

# Content



## 0. Introduction

## 1. Regression

1.1 Multivariate Linear Regression (curve fitting)

1.2 Regularization (Lagrange multiplier)

1.3 Logistic Regression (Fermi-Dirac distribution)

1.4 Support Vector Machine (high-school geometry)

## 2. Dimensionality Reduction/feature extraction

2.1 Principal Component Analysis (order parameters)

2.2 Recommender Systems

2.3 Clustering (phase transition)

# Content



## 3. Neural Networks

3.1 Biological neural networks

3.2 Mathematical representation

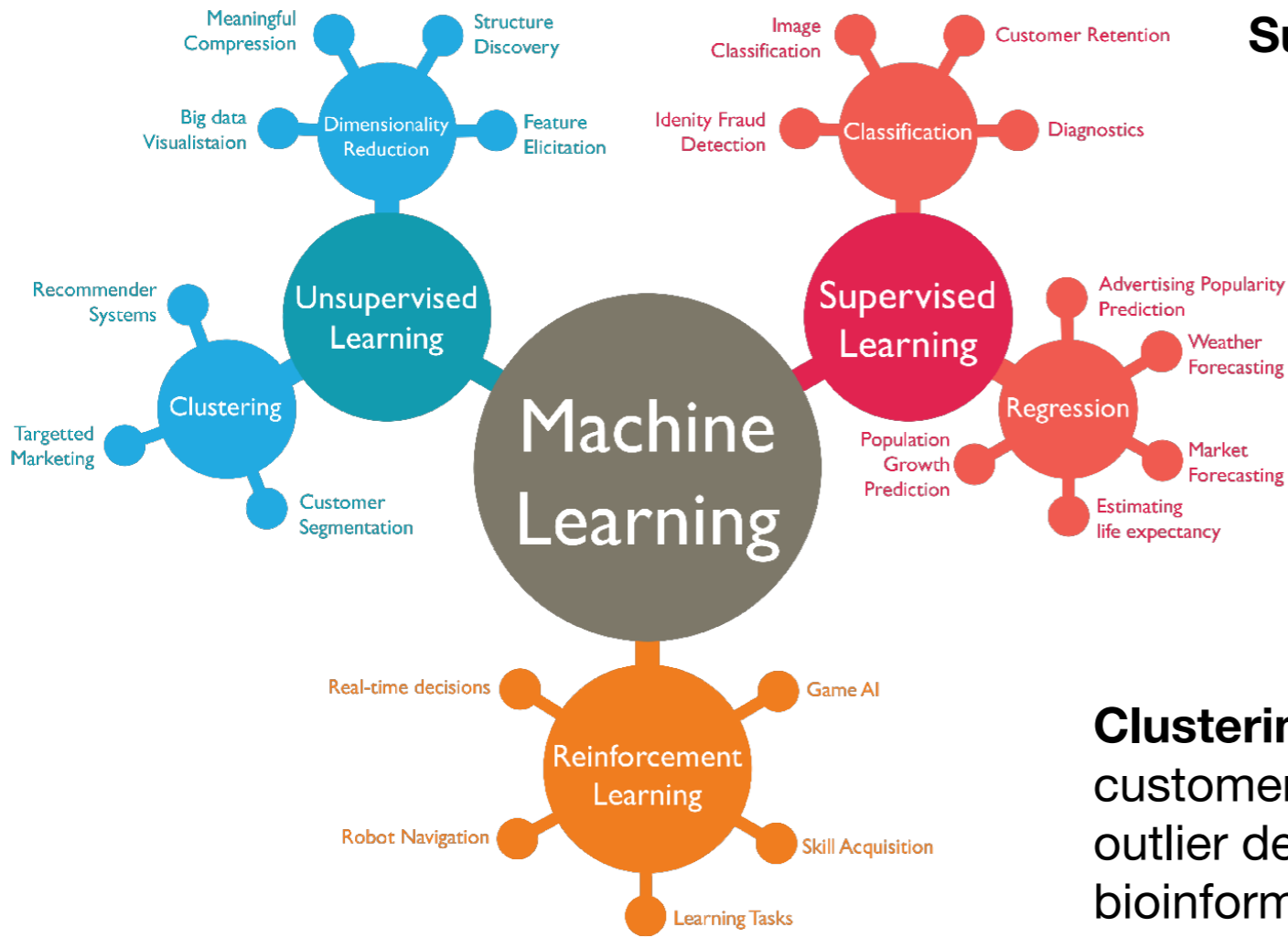
3.3 Factoring biological ingredient

3.4 Feed-forward neural networks

3.5 Learning algorithm

3.6 Universal Approximation Theorem

# AI & Machine Learning Basics



## Supervised Learning: Classification & Regression

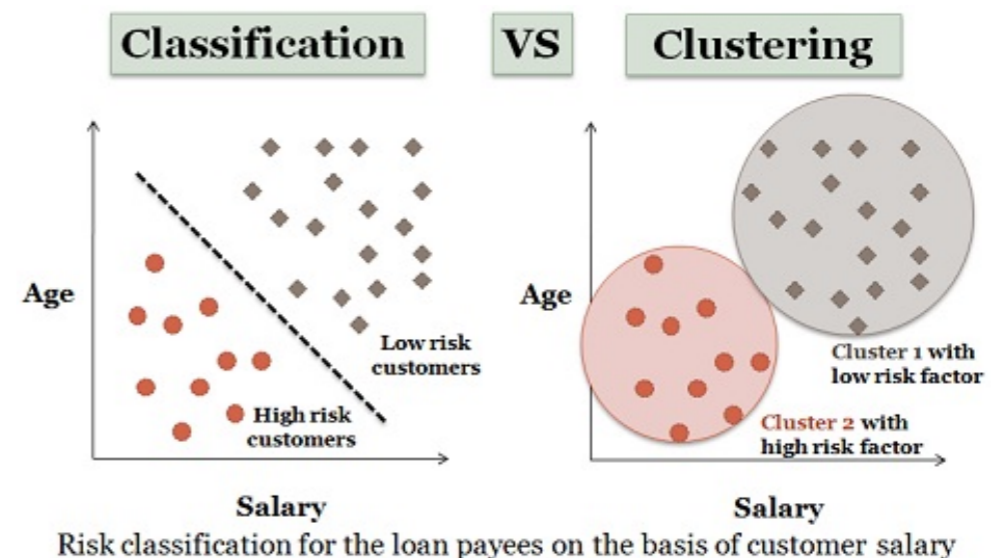
Labeled dataset  
 Input → machine/model → Output  
 Correct outputs are provided by the supervisor

## Unsupervised Learning: only have input data

Unlabelled dataset  
 Find regularities from the input

### Clustering:

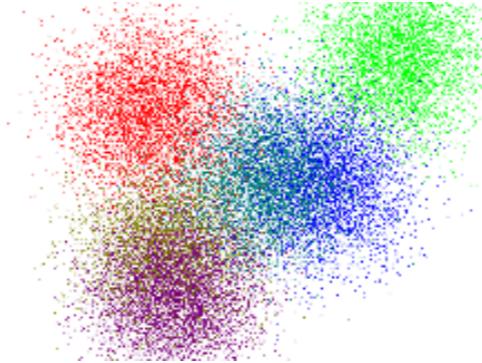
customer segmentation, customer relationship management, outlier detection; Image compression  
 bioinformatics: DNA, RNA, amino acids, Motif, Proteins, sequence alignments



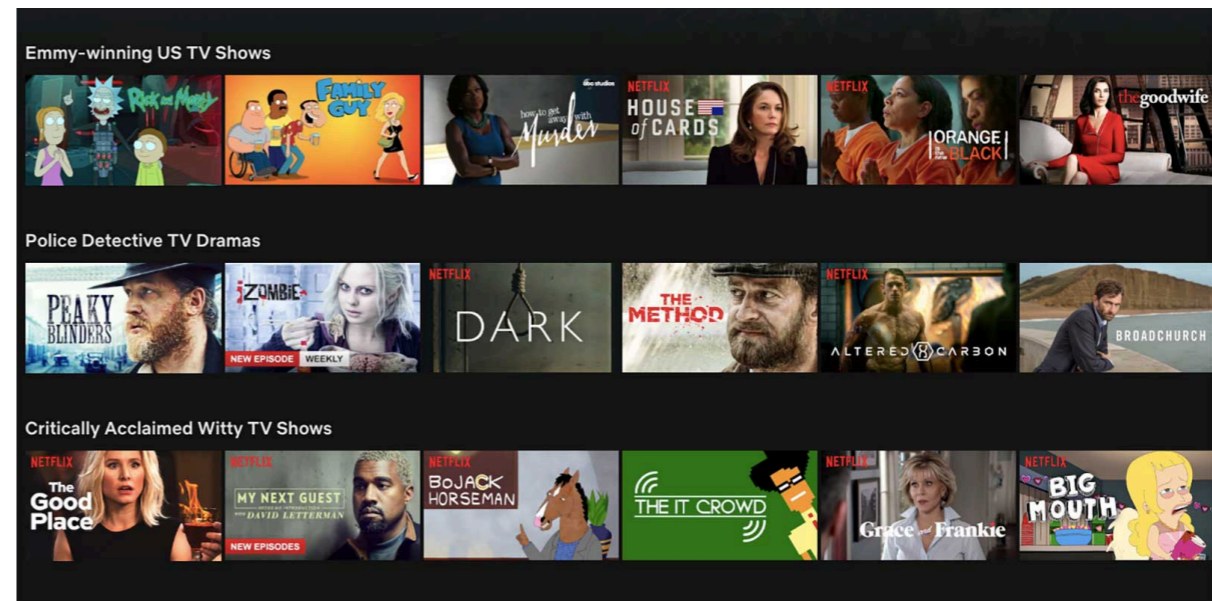
# Clustering

- 📍 Grouping of data points

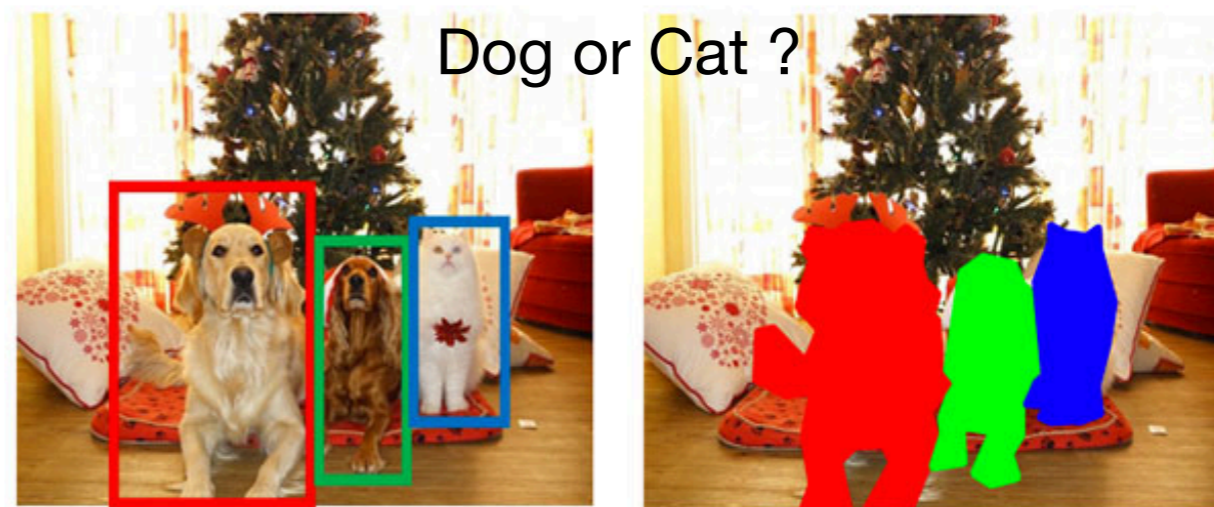
“Clustering” literally means grouping similar things together



- 📍 Recommendation Engines



- 📍 Image Segmentation



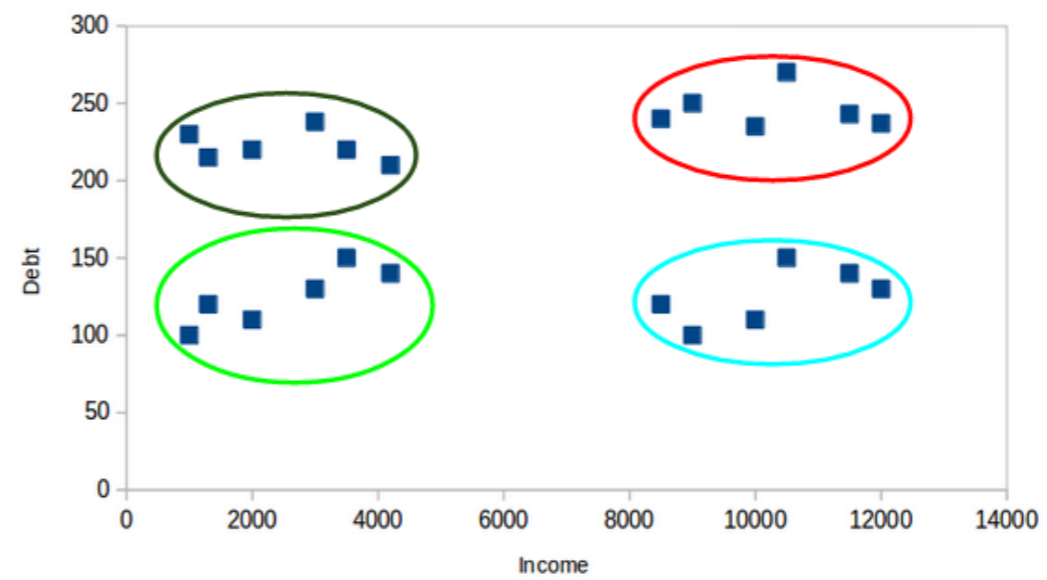
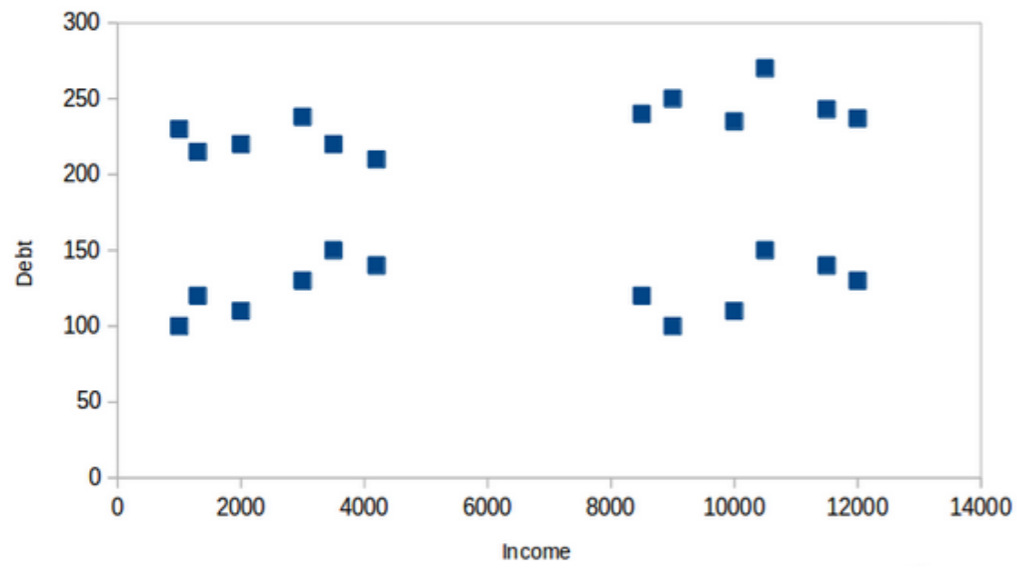
Good references:

<https://www.analyticsvidhya.com/blog/2019/08/comprehensive-guide-k-means-clustering/>

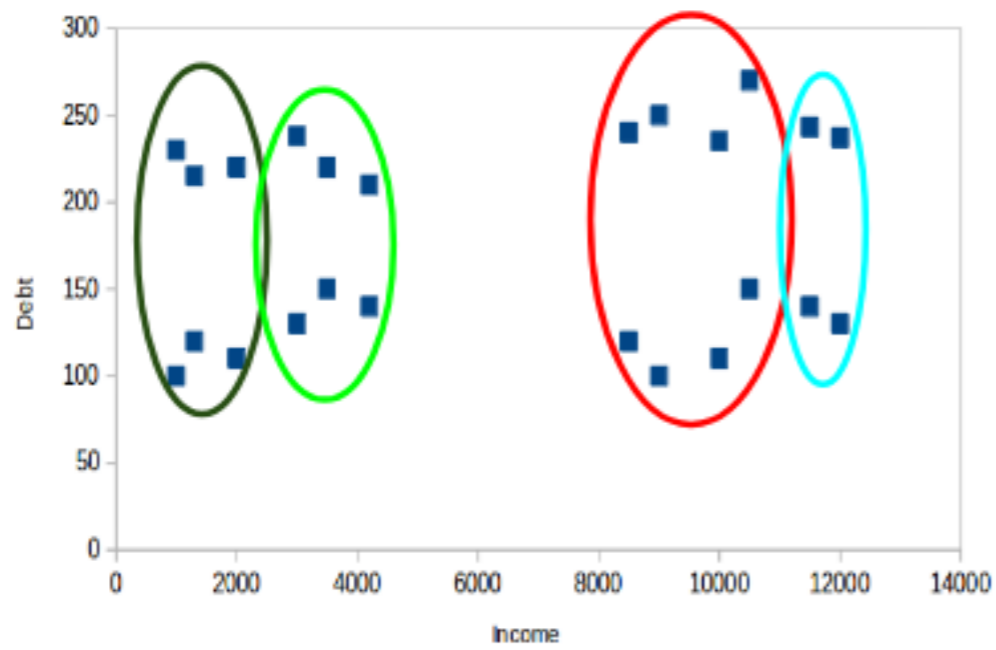
<https://towardsdatascience.com/k-means-clustering-from-a-to-z-f6242a314e9a>

# Clustering

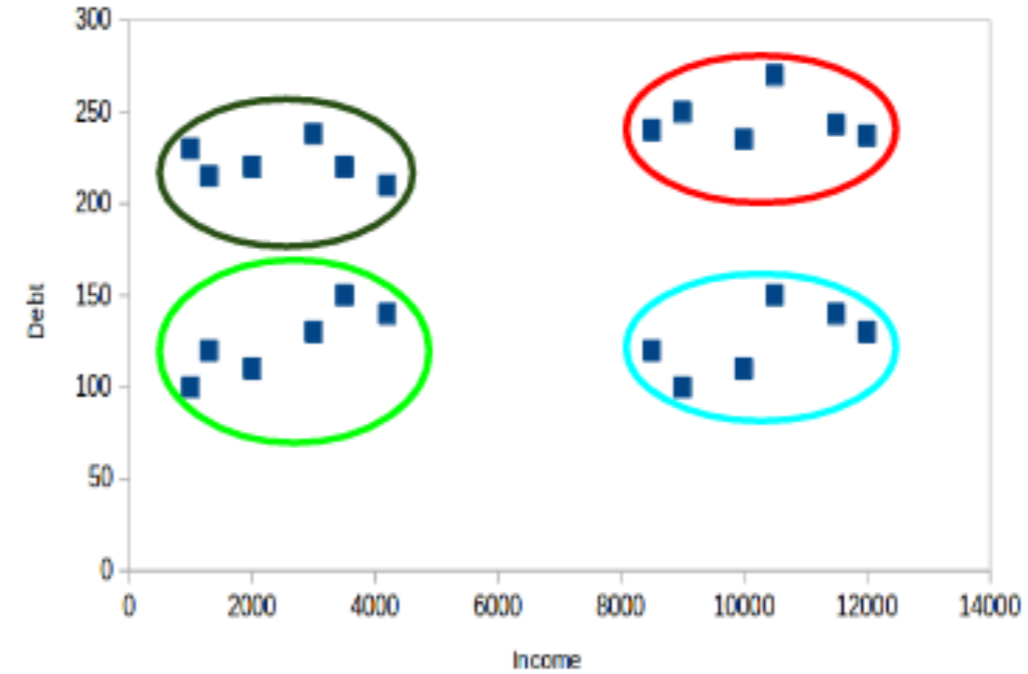
All the data points in a cluster should be similar to one another



The data points from different clusters should be as different as possible



Case - I

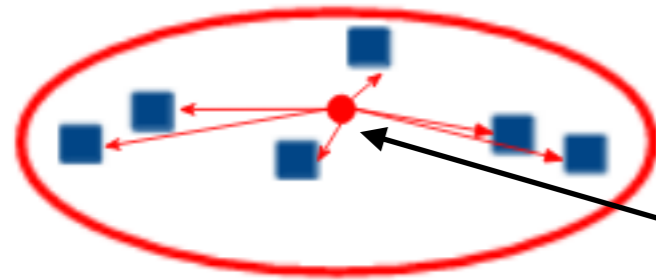


Case - II

# Evaluation Metrics for Clustering

► Inertia: Sum of intracluster distances

The lesser the inertia value, the better the cluster is



Inertia: Sum of intracluster distances

$$\sqrt{\sum_{i=1}^m |\vec{x}_i - \vec{c}_i|^2}$$

Centroid

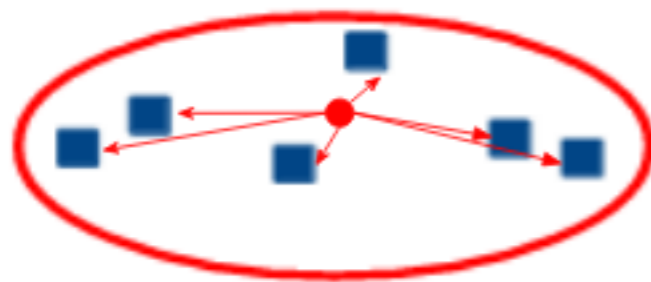
Intra cluster distance

► Dunn Index:

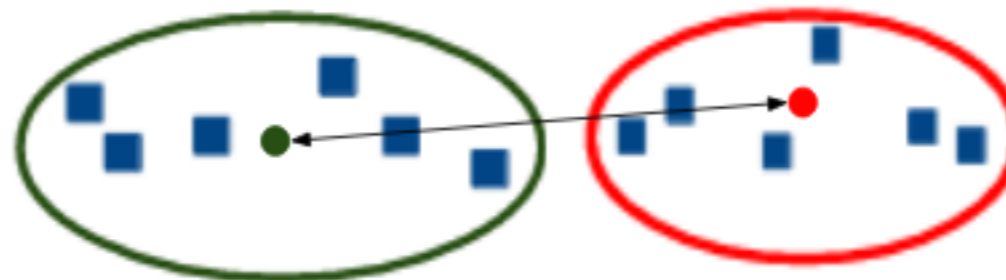
$$\text{Dunn Index} = \frac{\min(\text{Inter cluster distance})}{\max(\text{Intra cluster distance})}$$

Clusters are far apart

Clusters are compact



Intra cluster distance



Inter cluster distance

# K-Means Clustering

Centroid-based or distance-based algorithm, minimise the sum of distances

► Step 1: Choose the number of clusters  $k$

take  $k=2$

► Step 2: Select  $k$  random points from the data as centroids

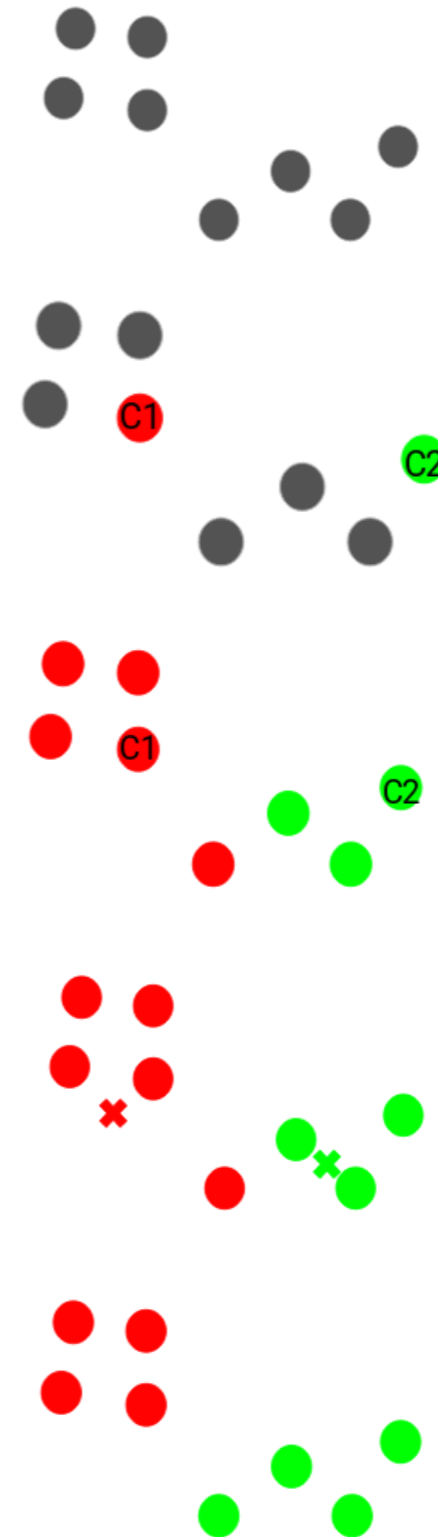
► Step 3: Assign all the points to the closest cluster centroid

► Step 4: Recompute the centroids of newly formed clusters

► Step 5: Repeat steps 3 and 4

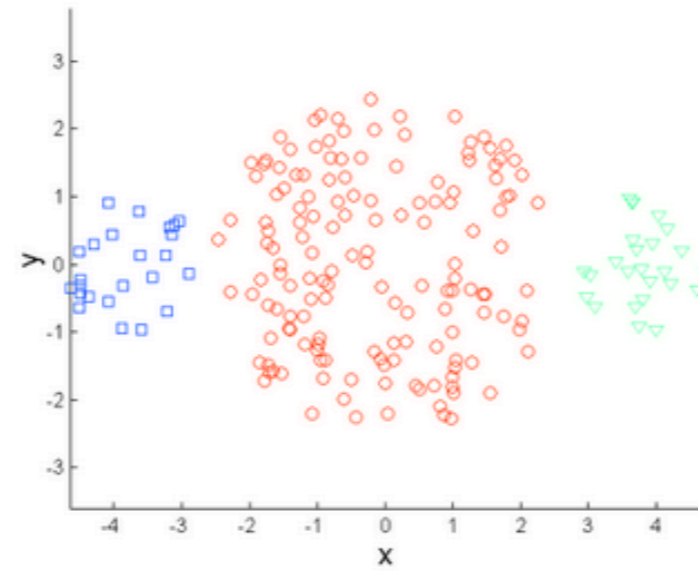
## Stopping Criteria

1. Centroids of newly formed clusters do not change
2. Points remain in the same cluster
3. Maximum number of iterations are reached

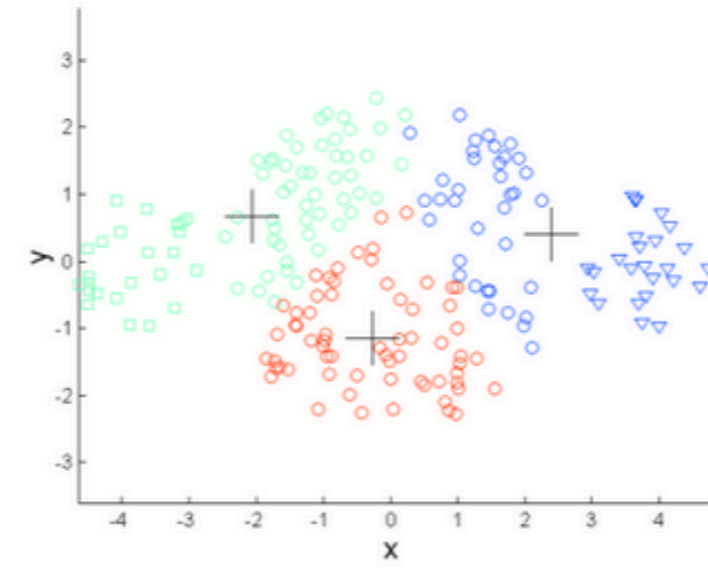


# Challenges with the K-Means Clustering

The size of clusters is different

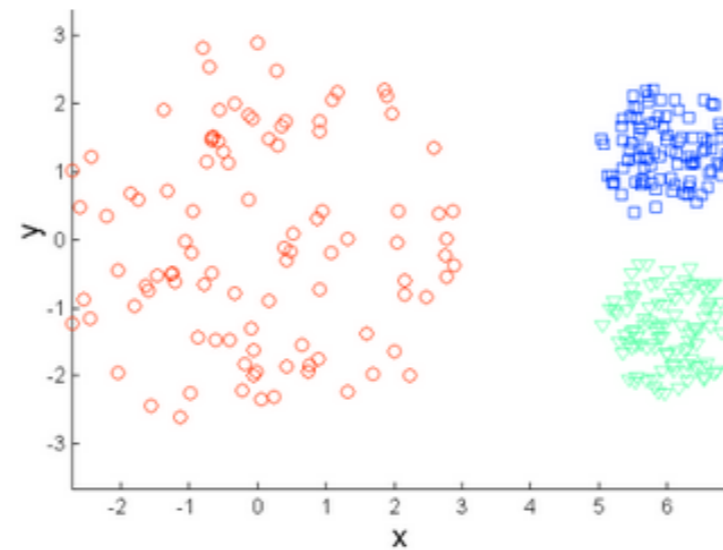


Original Points

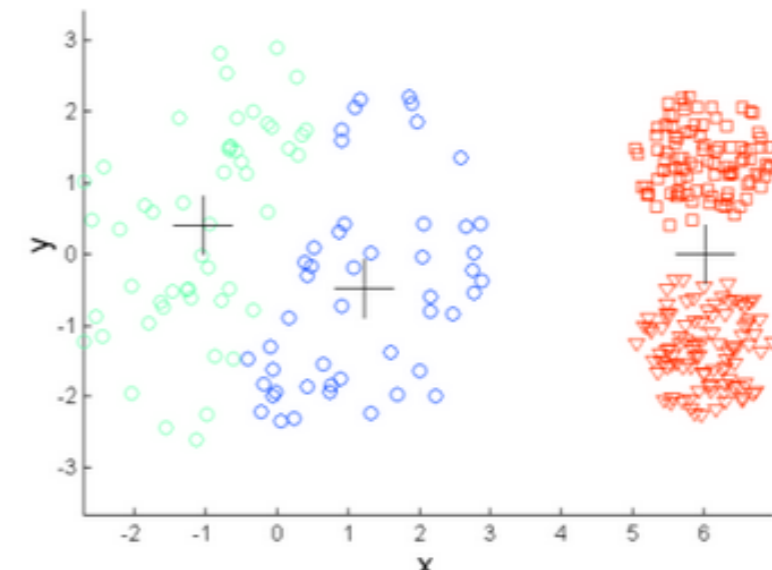


K-means (k = 3)

The densities of the original points are different



Original Points



K-means (k = 3)



# K-Means++ Clustering

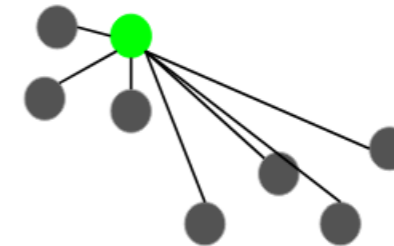
Specifies a procedure to initialise the cluster centres before moving forward with k-means, take  $k=3$

- ▶ Step 1: randomly pick a data point as a cluster centroid

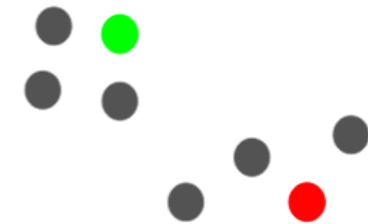
(not all the centroids but one)



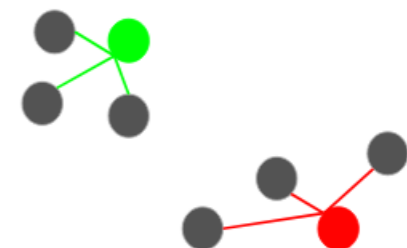
- ▶ Step 2: calculate the distance of each data point with this centroid



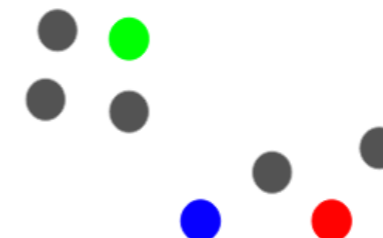
- ▶ Step 3: the next centroid is the one whose distance is the farthest from the current centroid



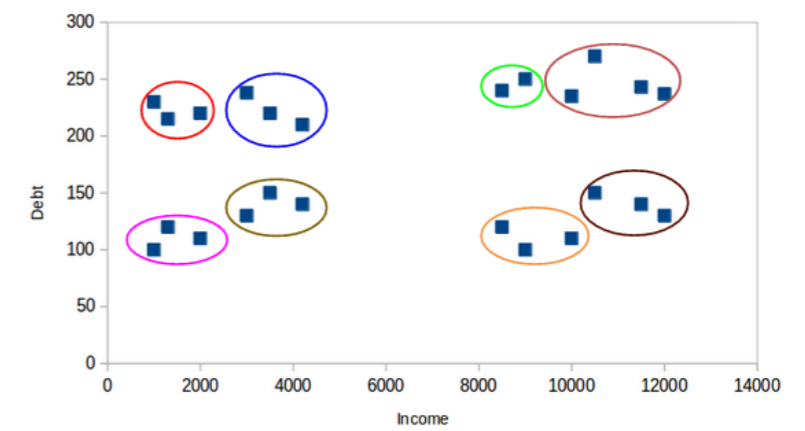
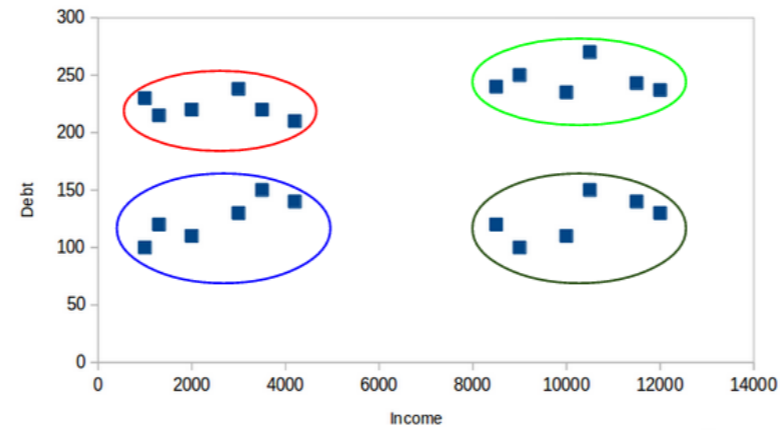
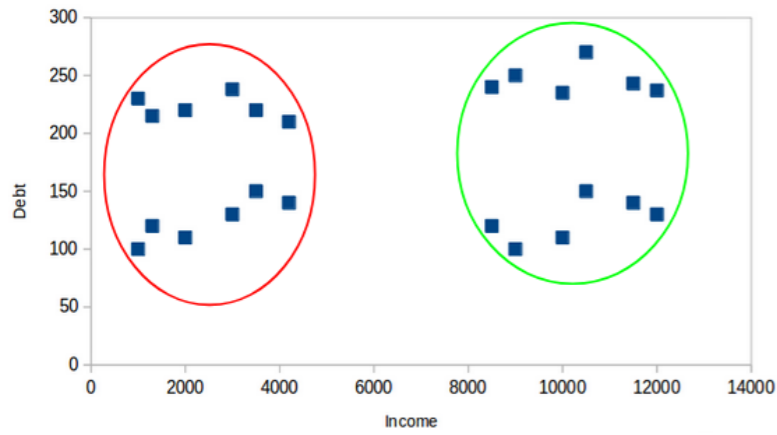
- ▶ Step 4: take the distance of each point from its closest centroid and the point having the largest distance will be selected as the next centroid



- ▶ Step 5: continue with the K-means after initialising the centroids



# How to choose the right number of clusters



 Elbow curve, x-axis represent the number of clusters and y-axis the evaluation metric

