



COMP3055

Machine Learning

Topic 8 – Data Clustering

Ying Weng
2024 Autumn

Supervised VS Unsupervised Learning

- Supervised learning
 - Learns a function that maps an input to an output based on example input-output pairs.
 - Training data is labeled.
- Unsupervised learning
 - Learns from test data that has not been labeled.
 - Learns relationships between elements in a data set and classify the raw data without "help."
 - Typical application includes data clustering.

Motivating Problems

- ★ A true colour image – 24bits/pixel, R – 8 bits, G – 8 bits, B – 8 bits



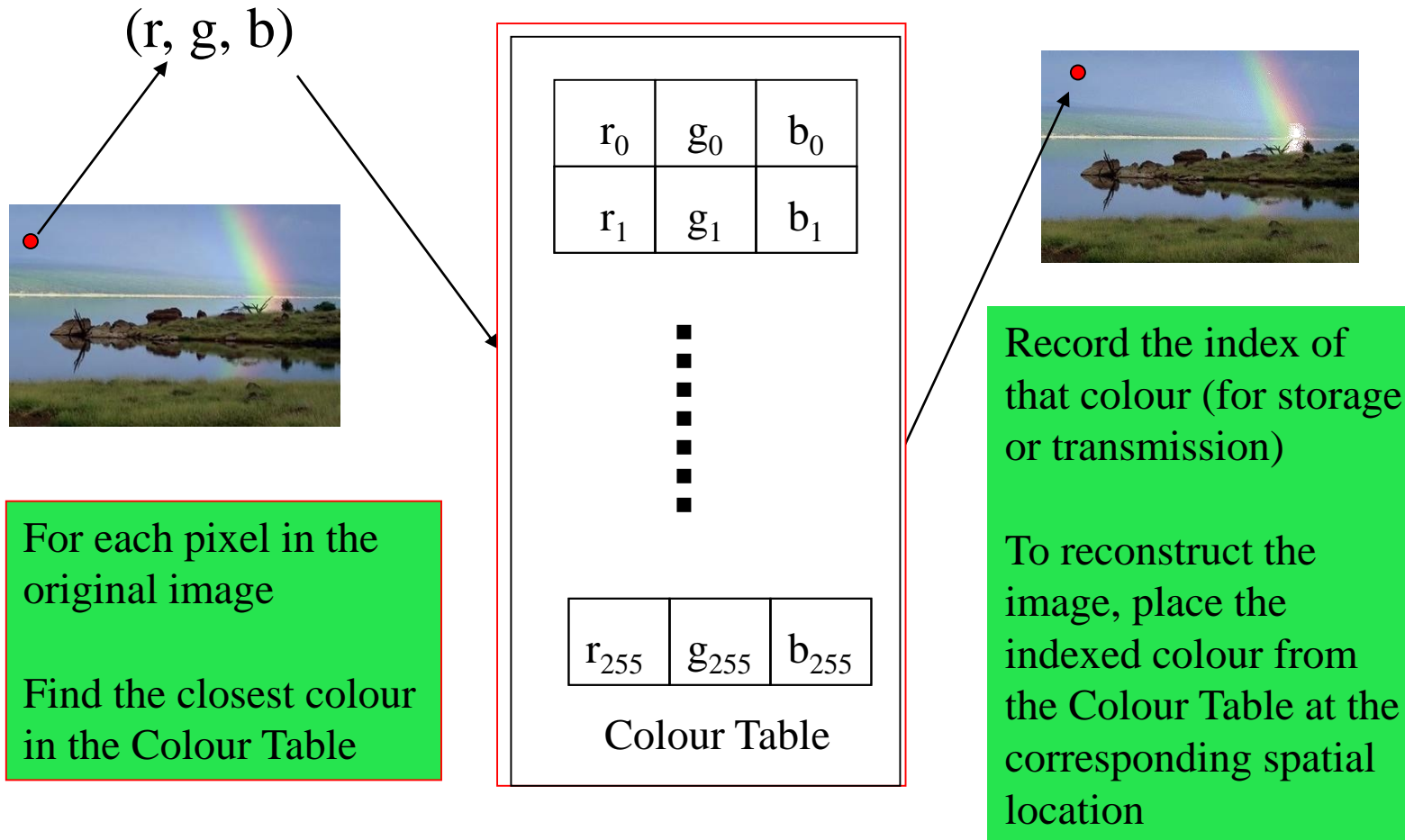
16777216 possible colours

- ★ A gif image - 8bits/pixel

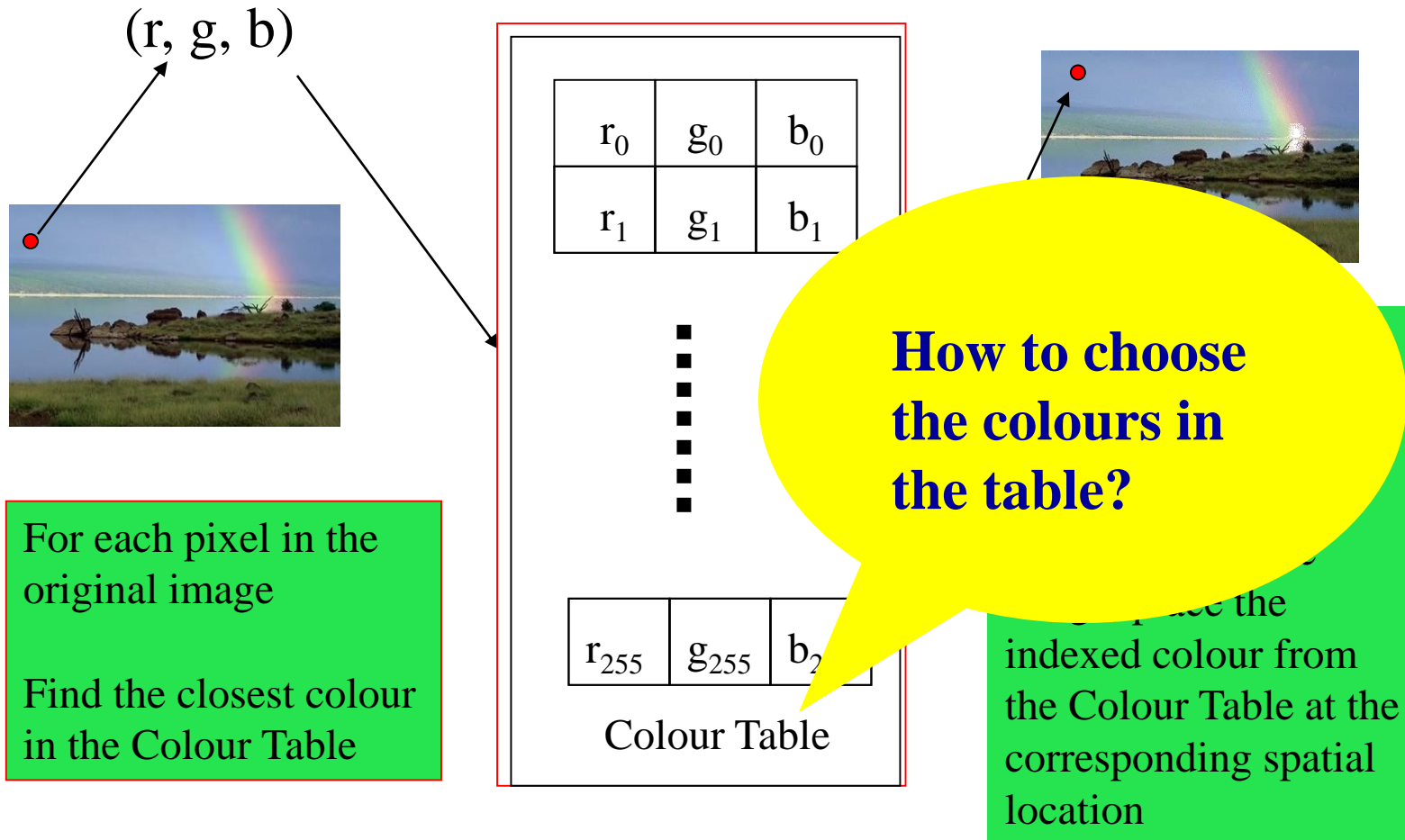


256 possible colours

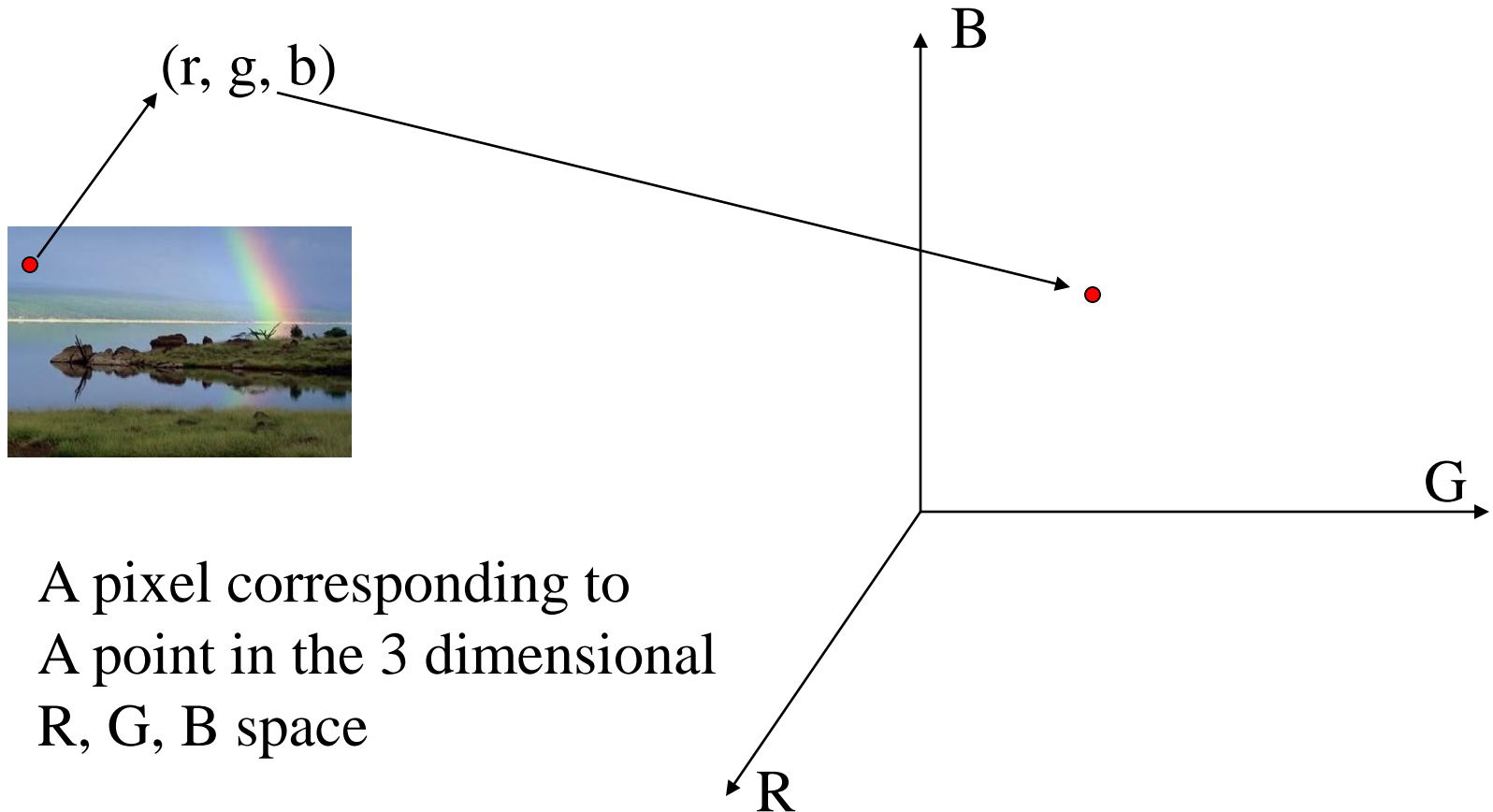
Motivating Problems



Motivating Problems



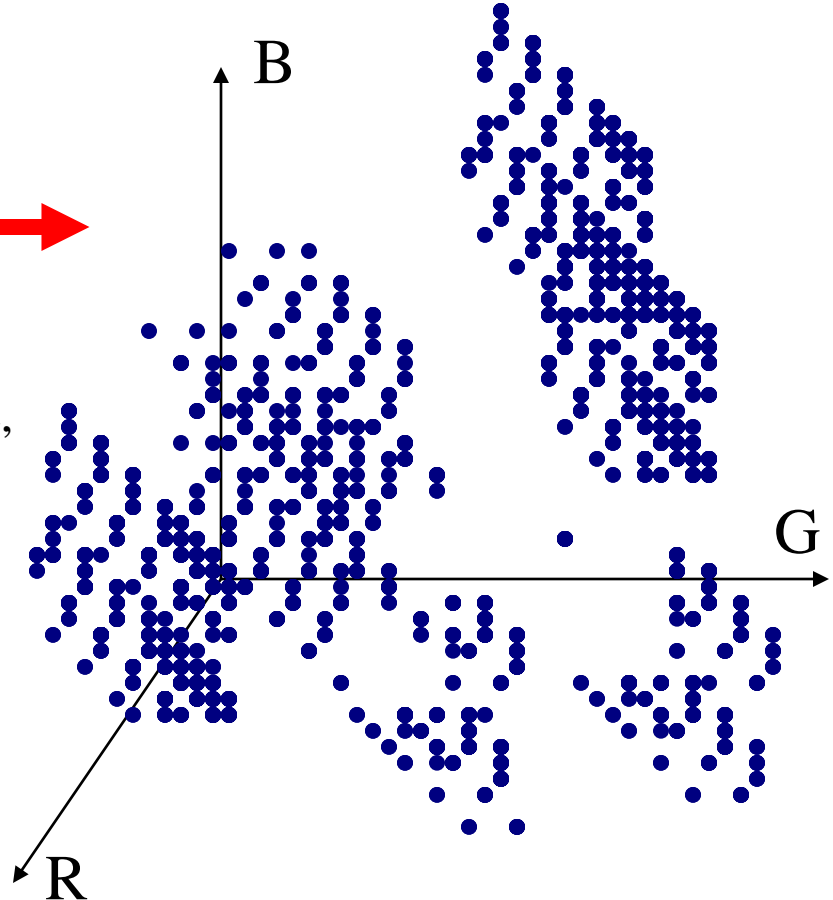
Motivating Problems



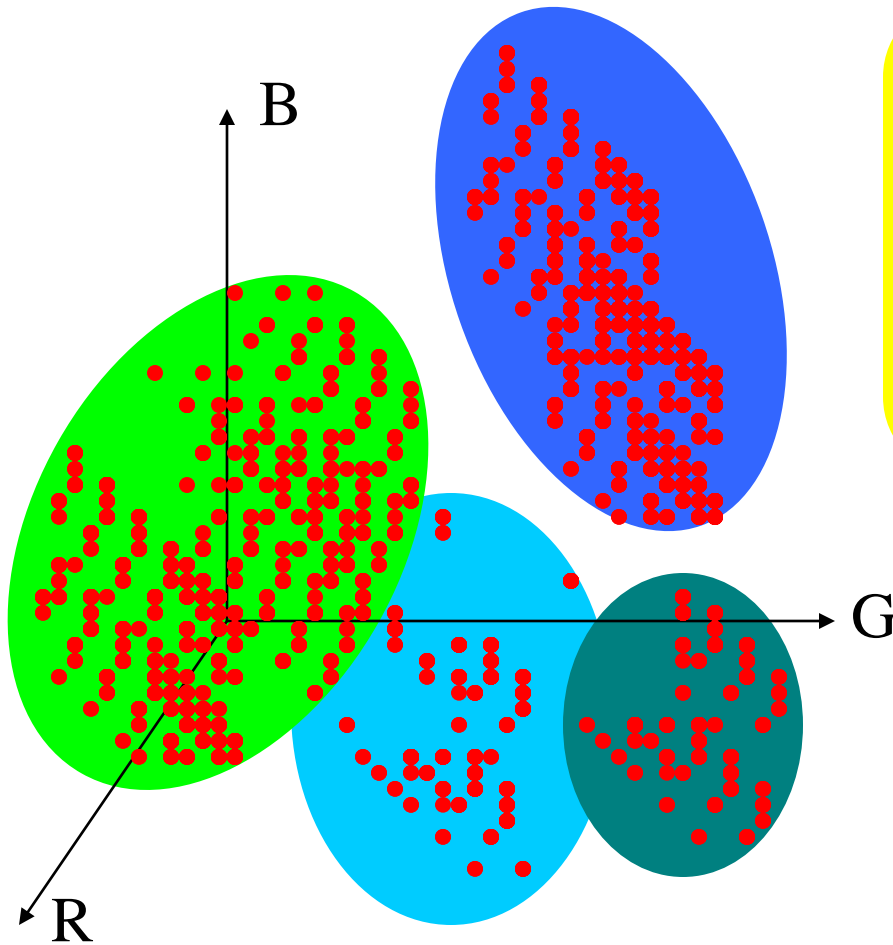
Motivating Problems



Map all pixels into the R, G, B space,
“clouds” of pixels are formed

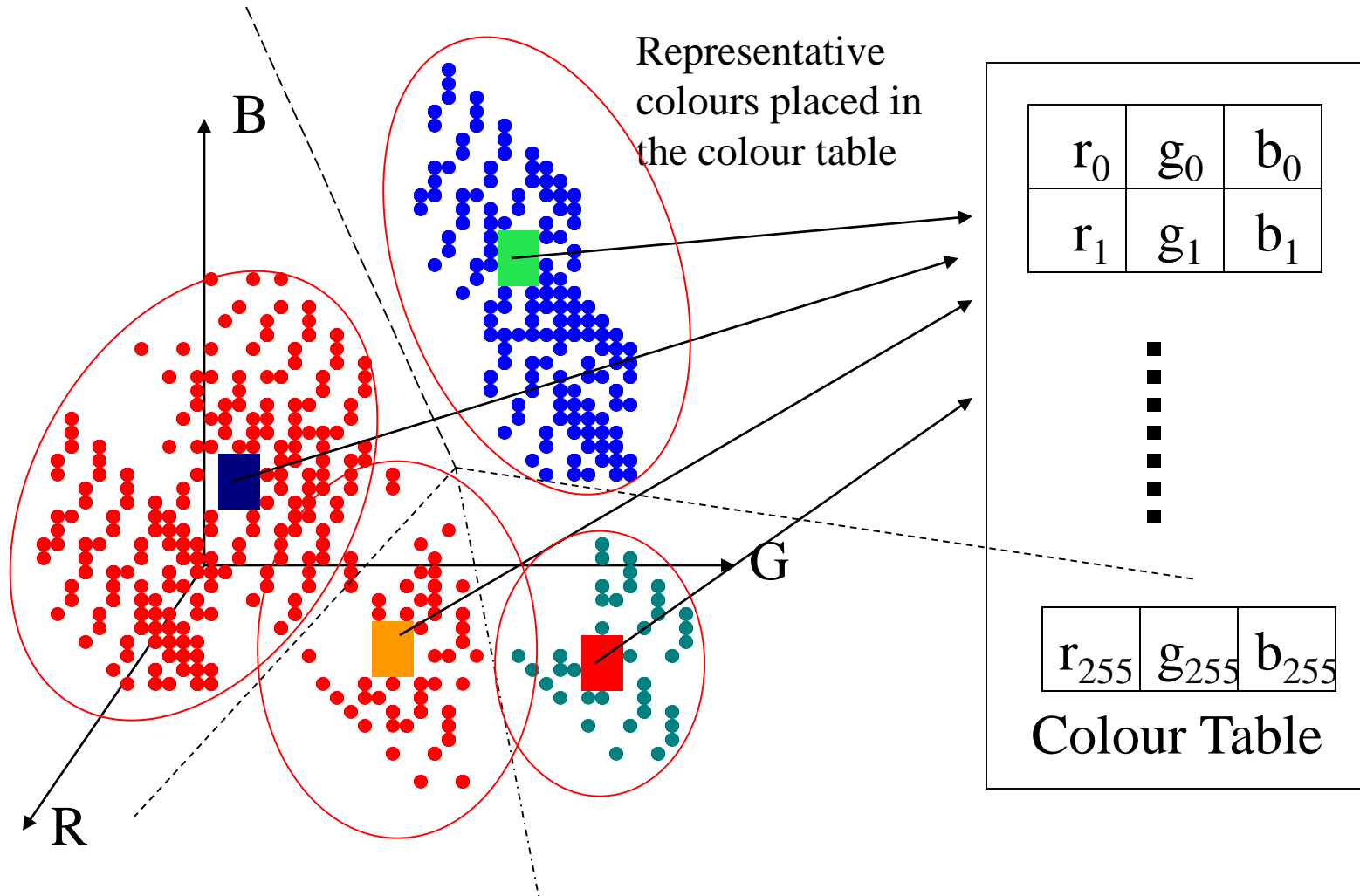


Motivating Problems



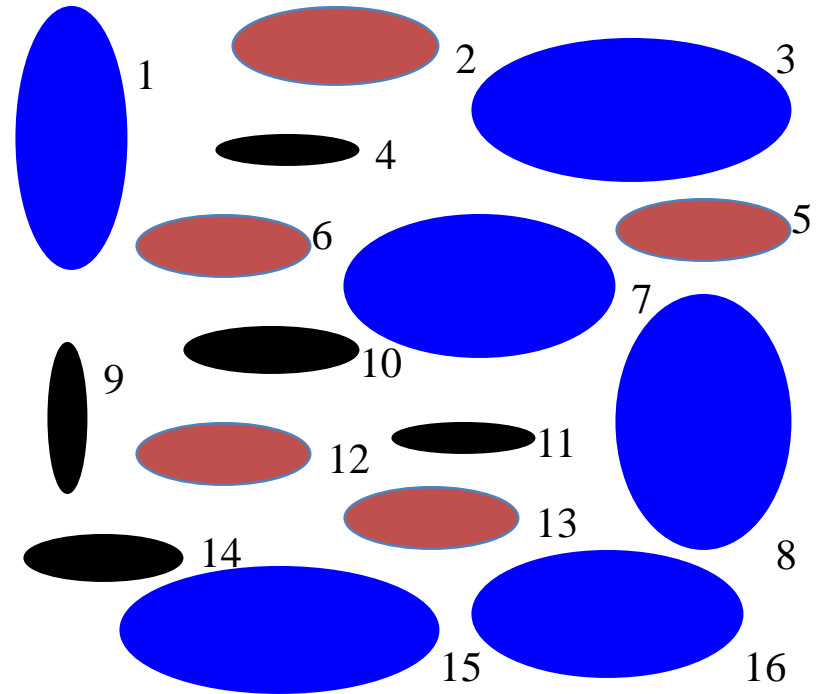
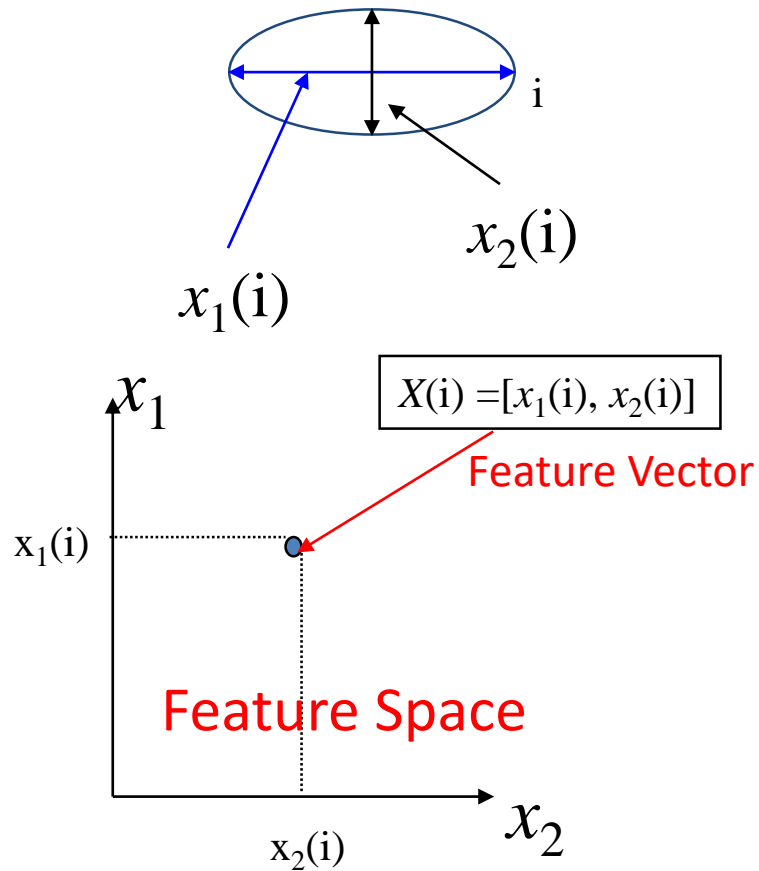
Group pixels that are close to each other, and replace them by one single colour

Motivating Problems



Motivating Example

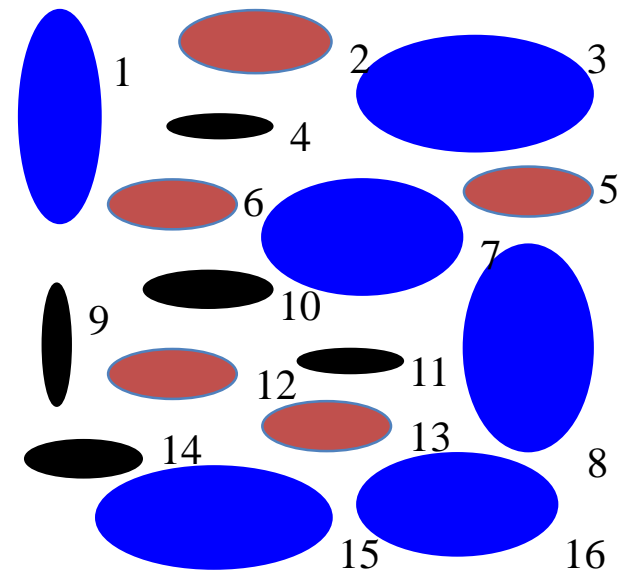
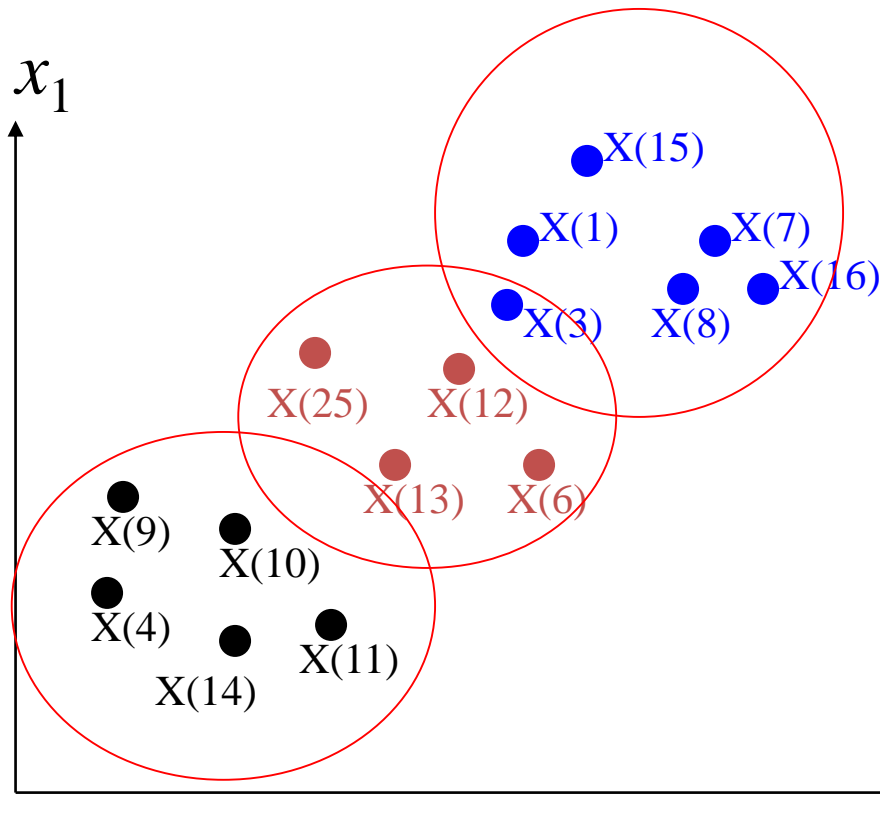
- Classify objects (Oranges, Potatoes) into large, middle, small sizes



Elliptical blobs (objects)

Motivating Example

- ★ From **Objects** to **Feature Vectors** to **Points** in the **Feature Space**



Elliptical blobs (objects)

Motivation of Clustering

- Patterns within a valid cluster are *more similar to each other* than they are to a pattern belonging to a different cluster.
- In clustering, the problem is to group a given collection of *unlabeled patterns* into meaningful clusters. Clustering is a *data driven method*, the clusters are obtained solely from the data.
- Clustering could be used in the field of pattern-analysis, grouping, decision-making, and machine-learning situations, including data mining, document retrieval, image segmentation

K-Means

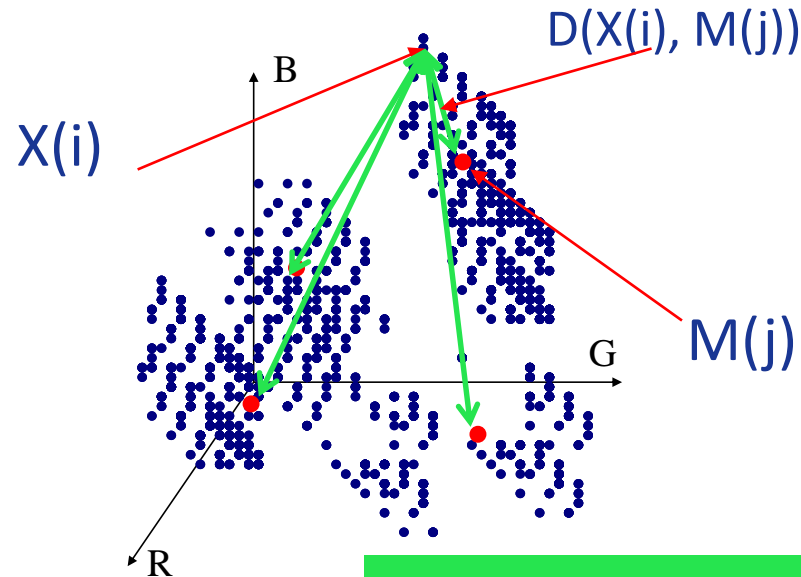
- ★ An algorithm for partitioning (or clustering) N data points into K disjoint subsets S_j containing N_j data points
 - ★ Define, $X(i) = [x_1(i), x_2(i), \dots, x_n(i)]$, $i = 1, 2, \dots, N$, as N data points
 - ★ We want to cluster these N points into K subsets, or K clusters, where K is pre-set
 - ★ For each cluster, we define $M(j) = [m_1(j), m_2(j), \dots, m_n(j)]$, $j=1, 2, \dots, K$, as its prototype or cluster centroids
 - ★ Define the distance between data point $X(i)$ and cluster prototype $M(j)$ as

$$D(X(i), M(j)) = \|X(i) - M(j)\| = \sqrt{\sum_{l=1}^n (x_l(i) - m_l(j))^2}$$

K-Means

- ★ A data point $X(i)$ is assigned to the j th cluster, $C(j)$, $X(i) \in C(j)$, if following condition holds

$$D(X(i), M(j)) \leq D(X(i), M(l)) \quad \text{for all } l = 1, 2, \dots, k$$



Minimum distance classifier

K-Means Algorithm

Step 1

- ★ Arbitrarily choose from the given sample set k initial cluster centres,

$$M^{(0)}(j) = [m^{(0)}_1(j), m^{(0)}_2(j), \dots, m^{(0)}_n(j)] \quad j = 1, 2, \dots, K,$$

e.g., the first K samples of the sample set
or can also be generated randomly

Set $t = 0$ (t is the iteration index)

K-Means Algorithm

Step 2

- ★ Assign each of the samples $X(i) = [x_1(i), x_2(i), \dots, x_n(i)]$, $i = 1, 2, \dots, N$, to one of the clusters according to the distance between the sample and the centre of the cluster:

$$\begin{aligned} X(i) &\in C^{(t)}(j) \\ \text{if } D(X(i), M^{(t)}(j)) &\leq D(X(i), M^{(t)}(l)) \\ \text{for all } l &= 1, 2, \dots, k \end{aligned}$$

K-Means Algorithm

Step 3

Update the cluster centres to get

$$M^{(t+1)}(j) = [m^{(t+1)}_1(j), m^{(t+1)}_2(j), \dots, m^{(t+1)}_n(j)] ; j = 1, 2, \dots, K$$

according to

$$M^{(t+1)}(j) = \frac{1}{N_j^{(t)}} \sum_{X(i) \in C^{(t)}(j)} X(i)$$

$N_j^{(t)}$ is the number of samples in $C^{(t)}_j$

K-Means Algorithm

Step 4

- Calculate the error of approximation

$$E(t) = \sum_{j=1}^K \sum_{X(i) \in C^{(t)}(j)} \|X(i) - M^{(t)}(j)\|$$

K-Means Algorithm

Step 5

- If the terminating criterion is met, then stop, otherwise

Set $t = t+1$

Go to Step 2.

K-Means Algorithm

Stopping criteria

- The K-means algorithm can be stopped based on following criteria

1. The errors do not change significantly in two consecutive epochs

$|E(t) - E(t-1)| < \varepsilon$, where ε is some preset small value

2. No further change in the assignment of the data points to clusters in two consecutive epochs.
3. It can also stop after a fixed number of epochs regardless of the error

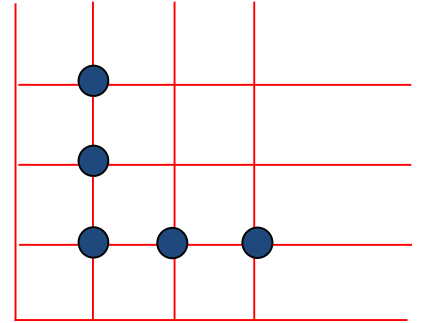
K-Means Algorithm

A worked example to see how it works exactly

Five 2-dimensional data points:

$(1, 1)$, $(2, 1)$, $(3, 1)$, $(1, 2)$, $(1, 3)$

Cluster them into two clusters and find the cluster centres



K-Means Algorithm

A worked example to see how it works exactly

(1) Euclidean distance to $m_1^0(1,2)$

$$\sqrt{(1-1)^2 + (1-2)^2} = 1$$

$$\sqrt{(2-1)^2 + (1-2)^2} = \sqrt{2} = 1.41$$

$$\sqrt{(3-1)^2 + (1-2)^2} = \sqrt{5} = 2.24$$

$$\sqrt{(1-1)^2 + (2-2)^2} = 0$$

$$\sqrt{(1-1)^2 + (3-2)^2} = 1$$

Euclidean distance to $m_2^0(3,1)$

$$\sqrt{(1-3)^2 + (1-1)^2} = 2$$

$$\sqrt{(2-3)^2 + (1-1)^2} = 1$$

$$\sqrt{(3-3)^2 + (1-1)^2} = 0$$

$$\sqrt{(1-3)^2 + (2-1)^2} = \sqrt{5} = 2.24$$

$$\sqrt{(1-3)^2 + (3-1)^2} = \sqrt{8} = 2.83$$

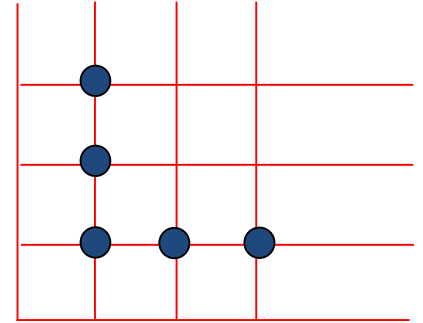
$\therefore C_1$

$\therefore C_2$

$\therefore C_2$

$\therefore C_1$

$\therefore C_1$



$$C_1 \text{ class: } (1,1), (1,2), (1,3) \Rightarrow m_1^{(0+1)}: \quad \frac{1}{3} \sum_{i=1}^3 x_i = \frac{1+1+1}{3} = 1$$

$$\frac{1}{3} \sum_{i=1}^3 y_i = \frac{1+2+3}{3} = 2$$

$$C_2 \text{ class: } (2,1), (3,1) \Rightarrow m_2^{(0+1)}: \quad \frac{1}{2} \sum_{i=1}^2 x_i = \frac{2+3}{2} = 2.5$$

$$\frac{1}{2} \sum_{i=1}^2 y_i = \frac{1+1}{2} = 1$$

(2) Euclidean distance to $m_1^1(1,2)$

$$1$$

$$\sqrt{2} = 1.41$$

$$\sqrt{5} = 2.24$$

$$0$$

$$1$$

Euclidean distance to $m_2^1(2.5,1)$

$$\sqrt{(1-2.5)^2 + (1-1)^2} = 1.5$$

$$\sqrt{(2-2.5)^2 + (1-1)^2} = 0.5$$

$$\sqrt{(3-2.5)^2 + (1-1)^2} = 0.5$$

$$\sqrt{(1-2.5)^2 + (2-1)^2} = \sqrt{3.25} = 1.80$$

$$\sqrt{(1-2.5)^2 + (3-1)^2} = \sqrt{6.25} = 2.5$$

$\therefore C_1$

$\therefore C_2$

$\therefore C_2$

$\therefore C_1$

$\therefore C_1$

$\therefore C_1 \text{ class: } (1,1), (1,2), (1,3)$

$C_2 \text{ class: } (2,1), (3,1)$

K-Means Algorithm

- What is the algorithm doing exactly?
 - It tries to find the centre vectors $M(j)$'s that optimize the following cost function

$$E = \sum_{j=1}^K \sum_{X(i) \in C(j)} \|X(i) - M(j)\|$$

K-Means Algorithm

- Some remarks
 - Is a gradient descent algorithm, trying to minimize a cost function E
 - In general, the algorithm does not achieve a global minimum of E over the assignments
 - Sensitive to initial choice of cluster centers. Different starting cluster centroids may lead to different solution
 - Is a popular method, many more advanced methods derived from this simple algorithm

Any Questions?

