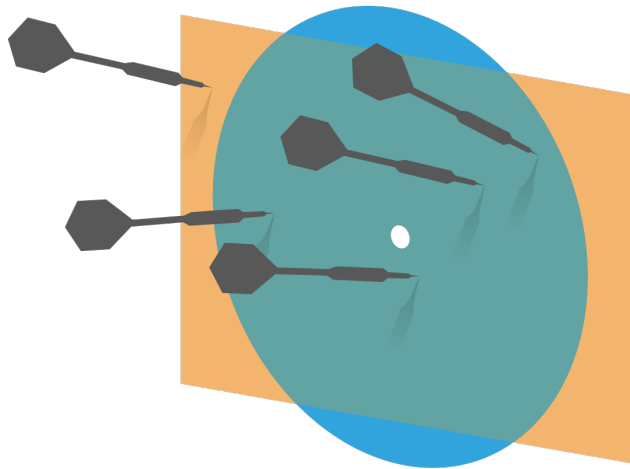


Probability and Statistics for Computer Science



Credit: wikipedia

“Unsupervised learning is arguably more typical of human and animal learning...”--- Kelvin Murphy, former professor at UBC

Last time

- ✱ Linear Regression (II)
- ✱ Nearest Neighbor Regression

$$v^T v = \|v\|^2$$

$$v^T = [v_1 \quad v_2 \quad \dots \quad v_d]$$

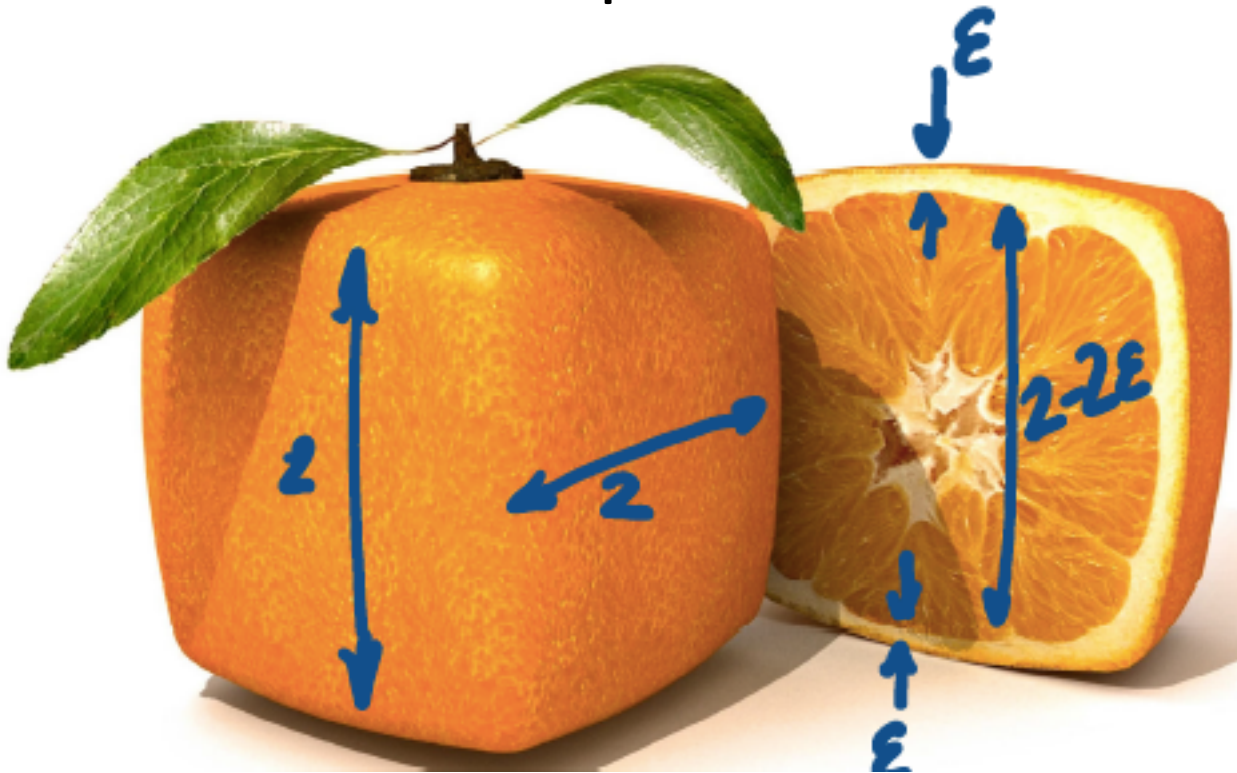
$$v^T v = v_1^2 + v_2^2 + \dots + v_d^2 = \|v\|^2$$

Objectives

- ✱ The curse of dimensionality
- ✱ Multivariate normal distribution
- ✱ Unsupervised learning
- ✱ Clustering (I)
K - means

First let's take a look at a 3D object

Is there more fruit than peel?



First take a look at a 3D object

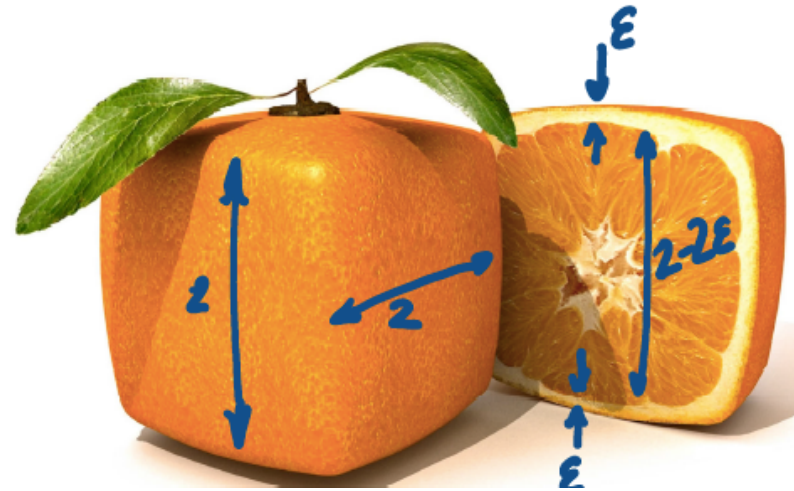
Is there more fruit or more peel?

Total Volume: 2^3

Vol. of fruit: $(2-2\varepsilon)^3$

Vol. of peel: $2^3 - (2-2\varepsilon)^3$

Fraction of peel: $1 - (1-\varepsilon)^3$



If $\varepsilon = 0.05$ fraction of peel ≈ 0.143

What if we have a d -dimensional orange?

Is there always more fruit?

A. YES

☒ B. NO

In arbitrary d-dimension

- ✱ Total amount of orange

$$2^d$$

peel

- ✱ Amount of ~~fruity~~ part

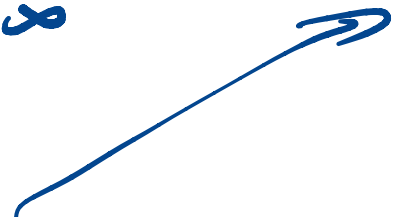
$$2^d - (2 - 2\varepsilon)^d$$

- ✱ Fraction of orange that is peel

$$\rightarrow 1 - (1 - \varepsilon)^d$$

$$\varepsilon \rightarrow \overbrace{0.001}$$

$$\lim_{d \rightarrow \infty} 1 - (1 - \varepsilon)^d = 1$$



$$d = \infty$$

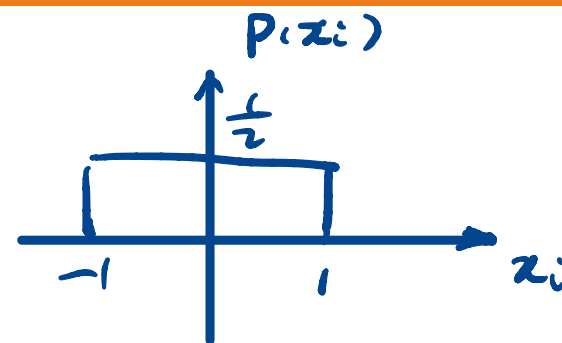
The curse of dimensions

- ✱ If a dataset is uniformly distributed in a high-dimensional cube (or other shape), majority of data is far from the origin.
- ✱ The above can be roughly proved by calculating the expected distance from the origin

The Expected distance from the origin in d-dimensional cube

$$E[\mathbf{x}^T \mathbf{x}] = E\left[\sum_{i=1}^d x_i^2\right] = \sum_{i=1}^d E[x_i^2]$$

$$= \sum_{i=1}^d \int_{cube} x_i^2 P(\mathbf{x}) d\mathbf{x}$$



Assuming the independence of each x_i

$$P(\mathbf{x}) = P(x_1)P(x_2)\dots P(x_d)$$

$$\int_{-\infty}^{+\infty} P(x_i) dx_i = 1$$

The general law of continuous probability density

$$\Rightarrow E[\mathbf{x}^T \mathbf{x}] = \sum_{i=1}^d \int_{-1}^1 x_i^2 P(x_i) dx_i$$

A lot of data is far from the origin.

- ✳ On average, data points are $d/3$ away from the origin (using square of distance)

$$\begin{aligned} E[\mathbf{x}^T \mathbf{x}] &= \sum_{i=1}^d \int_{-1}^1 x_i^2 P(x_i) dx_i \\ &= \sum_{i=1}^d \frac{1}{2} \int_{-1}^1 x_i^2 dx_i \\ &= \frac{d}{3} \end{aligned}$$

$$P(x_i)$$

$$= \begin{cases} \frac{1}{2} & x_i \in (-1, 1) \\ 0 & \text{otherwise} \end{cases}$$

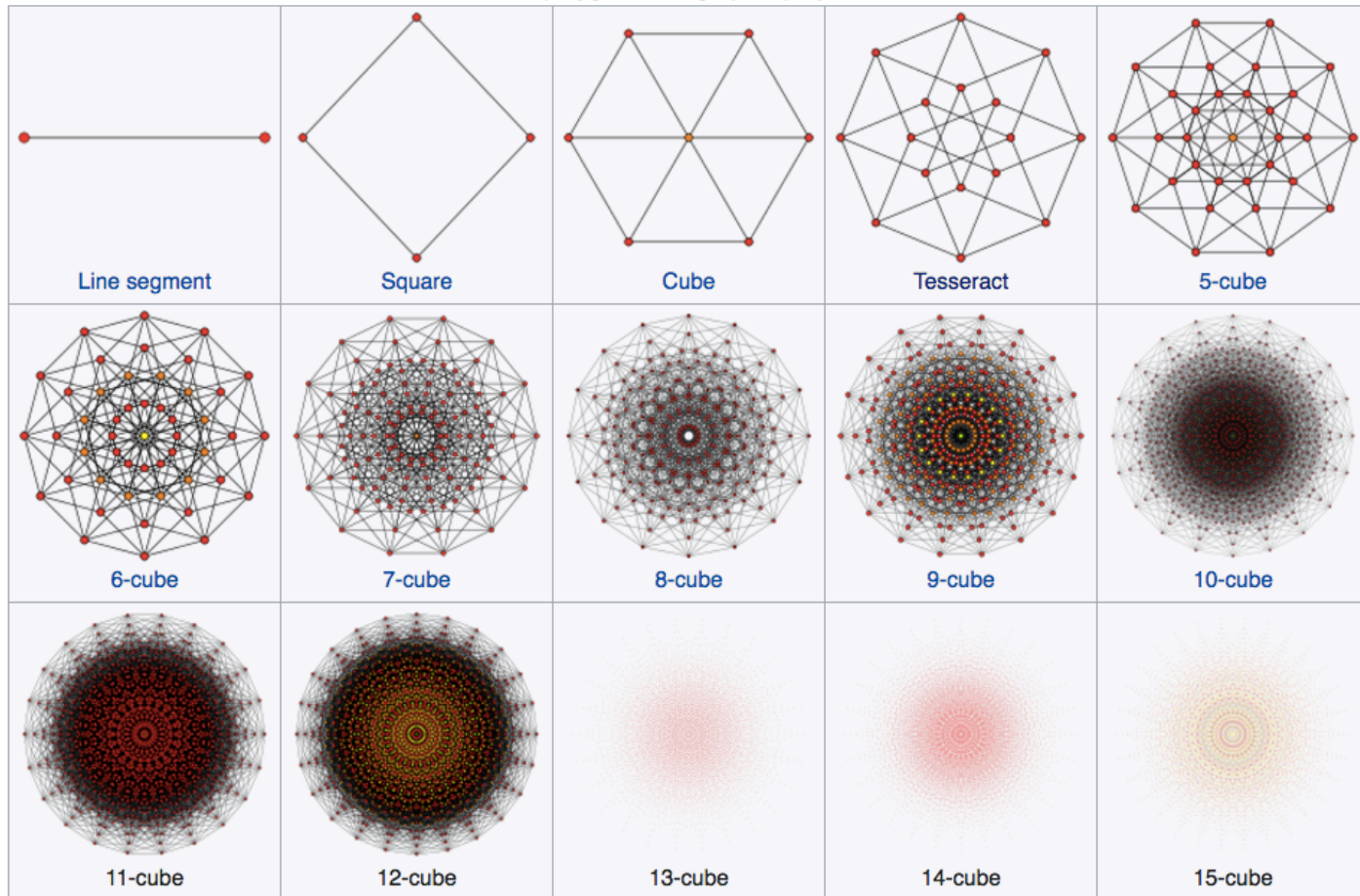
$$\frac{x^3}{3} \Big|_{-1}^1 = \frac{1}{3} - \left(-\frac{1}{3}\right) = \frac{2}{3}$$

What do high-dimensional cubes look like?



What do high-dimensional cubes look like?

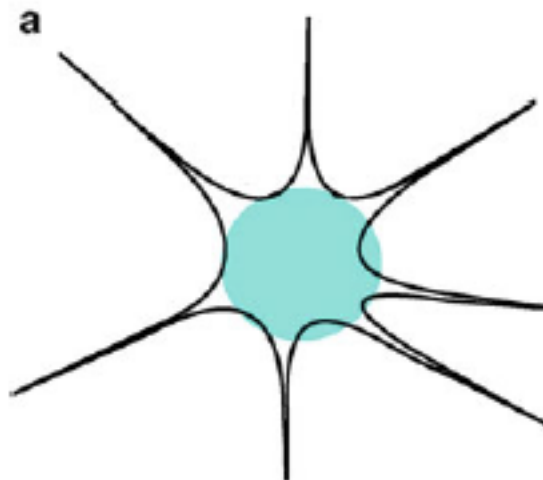
Petrie polygon Orthographic projections



Credit:
Wiki

What does a convex object K in high dimensions look like?

The spikes are outliers in high dimension



A general convex set

Credit: G. Pfander editor, “Sampling theory, a Renaissance”

With this scaling, most of the volume of K is located around the Euclidean sphere of radius \sqrt{n} . Indeed, taking traces on both sides of the second equation in (1.2), we obtain

$$\mathbb{E} \|X\|_2^2 = n.$$

Therefore, by Markov’s inequality, at least 90% of the volume of K is contained in a Euclidean ball of size $O(\sqrt{n})$. Much more powerful concentration results are known—the bulk of K lies very near the sphere of radius \sqrt{n} and the outliers have exponentially small volume. This is the content of the two major results in high-dimensional convex geometry, which we summarize in the following theorem.

Distance between points grows with increasing dimensions

$$\begin{aligned} E[d(u, v)^2] &= E[(u - v)^T (u - v)] \\ &= E[\underbrace{u^T u}] + E[\underbrace{v^T v}] - 2E[\cancel{u^T v}] \end{aligned}$$

$$= \frac{d}{3} + \frac{d}{3} - 0$$

u, v are
orthogonal

$$= \frac{2}{3}d$$

$$u^T v = 0$$

High dimensional histogram of a data set is unhelpful



- ✱ Most bins will be empty
- ✱ Some bins will have single data
- ✱ Very few will have more than one data point

Dealing with high dimensional data

- ✱ Collect as much data as possible
- ✱ Cluster data into blobs/cluster
- ✱ Fit each blob with simple probability model

Multivariate normal distribution

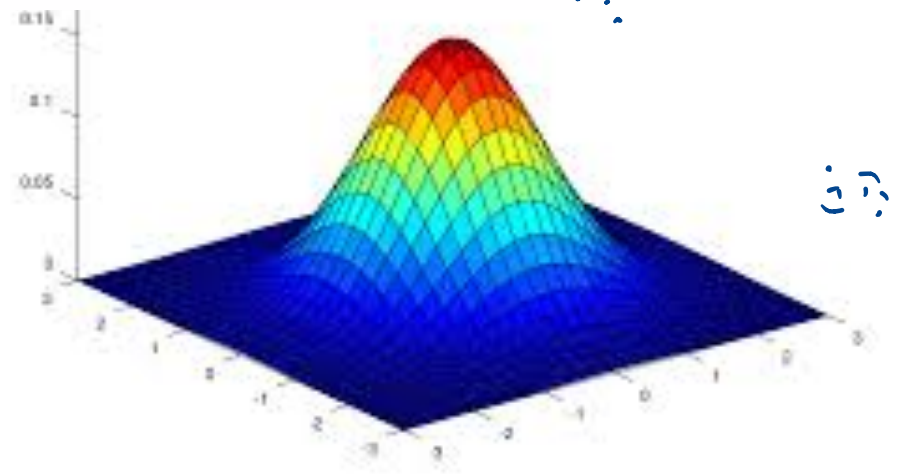
- ✱ Extension of the normal distribution to multiple dimensions $p(x) = \frac{e^{-\frac{(x-\mu)^2}{2\sigma^2}}}{\sqrt{2\pi}\sigma}$

- ✱ Bivariate normal distribution looks like this:

$$f(x, y) = \frac{1}{2\pi\sigma_X\sigma_Y\sqrt{1-\rho^2}} e^{-\frac{1}{2(1-\rho^2)} \left[\left(\frac{x-\mu_X}{\sigma_X} \right)^2 - 2\rho \left(\frac{x-\mu_X}{\sigma_X} \right) \left(\frac{y-\mu_Y}{\sigma_Y} \right) + \left(\frac{y-\mu_Y}{\sigma_Y} \right)^2 \right]}$$

$$-1 < \rho < 1 \quad \rho \neq \pm 1$$

$$\rho \rightarrow \text{corr}(x, y)$$



Multivariate normal probability density

- ✱ A multivariate normal random vector \mathbf{X} of dimension d has this pdf: *Mahalanobis dist.*

$$P(\mathbf{x}) = \frac{1}{\sqrt{(2\pi)^d |\Sigma|}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu})\right)$$

where

$$\boldsymbol{\mu} = E[\mathbf{x}]$$

$$\Sigma = E[(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^T]$$

Multivariate MLE

- Given a d-dimensional data set ($\{x\}$) we can fit a multivariate normal model using MLE $x_i, \text{ iid}$

$$P(x_i | \theta) = \frac{1}{\sqrt{(2\pi)^d |\Sigma|}} \exp\left(-\frac{1}{2}(x_i - \mu)^T \Sigma^{-1}(x_i - \mu)\right)$$

$$x_i \sim d \times 1$$

$$\theta = \{\mu, \Sigma\}$$

$$L(\theta) = \prod_i^N P(x_i | \theta)$$

$$\hat{\mu} = \frac{\sum x_i}{N} \sim d \times 1$$

$$\hat{\theta} = \{\hat{\mu}, \hat{\Sigma}\}$$

$$\hat{\theta} = \arg \max_{\theta} L(\theta)$$

$\hat{\Sigma}$ is the covariance matrix of $\{x\} \sim d \times d$

Unsupervised learning

- ✱ **Unsupervised learning** means knowledge discovery from the feature vectors **without labels**.
- ✱ Unsupervised learning may include:
 - ✱ Discovering **latent factors** *i.e. eigenvectors of covmat.*
 - ✱ Discovering **clusters**
 - ✱ Discovering **graph structure**
 - ✱ Matrix completion

x	y
0	1
0	1
0	1
?	1

Q. Is this true?

✱ **Principal Component Analysis** is an unsupervised learning method.

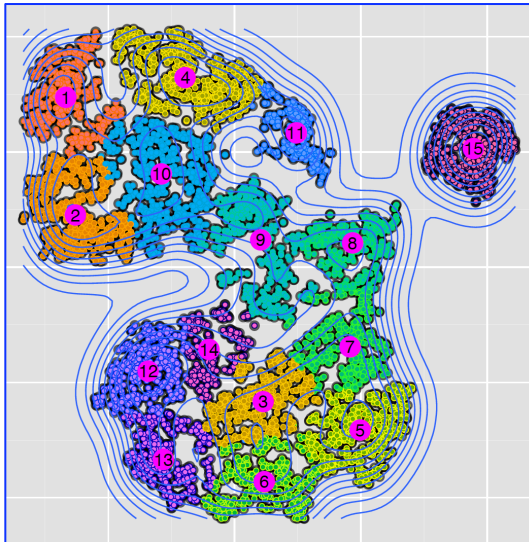
- A. TRUE
- B. FALSE

Dimension Reduction is unsupervised learning

- ✱ For example in **Principal Component Analysis**, no labels are assumed about the data.
- ✱ PCA discovers the latent factors--- the important eigenvectors of the covariance matrix

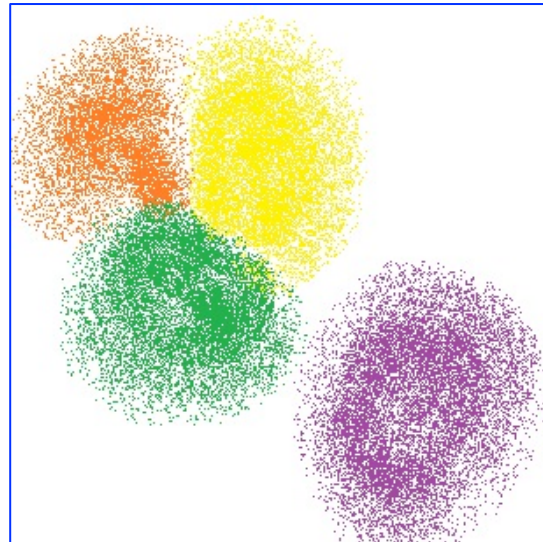
The family of unsupervised learning

Dimension reduction



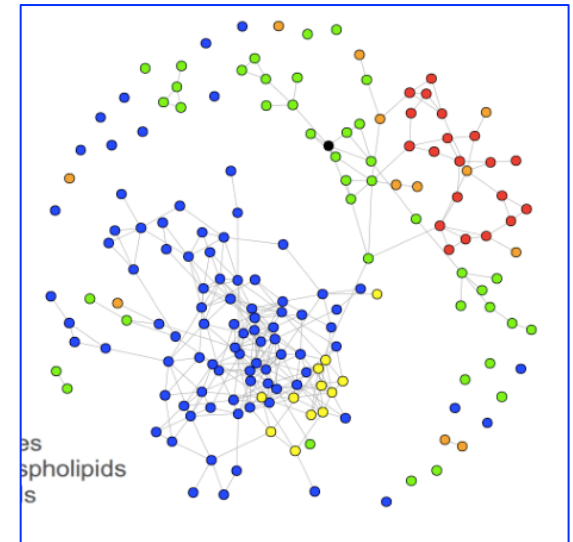
t-SNE

Clustering



K-means

Graph structure



Gaussian Graph model

• • • • •

Clustering as an unsupervised learning method

- ✱ Clustering identifies specific structure called **clusters**.
- ✱ In clustering **data is not labeled**. By identifying clusters, the method assigns cluster membership labels to data.
- ✱ A cluster is formed so that
 - ✱ Items within a cluster are “close” to each other
 - ✱ Items in different clusters are “far” from each other
 - ✱ Distance metric is important in clustering

Types of clustering method

- ✱ By input type:

- ✱ **Similarity based clustering:** input is $N \times N$ similarity/distance matrix

- ✱ **Feature based clustering:** input is $N \times D$ feature matrix

- ✱ By output type:

- ✱ **Hierarchical clustering**

- ✱ Top-down (divisive)

- ✱ Bottom-up (agglomerative)

- ✱ **Flat clustering:**

- ✱ Mixture models, K-means clustering, Spectral clustering...

N
0.5
 N
Iris data
 N [Sepal L Sepal W ...]

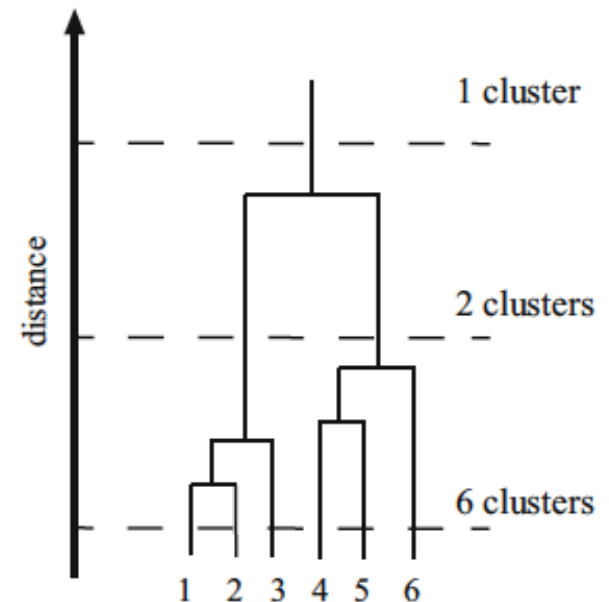
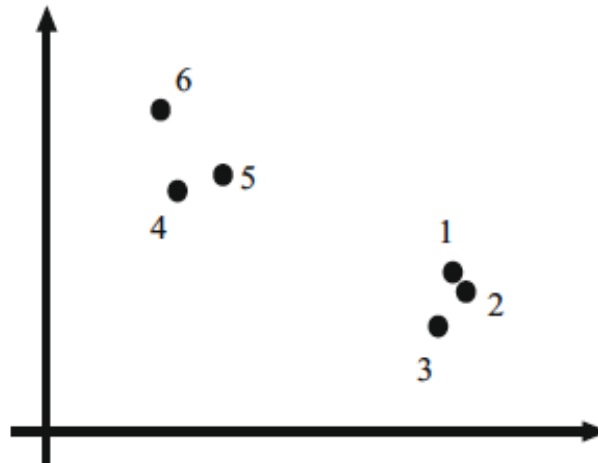
Hierarchical Clustering (I)

✧ Divisive clustering

- ✧ Treat the whole dataset as a single cluster
- ✧ Then split the data set recursively until you get a satisfactory clustering

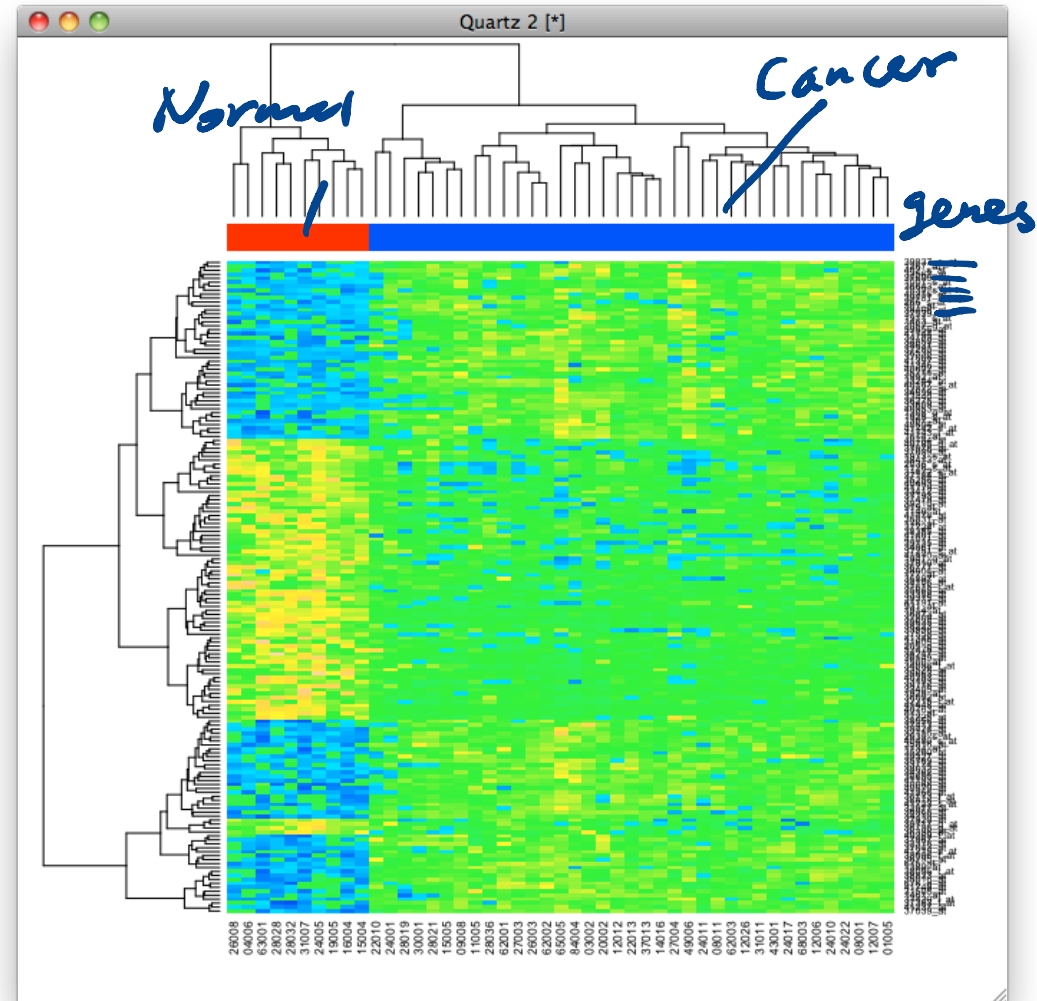
Hierarchical Clustering (II)

- ✱ Agglomerative clustering
 - ✱ Treat each data item as its own cluster
 - ✱ Then merge clusters until you get a satisfactory clustering
 - ✱ A “dendrogram” is created



Hierarchical Clustering example

- ✱ Agglomerative clustering of matrix of gene-tissue pairs of human samples.
- ✱ Columns are tissues; rows are genes
- ✱ Clustering is done for both directions

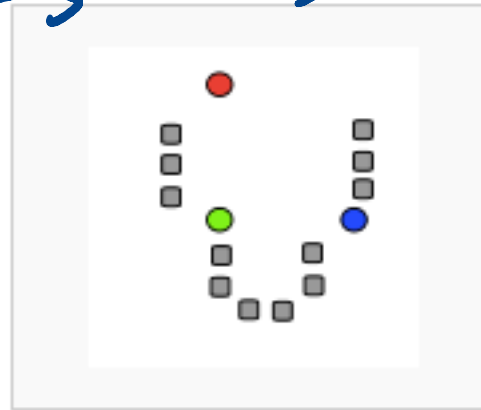


K-means clustering

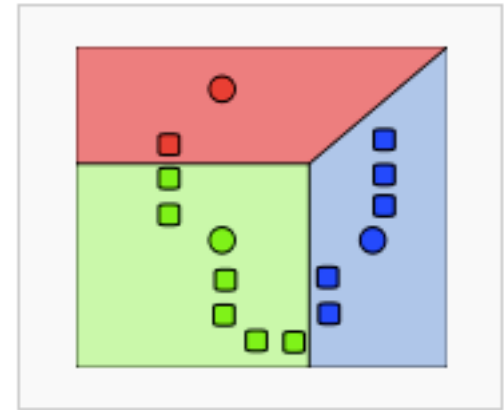
- * Pick a value k as the number of clusters
- * Select k random cluster centers
- * Iterate until convergence:
 - * Assign each data to the nearest center
 - * Update the center within the cluster

$k=3$

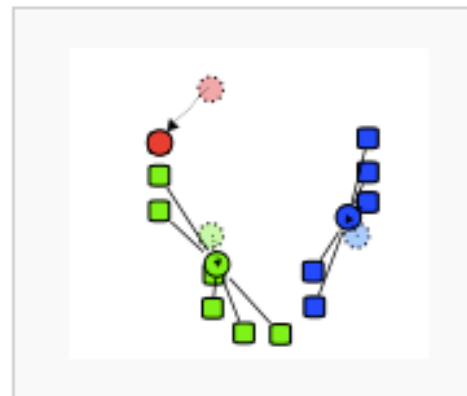
seed



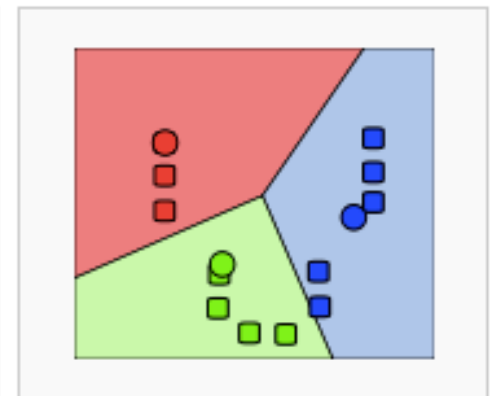
(1)



(2)



(3) Source:wikipedia



(4)

Q. What are the values of c_1 and c_2 ?

- Given a dataset $\{0, 2, 4, 6, 24, 26\}$, initialize the k-means clustering algorithm with 2 cluster centers $c_1 = 3$ and $c_2 = 4$. What are the values of c_1 and c_2 after **one** iteration of k-means?

3	4
0, 2	4, 6, 24, 26
$c_1' = \frac{0+2}{2} = 1$	$c_2' = \frac{4+6+24+26}{4} = 15$

Q. What are the values of c_1 and c_2 ?

- Given a dataset $\{0, 2, 4, 6, 24, 26\}$, initialize the k-means clustering algorithm with 2 cluster centers $c_1 = 3$ and $c_2 = 4$. What are the values of c_1 and c_2 after **two** iterations of k-means?

	1	15
	0 2 4 6	24 26
c_1^*	$\frac{0 + 2 + 4 + 6}{4} = 3$	$\frac{24 + 26}{2} = 25$

What does k-means do mathematically?

- ✱ It's a minimization of a cost function

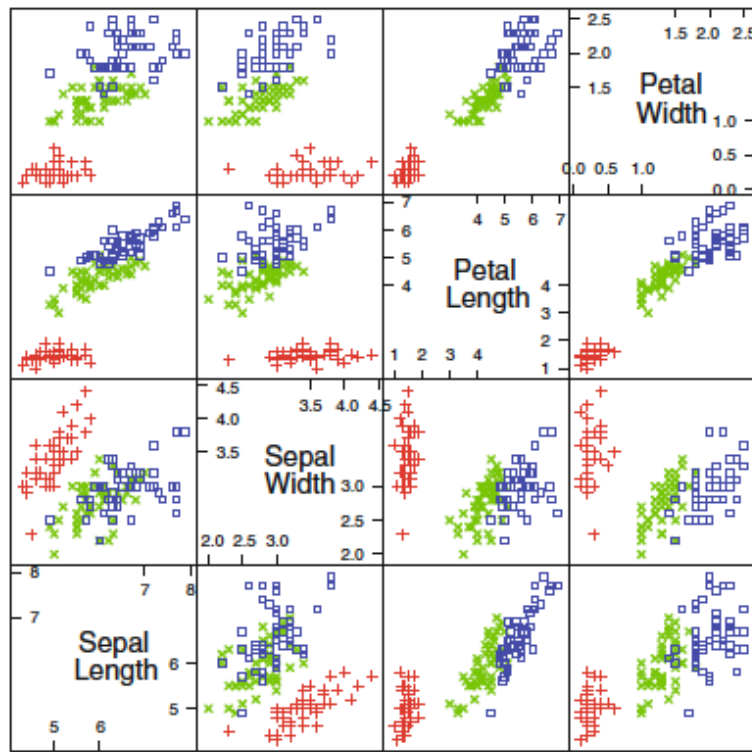
$$\begin{aligned}\phi(\delta, \mathbf{c}) &= \sum_{i,j} \delta_{i,j} [(\mathbf{x}_i - \mathbf{c}_j)^T (\mathbf{x}_i - \mathbf{c}_j)] \\ &= \sum_i^N \sum_j^k \delta_{i,j} \|\mathbf{x}_i - \mathbf{c}_j\|^2 \quad \delta_{i,j} = \begin{cases} 1 & \text{if } \mathbf{x}_i \in \text{cluster } j \\ 0 & \text{otherwise} \end{cases}\end{aligned}$$

- ✱ Cost is defined by the sum of squared distances of each data point from its cluster center

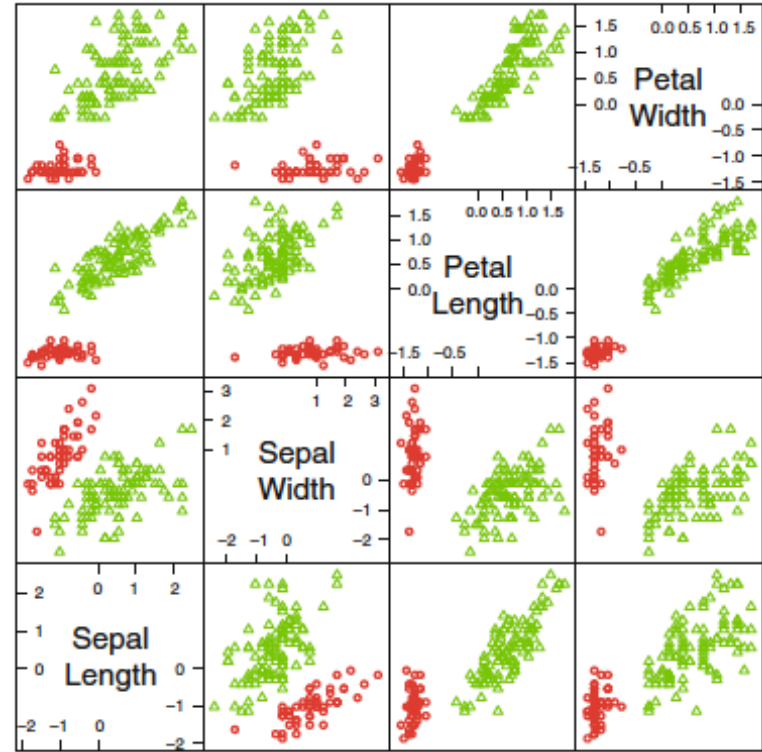
K-means clustering example: Iris

$k=2$

True labels

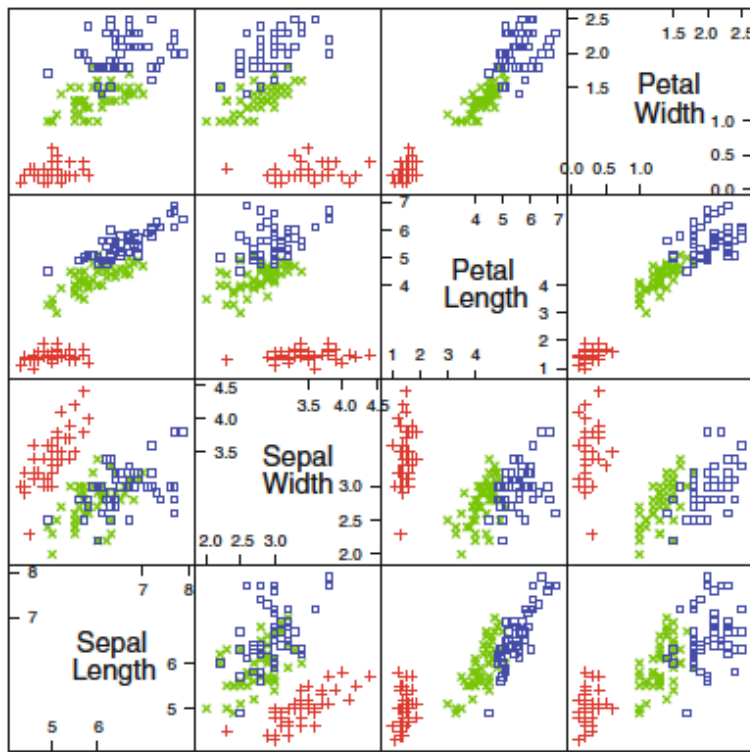


2 clusters

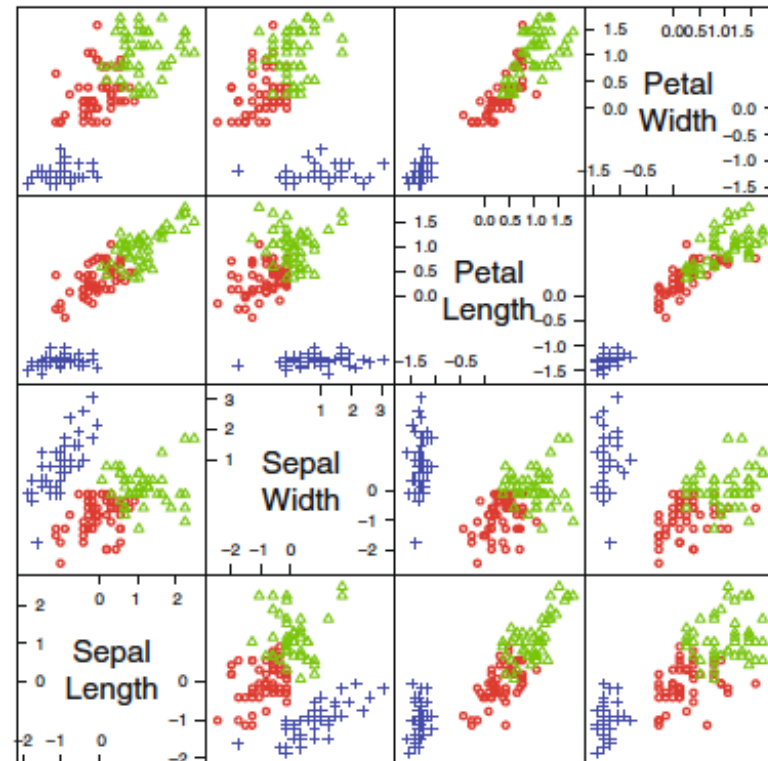


K-means clustering example: Iris

True labels

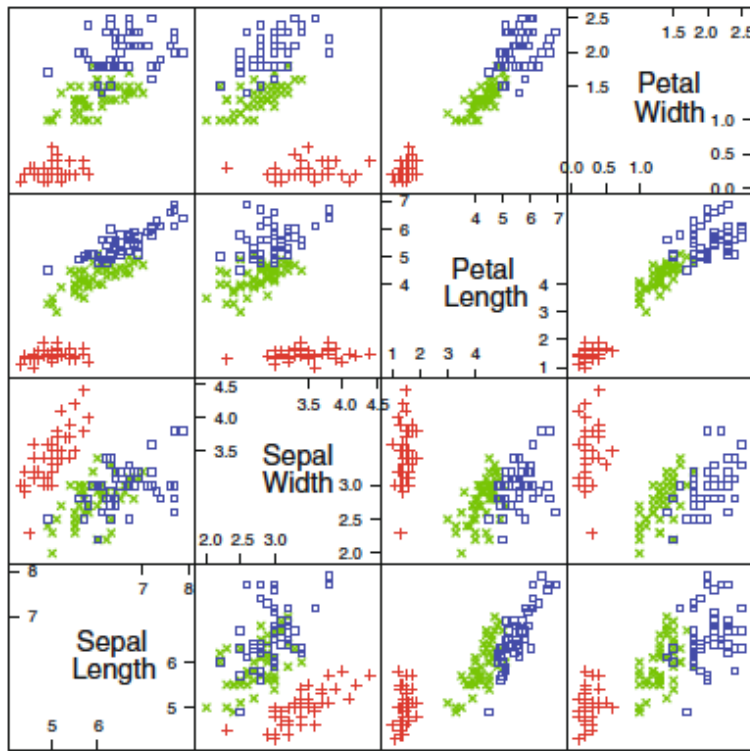


3 clusters

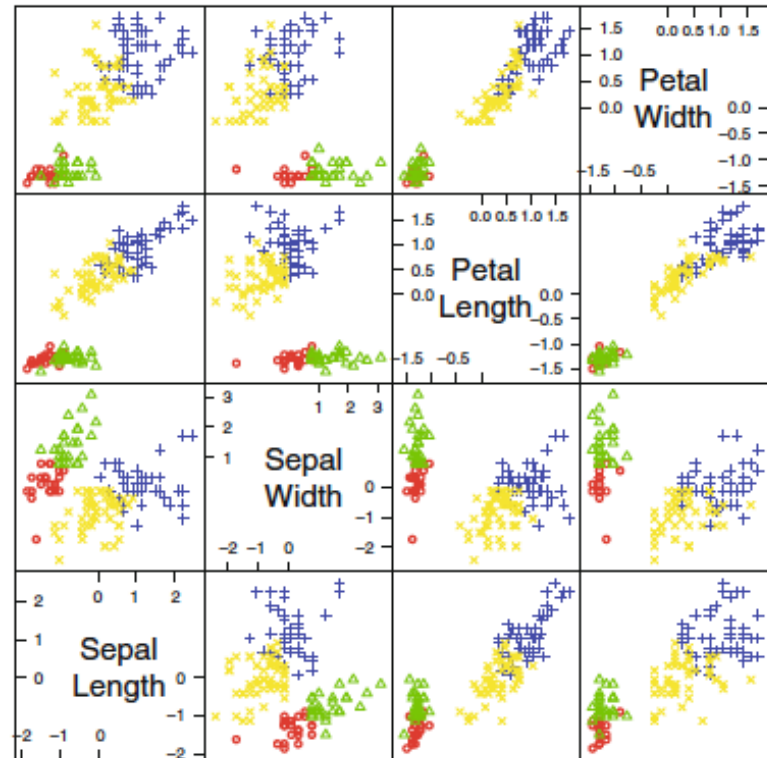


K-means clustering example: Iris

True labels



4 clusters



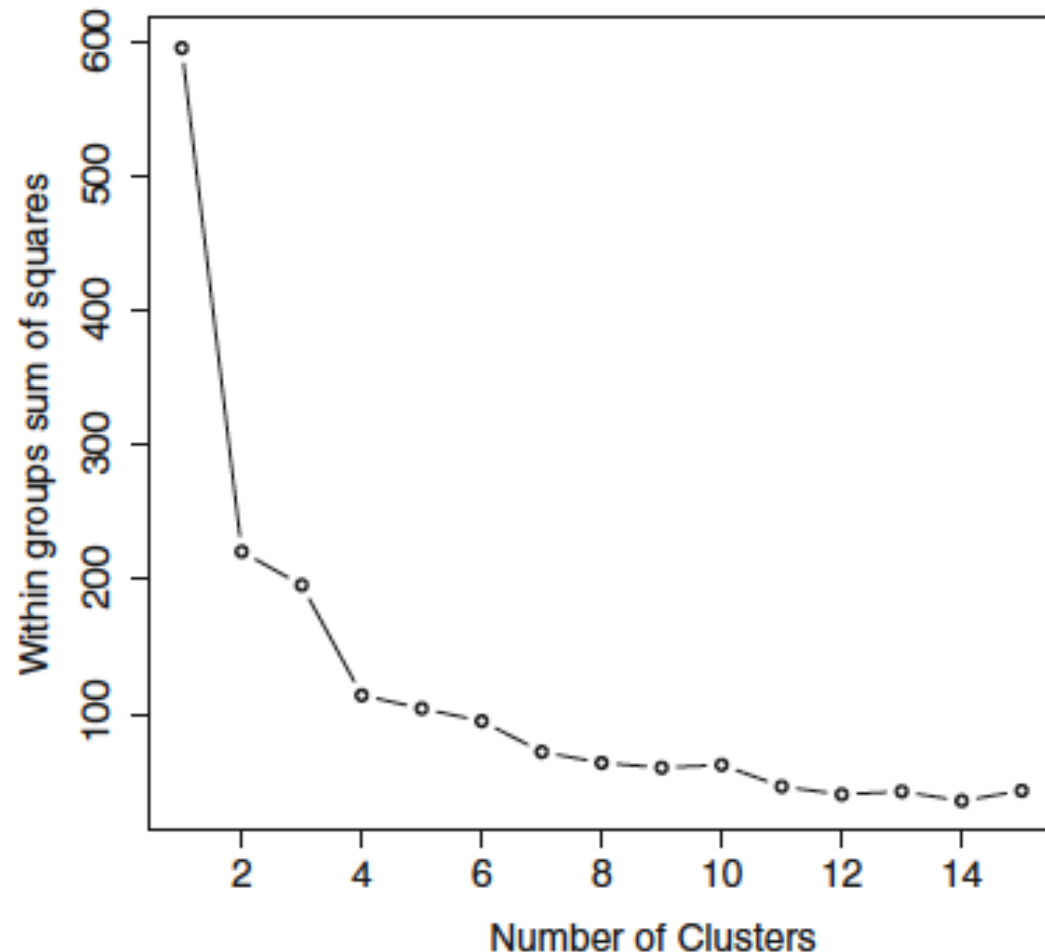
How to choose the value of k ?

- ✱ Sometimes we have the knowledge from the data set.
- ✱ Sometimes we have some other natural way to choose k .
- ✱ Otherwise given the cost function, we may perform clustering for many k values and choose k from the knee of the cost function empirically.

Choose k from the cost function curve

Which is best?
Still depends on
the application

Usually we want
fewer clusters.



Some variants of k-means clustering

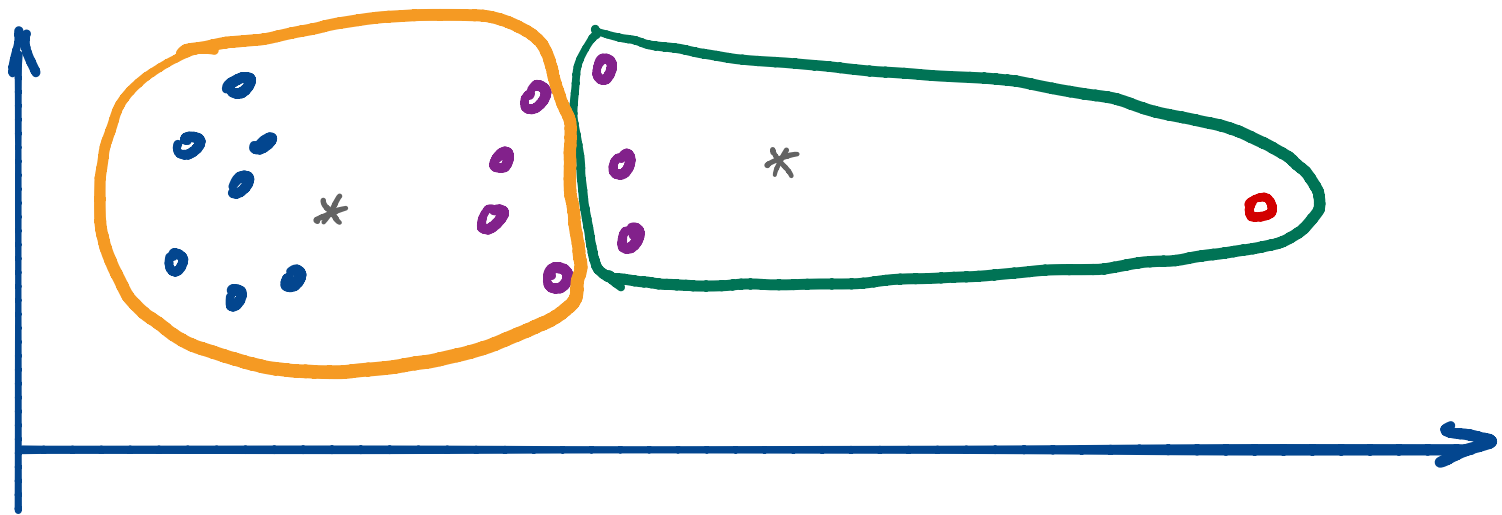
- ✱ Soft assignment allows some data items to belong to multiple clusters with weights associated with each cluster
- ✱ Hierarchical k-means speeds up clustering for very large datasets
- ✱ K-medoids allows clustering of data that cannot be averaged

Q. What is different between a hierarchical clustering (hc) and k-means?

- A. HC produces dendrogram while k-means results in only flat clusters.
- B. HC doesn't need to choose number of clusters while k-means needs that step.
- C. HC has higher order time complexity than k-means
- ☒ D. All the above.

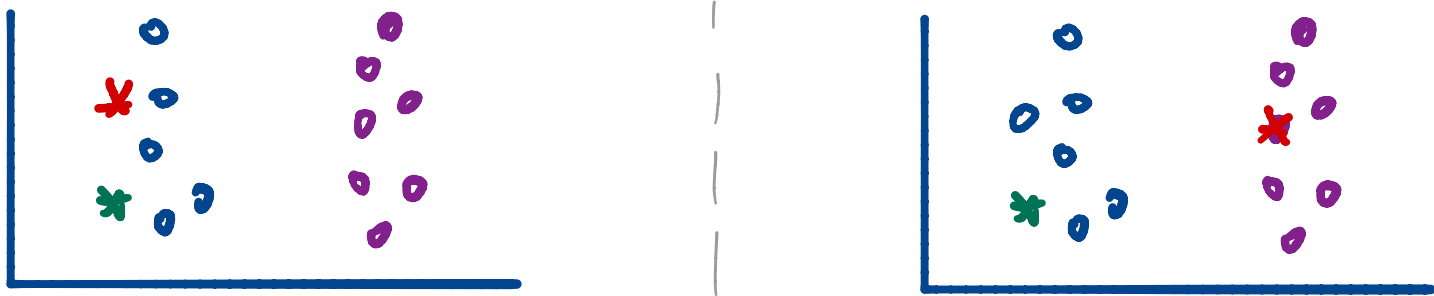
Some issues with k-means clustering

- ✱ Sensitive to outlier: example



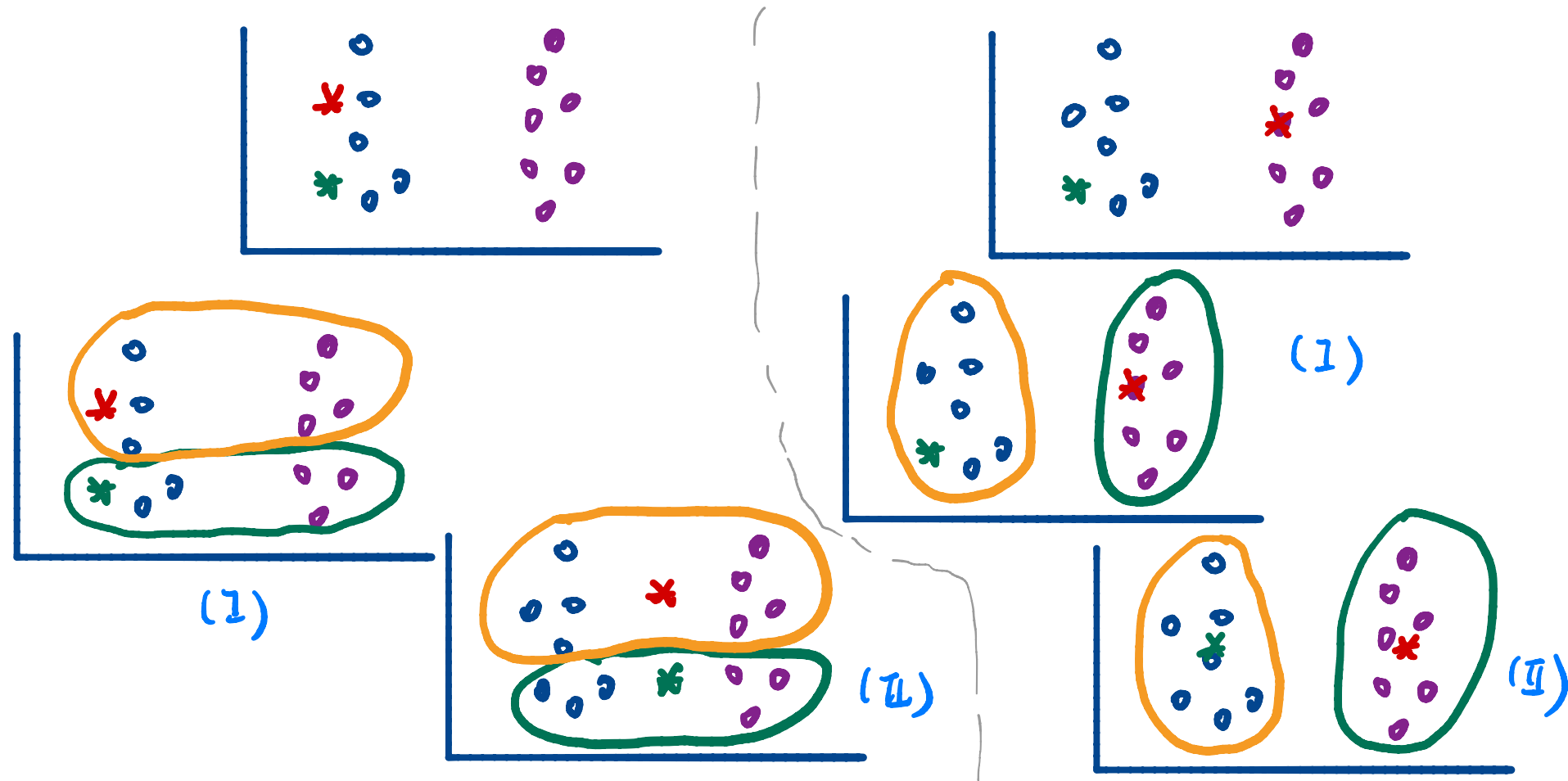
Some issues with k-means clustering

- ✱ Sensitive to the seeds (example)



Some issues with k-means clustering

✱ Sensitive to the seeds (example)



Assignments

- ✱ Read Chapter 11 of the textbook
- ✱ Week 13 Module, Quiz
- ✱ Final Project
- ✱ Happy Thanksgiving!



Additional References

- ✱ Robert V. Hogg, Elliot A. Tanis and Dale L. Zimmerman. “Probability and Statistical Inference”
- ✱ Kelvin Murphy, “Machine learning, A Probabilistic perspective”

See you next time

*See
You!*

