

K-Means Clustering

k-means clustering is a method of vector quantization, originally from signal processing, that aims to partition **n observations** into **k clusters** in which each observation belongs to the cluster with the nearest **mean** (cluster centers or cluster **centroid**), serving as a prototype of the cluster.

https://en.wikipedia.org/wiki/K-means_clustering

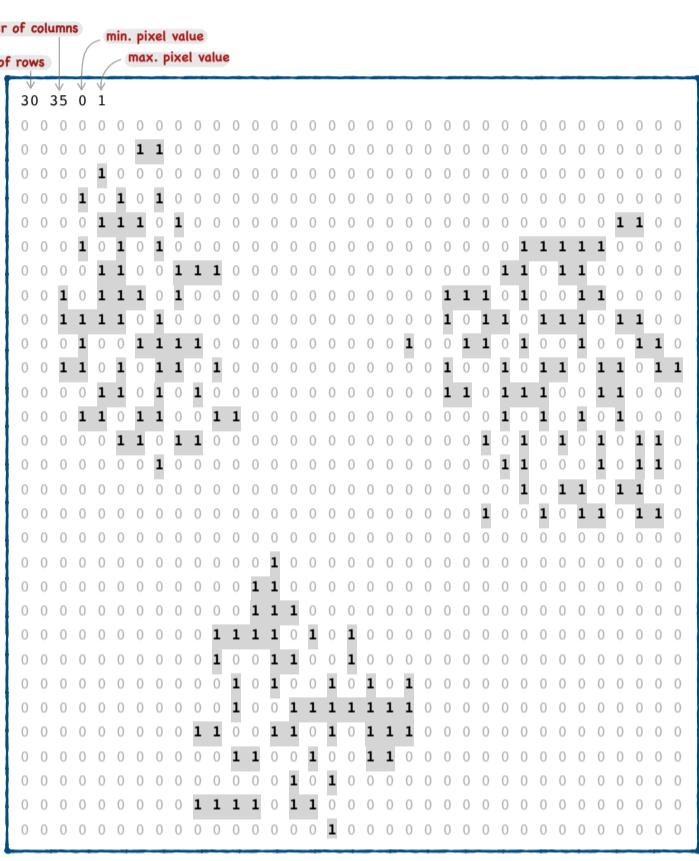
K-means Clustering is a clustering algorithm that works on data that has not been labeled yet. This means that: K-means Clustering is an **Unsupervised Learning** algorithm. K-means simply partitions the given dataset into groups (clusters) where each group has different features. K refers to the total number of clusters existing in the given dataset. The boundaries of each cluster are determined by the euclidean proximity between each of the clusters members and the cluster's centroid.

K-means Clustering works best for 2-Dimensional numerical data. In real applications, feature vectors (i.e. data members) usually have multiple dimensions, hence, dimensionality reduction has to be performed before K-means can be applied.

Task:

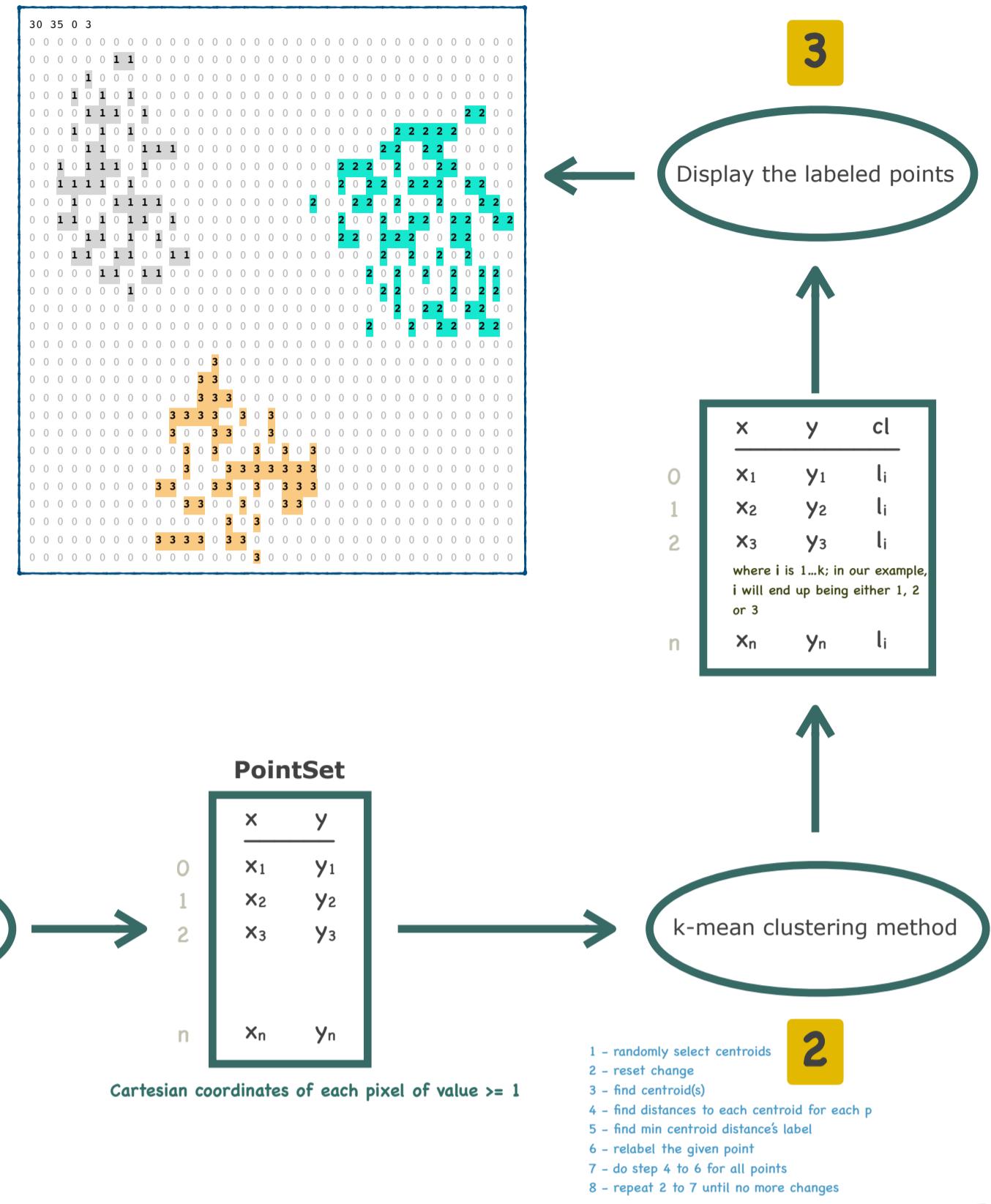
Given the following binary image (2-dimensions) containing 3 objects, where the background pixels have the value of 0 and the objects the value of 1, detect the objects by homogeneously labeling their constituent pixels. E.g.: label ALL pixels belonging to:

object-1 as 1
object-2 as 2
object-3 as 3



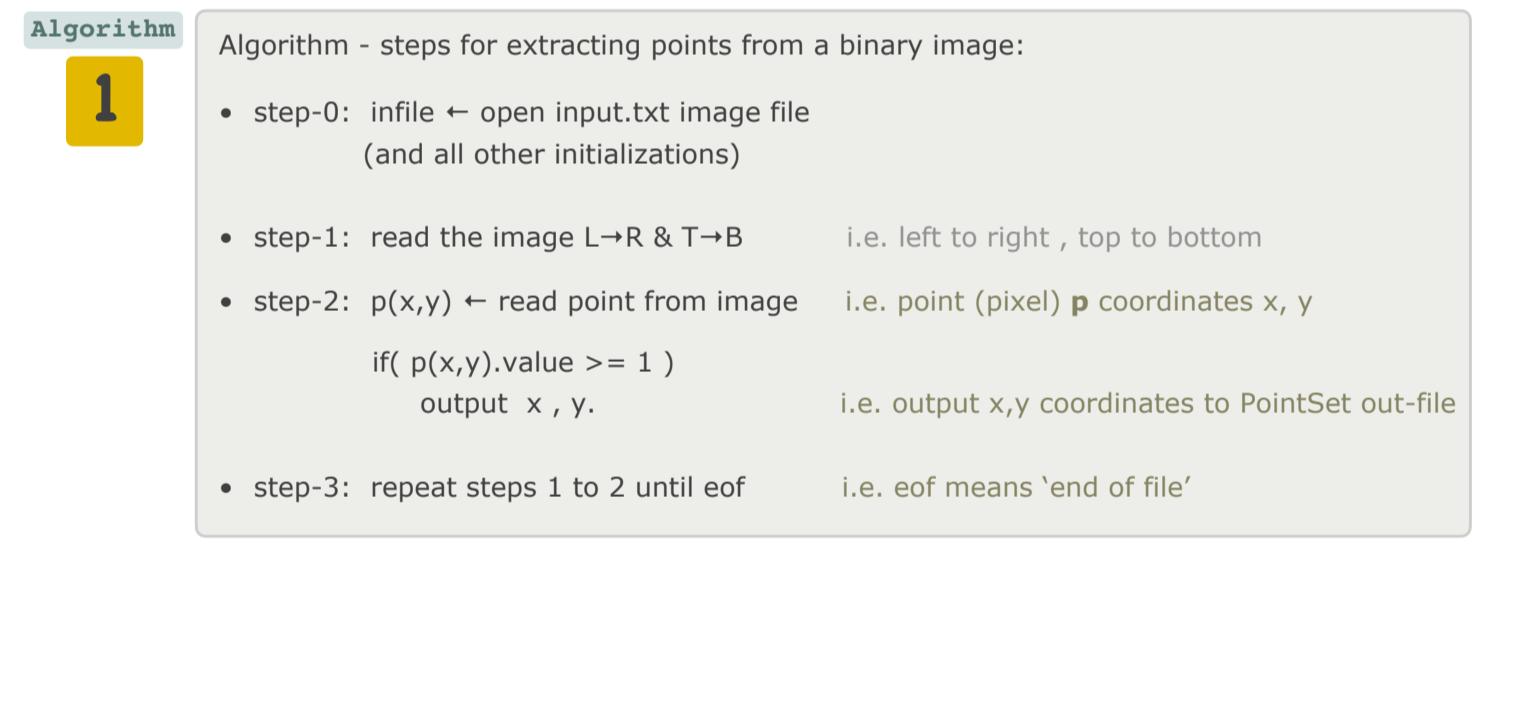
Extract points from the binary image.

1



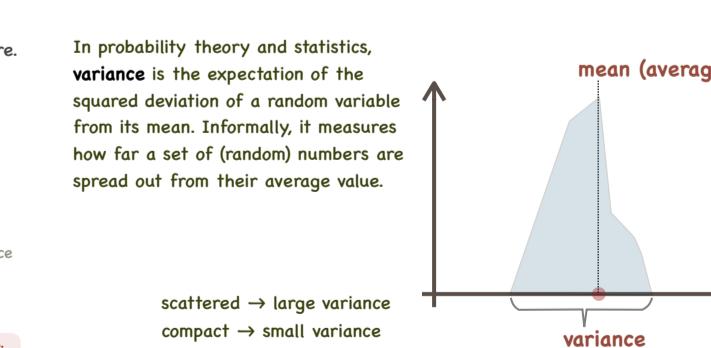
Extract the x,y-coordinates of ALL the points belonging to objects (i.e. pixels=1) from the above binary image into an output.txt file.

extracting points



Sometimes there will be an oscillating point that toggles e.g. blue to burgundy to blue. Those points are very rare, increasing the precision of distance measuring (e.g. float → double) should fix that.

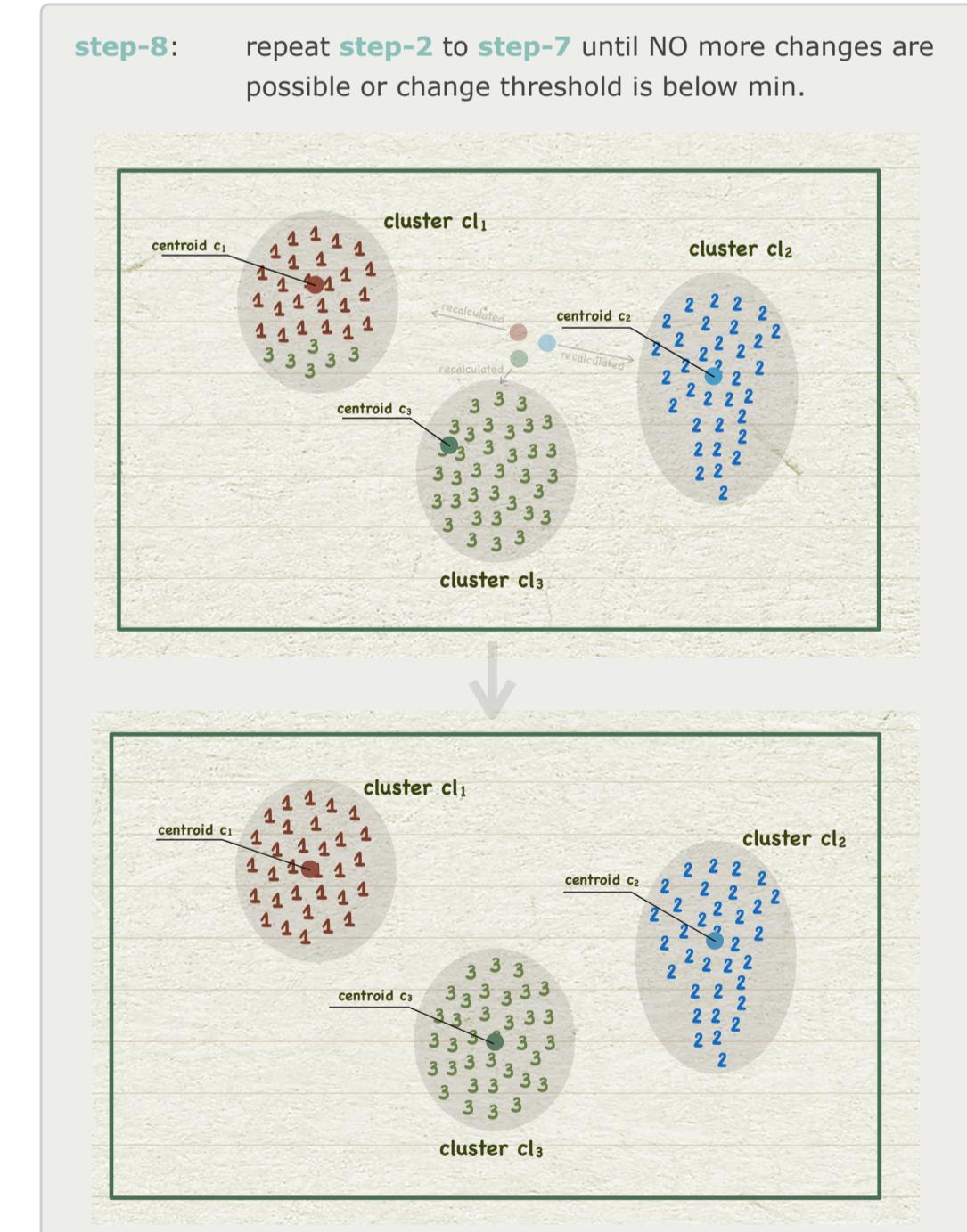
standard deviation σ is a measure of how spread out the numbers are.
 $\sigma = \sqrt{\text{variance}}$
variance is the average of the squared differences from the mean.
steps for computing variance:
Step1 work out the Mean (the simple average of the numbers)
Step2 $n \leftarrow$ get a number
Step3 $\text{result} \leftarrow (n - \text{Mean})^2$ // calculate the squared difference
repeat step 2 until all numbers are processed
Step4 compute variance:
$$\text{variance} = \frac{\text{result} + \text{result}_2 + \dots + \text{result}_n}{n}$$



In probability theory and statistics, variance is the expectation of the squared deviation of a random variable from its mean. Informally, it measures how far a set of (random) numbers are spread out from their average value.

What if we do NOT know K?

- in step0 k (unknown) ← start with 1
- in step7 compute variance for each k-cluster then sum-up
- in step8 repeat steps 2 to 7 until variance is acceptable (to the user), otherwise k++.



2 k-mean clustering

Algorithm steps for k-mean clustering:

step-0: instantiations (allocation), initializations

pointSet ← open input image file and extract all points

$n \leftarrow$ number of points (count the points when extracting)

$k \leftarrow$ given by the user

change ← 0

2 alternative approaches:
1. define "change" based on the Centroids travel distance. If the Centroids travel below some threshold (maybe 1 pixel), then no more changes are done.

2. define "change" as the variation in the "total distance" from all labels (in a cluster) to their Centroid. The smallest total distance wins. This means that the final state of each run i.e. the solved results as well as their corresponding "total distance" will have to be stored.

$cl_1, cl_2, \dots, cl_k \leftarrow$ clusters

$c_1, c_2, \dots, c_k \leftarrow$ centroids

step-1: randomly select k different points from among the points $p(x_i, y_i)$ in the pointSet and appoint them as starting centroids., where $1 \leq i, j \leq n$

$c_1, c_2, \dots, c_k \leftarrow$ some random $p(x_i, y_i)$

go to step-4

step-2: Reset changes

change ← 0

step-3: Find centroids

for each cluster $cl_1..k$ compute centroid $c_i(x(cl_i), y(cl_i))$, where $1 \leq i \leq n$.

$$x(c_i) = \frac{\sum x_i}{\sum 1}$$

$$y(c_i) = \frac{\sum y_i}{\sum 1}$$

e.g. centroid $c_1(x(c_1), y(c_1))$ think: the Center of Mass of ALL points labeled 1

Def.

Centroid: (the center of mass) in mathematics and physics, the centroid or geometric center of a plane figure is the arithmetic mean position of all the points in the figure. Informally, it is the point at which a cutout of the shape could be perfectly balanced on the tip of a pin.

extra dimensions

Note that: our example computes only in 2-dimensions x and y. If we want to increase the number of dimensions (e.g. to 7-dimensions because our feature vector has 7 attributes), we must compute the centroid for each of those extra dimensions i.e. $z(c_i)$, $K(c_i)$, $m(c_i)$... Also we extend the distance calculation by those same extra dimensions i.e.
$$\text{distance} = \sqrt{(z_1 - z_2)^2 + (K_1 - K_2)^2 + (m_1 - m_2)^2}$$

