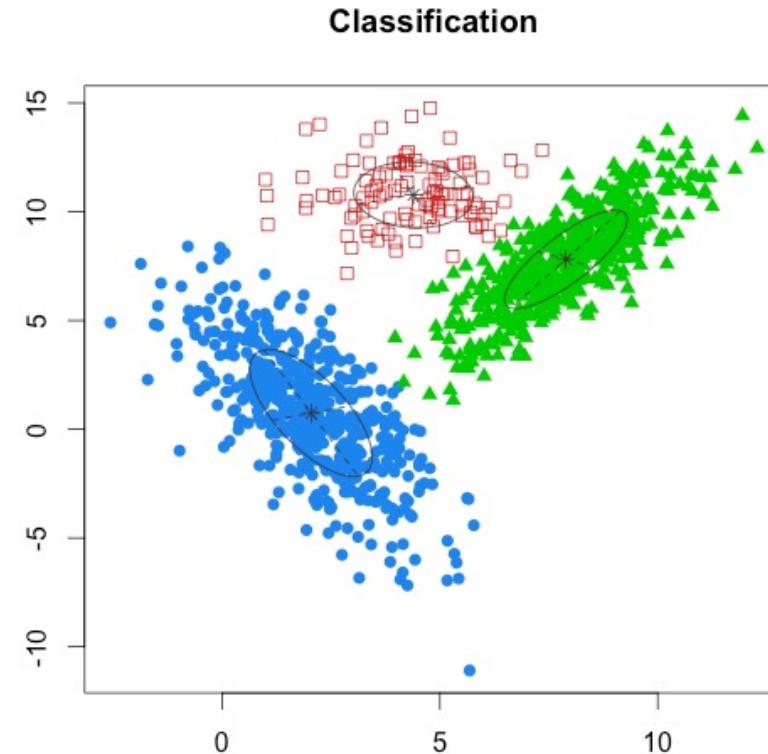
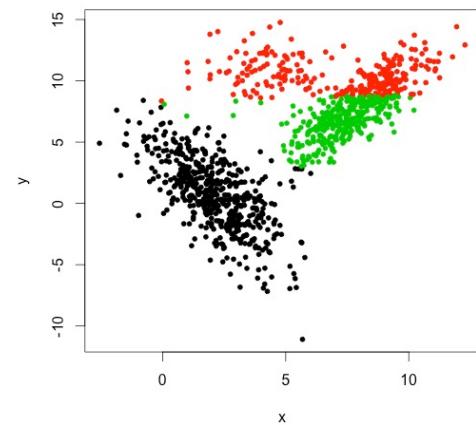
# Distribution-based clustering in practice

epcc



# Distribution-based clustering

- Handles covariance of features
  - No need to normalise data



# Limitations of distribution-based clustering

epcc

- Bad for density-based clusters that don't match distribution model

