# K - Means Clustering

- can be on line, graph, heatmap.

- using computer to identify / group data into meaningful clusters

## K - Means clustering steps

1) Select number of clusters to identify (k)

2) Randomly select k data points to start (initial centroid)

3) Select selected data point & measure distance from all k centroids

4) Assign point to nearest centroid

5) Repeat until all points allocated

6) Calculate mean of each cluster

7) Reposition centroid to mean of each cluster

   ↳ reallocate points according to new centroid

   ↳ stop once centroid converges.

8) Assess quality of clustering by measuring variance within clusters
   ↑

Repeat entire process with different starting points.

## Determining ideal k

- As k increases, total variation ↓

- if k = N, total variation = 0

* consider reduction in variance for k

  ↳ Elbow plot to determine largest reduction in variance

## Multi - Dimension K - Mean clustering

- Key idea behind K - Means clustering is finding Euclidean distance

  Suppose 3 dimensions:

  1) Randomly choose k centroids $(x_1, y_1, z_1)$

  2) Calc. dist. b/w pt & centroid $\sqrt{(x_2-x_1)^2 + (y_2-y_1)^2 + (z_2-z_1)^2}$

  3) compute mean of cluster as new centroid : $\left( \dfrac{\Sigma x_i}{N_1}, \dfrac{\Sigma y_i}{N_1}, \dfrac{\Sigma z_i}{N_1} \right)$

  4) Repeat 2 & 3 until centroid converge.

  5) Evaluate Within-Cluster Sum of Squares (Inertia)

     ↳ used in Elbow method.

  cluster
  data points.